

# Kuhn-Tucker and Multiple Discrete-Continuous Extreme Value Model Estimation and Simulation in R: The `rmdcev` Package

by Patrick Lloyd-Smith

**Abstract** This paper introduces the package `rmdcev` in R for estimation and simulation of Kuhn-Tucker demand models with individual heterogeneity. The models supported by `rmdcev` are the multiple-discrete continuous extreme value (MDCEV) model and Kuhn-Tucker specification common in the environmental economics literature on recreation demand. Latent class and random parameters specifications can be implemented and the models are fit using maximum likelihood estimation or Bayesian estimation. The `rmdcev` package also implements demand forecasting and welfare calculation for policy simulation. The purpose of this paper is to describe the model estimation and simulation framework and to demonstrate the functionalities of `rmdcev` using real datasets.

## Introduction

Individual choice contexts are often characterized by both extensive (i.e. what alternative to choose) and intensive (i.e. how much of an alternative to consume) margins (Bhat, 2008). These multiple discrete-continuous (MDC) choice situations are pervasive, arising in transportation, marketing, health, and decisions regarding environmental resources (Bhat and Pinjari, 2014). The Kuhn-Tucker (KT) modelling framework is often employed to analyze these MDC situations and substantial progress has been made in improving these econometric modeling structures (von Haefen and Phaneuf, 2005; Bhat and Pinjari, 2014). Despite the large potential applications for KT models, there remains a gap between this potential and actual examples of these models being used. One of the reasons cited for the lack of widespread use of KT models is that estimating and simulating these models is challenging. The explanations of methods used to work with these models are spread across many papers and few user friendly software tools are available. The purpose of this paper is to present a unified account for KT estimation and simulation alongside computer code for easy and efficient implementation.

This paper presents an overview of the R package `rmdcev` which can estimate and simulate KT demand models with discrete or continuous unobserved individual heterogeneity.<sup>1</sup> The common starting point for all KT models is the individual's constrained optimization problem and exploiting the resulting KT first order conditions in estimation. The most popular empirical KT modelling framework is the multiple-discrete continuous extreme value (MDCEV) model as first introduced by Bhat (2008). A separate stream of literature in the environmental economics on recreation demand has developed a closely related set of models and use the term KT to describe the models. In this paper, we use KT to describe the general modelling framework, MDCEV to describe the Bhat (2008) specifications, and KT-EE to describe the environmental economics literature KT specification (von Haefen et al., 2004). One of the main differences between the MDCEV and KT-EE frameworks is how alternative-specific attributes enter the utility function, a point we describe in the paper.

Incorporating preference heterogeneity has been an important advancement in choice modeling. Both the MDCEV and KT-EE specifications can be estimated to incorporate unobserved preference heterogeneity by assuming continuous distributions using random parameters or using a latent class (LC) specification assuming a discrete distribution where people can be divided into distinct segments. The models in `rmdcev` can be fit using maximum likelihood estimation or Bayesian estimation. Besides estimation, the `rmdcev` package also implements demand forecasting and welfare calculation for policy simulation. The two main functions in the `rmdcev` are `mdcev` used to estimate all model specifications and `mdcev.sim` used to simulate both demand and welfare implications. `rmdcev` is available from the Comprehensive R Archive Network (CRAN) at <https://CRAN.R-project.org/package=rmdcev> as well as from GitHub at <https://github.com/plloydsmith/rmdcev>.

While there are several R packages available to estimate discrete choice data such as `apollo` (Hess and Palma, 2019), `mlogit` (Croissant, 2019), and `gmnl` (Sarrias and Daziano, 2017)<sup>2</sup>, there are limited options for users interested in estimating and simulating KT models. In addition to `rmdcev`, the

<sup>1</sup>This paper uses version 1.2.0 of the `rmdcev` package.

<sup>2</sup>Sarrias and Daziano (2017) provides a good overview of the different R packages available to estimate discrete choice models

**apollo** package developed by Stephane Hess and David Palma at the Choice Modelling Centre in Leeds provides a flexible modelling platform for estimating MDCEV models and simulating demand behaviour (Hess and Palma, 2019). **apollo** estimates a full suite of choice models including discrete choice models and is thus more comprehensive and flexible than **rmdecv**. The main advantages for KT modeling in using the **rmdecv** is that it 1) provides functions for calculating welfare implications of policy scenarios, 2) allows the estimation and simulation of the KT formulation used in environmental economics (von Haefen and Phaneuf, 2005), 3) uses the Stan program (Carpenter et al., 2017) for Bayesian estimation and thus the user has access to specialized postestimation commands, and 4) is primarily coded in C++ and thus around 20 times faster. The main limitations **rmdecv** compared to **apollo** is that it 1) only estimates model specifications with an outside good that is always consumed whereas **apollo** can estimate models without an outside good, 2) users have more control over which particular parameters are fixed at their starting values and which are allowed to be random parameters, and 3) **apollo** allows users to estimate the multiple discrete continuous nested extreme value model and LC-random parameter MDCEV specifications.

The paper first introduces the conceptual framework underlying KT models and the connection to economic theory and welfare measures. Section 2 also describes the various empirical specifications for KT models. Section 3 introduces the **rmdecv** package focusing first on estimation before moving on to discuss how to conduct welfare and demand simulations. Section 4 provides conclusions of the paper.

## Models

### Conceptual framework

This section describes the underlying conceptual framework for KT models. Each individual  $i$  maximizes utility through the choice of the numeraire or outside good ( $x_{i1}$ ) and the non-numeraire alternatives ( $x_{ik}$ ) subject to a monetary or non-monetary budget constraint. We assume there is a numeraire good (i.e. essential Hicksian composite good) which is always consumed and has a price of one. The individual's maximization problem is

$$\begin{aligned} \max_{x_{ik}, x_{i1}} U(x_{ik}, x_{i1}) \\ \text{s.t. } y_i = \sum_{k=2}^K p_{ik} x_{ik} + x_{i1}, \quad x_{ik} \geq 0, \quad k = 2, \dots, K, \end{aligned} \quad (1)$$

where  $x_{ik}$  is the consumption level for alternative  $k$ ,  $x_{i1}$  is consumption of the numeraire,  $y_i$  is any arbitrary budget amount (e.g. annual income), and  $p_{ik}$  is the unit price of alternative  $k$ .

The resulting first-order KT conditions that implicitly define the solution to the optimal consumption bundles of  $x_{ik}$  and  $x_{i1}$  are

$$\begin{aligned} \frac{U_{x_{ik}}}{U_{x_{i1}}} &\leq p_{ik}, \quad k = 1, \dots, K, \\ x_{ik} \left[ \frac{U_{x_{ik}}}{U_{x_{i1}}} - p_{ik} \right] &= 0, \quad k = 1, \dots, K. \end{aligned} \quad (2)$$

For alternatives with positive consumption levels, the marginal rate of substitution between these alternatives and the numeraire good is equal to the price of the alternative. For unconsumed alternatives, the marginal rate of substitution between these alternatives and the numeraire good is less than the price of the alternatives. For the rest of the paper, we drop the subscript  $i$  for notational simplicity.

These first-order conditions can be used to derive Marshallian and Hicksian demands and welfare measures (von Haefen and Phaneuf, 2005). We assume that alternatives have non-price attribute  $q_k$  and the vector of  $k$  prices and attributes is denoted as  $p$  and  $q$ . The Hicksian compensating surplus ( $CS^H$ ) for a change in price and quality from baseline levels  $p^0$  and  $q^0$  to new 'policy' levels  $p^1$  and  $q^1$  is defined explicitly using an expenditure function

$$CS^H = y - e(p^1, q^1, \bar{U}, \theta, \varepsilon), \quad (3)$$

where  $\theta$  is the vector of structural parameters ( $\psi_k, \alpha_k, \gamma_k$ ),  $\varepsilon$  is a vector or matrix of unobserved heterogeneity, and  $\bar{U} = V(p^0, q^0, y, \theta, \varepsilon)$  and represents baseline utility.

## Multiple discrete-continuous extreme value model (MDCEV)

The **rmdecv** package implements the random utility specification of the MDCEV as introduced by [Bhat \(2008\)](#). The model specifications included in **rmdecv** always assume an outside good (i.e. the numeraire good that is always consumed by every individual). The general utility function is specified as

$$U(x_k, x_1) = \sum_{k=2}^K \frac{\gamma_k}{\alpha_k} \psi_k \left[ \left( \frac{x_k}{\gamma_k} + 1 \right)^{\alpha_k} - 1 \right] + \frac{\psi_1}{\alpha_1} x_1^{\alpha_1}, \quad (4)$$

where  $\gamma_k > 0$ ,  $\psi_k > 0$  and  $\alpha_k \leq 1$  for all  $k$  are required for this specification to be consistent with the properties of a utility function ([Bhat, 2008](#)). [Bhat \(2008\)](#) provides a detailed overview of the parameter interpretation and in brief

- The  $\psi_k$  parameters represent the marginal utility of consuming alternative  $k$  at the point of zero consumption (i.e. baseline marginal utility).
- The  $\gamma_k$  parameters are translation parameters that allow for corner solutions (i.e. zero consumption levels for alternatives) and also influence satiation. The lower the value of  $\gamma_k$ , the greater the satiation effect in consuming  $x_k$ .
- The  $\alpha_k$  parameters control the rate of diminishing marginal utility of additional consumption. If  $\alpha_k$  equal to one, then there is no satiation effects (i.e. constant marginal utility).

The ‘random utility’ element of the model is introduced into the baseline utility through a random error term as

$$\psi_k = \psi(z_k, \varepsilon_k) = \exp(\beta' z_k + \varepsilon_k), \quad (5)$$

where  $z_k$  is a set of variables that can include alternative-specific attributes and individual-specific characteristics, and  $\varepsilon_k$  is an error term that allows for the utility function to be random over the population. We assume an extreme value distribution that is independently distributed across alternatives for  $\varepsilon_k$  with an associated scale parameter of  $\sigma$ . For identification, we specify  $\psi_1 = e^{\varepsilon_1}$ .

To ensure the estimated utility function corresponds to economic theory we specify  $\gamma_k = \exp(\gamma_k^*)$  such that  $\gamma_k > 0$  and  $\alpha_k = \exp(\alpha_k^*) / (1 + \exp(\alpha_k^*))$  such that  $0 < \alpha_k < 1$ .  $\gamma_k^*$  and  $\alpha_k^*$  are estimated in the package and  $\gamma_k$  and  $\alpha_k$  are reported to the user. Similarly, we specify  $\sigma = \exp(\sigma^*)$ . Weak complementarity, which is required for deriving unique welfare measures ([Mäler, 1974](#)), is imposed in this specification by adding and subtracting one in the non-numeraire part of the utility function.

While the most general form of the MDCEV model includes  $\psi_k$ ,  $\gamma_k$ , and  $\alpha_k$  parameters for each alternative, [Bhat \(2008\)](#) discusses the identification concerns regarding estimating separate  $\gamma_k$  and  $\alpha_k$  parameters for each non-numeraire alternative. Typically only a subset of these parameters can be identified and there are three common utility function specifications:

1.  $\alpha$ -profile: set all  $\gamma_k$  parameters to 1.

$$U(x_k, x_1) = \sum_{k=2}^K \frac{1}{\alpha_k} \exp(\beta' z_k + \varepsilon_k) [(x_k + 1)^{\alpha_k} - 1] + \frac{\exp(\varepsilon_1)}{\alpha_1} x_1^{\alpha_1}. \quad (6)$$

2.  $\gamma$ -profile: set all non-numeraire  $\alpha_k$  parameters to 0.

$$U(x_k, x_1) = \sum_{k=2}^K \gamma_k \exp(\beta' z_k + \varepsilon_k) \ln \left( \frac{x_k}{\gamma_k} + 1 \right) + \frac{\exp(\varepsilon_1)}{\alpha_1} x_1^{\alpha_1}. \quad (7)$$

$$U(x_k, x_1) = \sum_{k=2}^K \exp(\beta' z_k + \varepsilon_k) \ln (\phi_k x_k + \gamma_k) + \frac{1}{\alpha_1} x_1^{\alpha_1}. \quad (8)$$

3. hybrid-profile: set all  $\alpha_k = \alpha_1 = \alpha$ .

$$U(x_k, x_1) = \sum_{k=2}^K \frac{\gamma_k}{\alpha} \exp(\beta' z_k + \varepsilon_k) \left[ \left( \frac{x_k}{\gamma_k} + 1 \right)^{\alpha} - 1 \right] + \frac{\exp(\varepsilon_1)}{\alpha} x_1^{\alpha}. \quad (9)$$

The likelihood function representing the model probability of the consumption pattern where  $M$

alternatives are chosen can be expressed as [Bhat \(2008\)](#)

$$P(x_1^*, x_2^* \dots x_M^*, 0, \dots, 0) = \frac{1}{\sigma^{M-1}} \left( \prod_{m=1}^M c_m \right) \left( \sum_{m=1}^M \frac{p_m}{c_m} \right) \left( \frac{\prod_{m=1}^M e^{V_m/\sigma}}{\left( \sum_{k=1}^J e^{V_k/\sigma} \right)^M} \right) (M-1)!, \quad (10)$$

where  $\sigma$  is the scale parameter and  $c_m = \frac{1-\alpha_m}{x_m+\gamma_m}$ . The  $V$  expression depend on what model specification is used:

1.  $\alpha$ -profile:  $V_k = \beta'z_k + (\alpha_k - 1) \ln(x_k + 1) - \ln(p_k)$  for  $k \geq 2$ , and  $V_1 = (\alpha_1 - 1) \ln(x_1)$ .
2.  $\gamma$ -profile:  $V_k = \beta'z_k - \ln\left(\frac{x_k}{\gamma_k} + 1\right) - \ln(p_k)$  for  $k \geq 2$ , and  $V_1 = (\alpha_1 - 1) \ln(x_1)$ .
3. hybrid-profile:  $V_k = \beta'z_k + (\alpha - 1) \ln\left(\frac{x_k}{\gamma_k} + 1\right) - \ln(p_k)$  for  $k \geq 2$ , and  $V_1 = (\alpha - 1) \ln(x_1)$ .

### Kuhn-Tucker model specifications in Environmental Economics (KT-EE)

The **rmdecv** package also implements the KT-EE specification ([von Haefen and Phaneuf, 2005](#)). The utility function in this specification is similar to the  $\gamma$ -profile of the MDCEV specification introduced above and is

$$U(x_k, x_1) = \sum_{k=2}^K \psi_k \ln(\phi_k x_k + \gamma_k) + \frac{1}{\alpha_1} x_1^{\alpha_1}, \quad (11)$$

where  $\phi_k > 0$ .<sup>3</sup>

An important difference between this KT formulation and the MDCEV models is the way weak complementarity is imposed. In this KT formulation, weak complementarity is imposed by only including alternative-specific attributes in the  $\phi_k$  parameter and not the  $\psi_k$  parameter.<sup>4</sup>

In this formulation, the estimating first-order conditions can be written as

$$\varepsilon_k \leq \frac{1}{\sigma} \left( -\beta's + \ln\left(\frac{p_k}{\phi_k}\right) + \ln(\phi_k x_k + \gamma_k) + (\alpha_1 - 1) \ln(y - p_k * x_k) \right), \quad \forall k, \quad (12)$$

and the resulting likelihood function as

$$P(x) = |J| \prod_k [\exp(-g_k(\cdot))/\sigma]^{1(x_k > 0)} \exp[-\exp(-g_k(\cdot))], \quad (13)$$

where  $|J|$  is the determinant of the Jacobian of transformation,  $g_k(\cdot)$  is the right hand side of Equation (12), and  $1(x_k > 0)$  is equal to one if  $x_k$  is positive and equal to zero if  $x_k$  is zero ([von Haefen and Phaneuf, 2005](#)). In previous implementations, the KT formulation used the computationally intensive numerical gradient approach to the calculation of the determinant of the Jacobian of transformation ([von Haefen and Phaneuf, 2005](#)).

The **rmdecv** package uses the compact structure of the determinant of the Jacobian as derived by [Bhat \(2008\)](#) and defined as

$$|J| = \frac{(1 - \alpha_1)}{x_1} \left[ \prod_m \frac{\phi_m}{\phi_m * x_m + \gamma_m} \right] \left[ x_1(1 - \alpha_1) + \sum_m \frac{(\phi_m * x_m + \gamma_m) * p_m}{\phi_m} \right], \quad (14)$$

where  $m$  denotes non-numeraire alternatives with positive consumption levels. Using this analytical gradient approach has the benefit of substantially speeding up estimation by around 70% relative to the numerical gradient approach.

In both the MDCEV and KT-EE specifications described above, the parameters  $(\beta, \alpha_k, \gamma_k, \phi_k, \sigma)$  are structural parameters that are assumed to be equal across the population which simplifies estimation. However, these fixed parameter specification is quite restrictive as they can only incorporate preference heterogeneity through interaction terms with observed individual characteristics. Without these interaction terms, the fixed specifications impose the assumption that all individuals have the same tastes for alternatives (i.e. preference homogeneity). This assumption is relaxed in the next two specifications which are able to accommodate both observed and unobserved preference heterogeneity.

<sup>3</sup>The environmental economics literature uses slightly different notation as typically  $\theta$  is used for  $\gamma$ ,  $\mu$  is used for  $\sigma$ , and  $\rho$  for  $\alpha_1$ . We change the notation slightly for consistency with the MDCEV model specifications.

<sup>4</sup>See [Herriges et al. \(2004\)](#) for more discussion on this point.

### Latent class (LC-KT) models

The latent class version of the KT model assumes that an individual belongs to a finite mixture of  $S$  segments each indexed by  $s$  ( $s = 1, 2, \dots, S$ ) (Sobhani et al., 2013; Kuriyama et al., 2010). Within each segment, the LC specification assumes preference homogeneity. We do not observe which segment an individual belongs to but we can attribute a probability  $\pi_{is}$  that individual  $i$  is a member of segment  $s$ . We impose that  $0 \leq \pi_{is} \leq 1$  and  $\sum_{s=1}^S \pi_{is} = 1$  through the use of the logit link function as

$$\pi_{is} = \frac{\exp(\delta'_s w_i)}{\sum_{s=1}^S \exp(\delta'_s w_i)}, \quad (15)$$

where  $w_i$  is a vector of individual characteristics and  $\delta_s$  is a vector of coefficients to be estimated. The  $\delta_s$  coefficients determine how the individual characteristics affect the membership of individual  $i$  in segment  $s$ . For identification, the  $\delta_1$  coefficients for the first segment are set to zero.

The likelihood function can be written as

$$P = \prod_i \pi_{is} P_{is}, \quad (16)$$

where  $P_{is}$  has the same form as Equations (10) and Equations (13) but is now class specific.

### Random parameters (RP-LC) models

The random parameter specification of the LC models assumes that the structural parameters  $\theta = (\beta, \alpha_k, \gamma_k)$  are not necessarily fixed but have an assumed distribution (Bhat, 2008). In **rmdecv**, parameters are distributed multivariate normal with a mean  $\bar{\theta}$  and variance covariance matrix  $\Sigma_\theta$  (von Haefen and Phaneuf, 2005). This structure allows for continuous preference heterogeneity and accommodates more flexible correlation patterns between alternatives in a similar fashion to the mixed logit model in discrete choice models. The  $\sigma$  scale parameter is always assumed to be a fixed parameter.

The most flexible model specification is to estimate the full variance covariance matrix and if there are  $Q$  parameters in  $\theta$  then there are  $Q(Q+1)/2$  unique variance covariance parameters to estimate in the correlated RP-MDCEV specification. An alternative is to assume the off-diagonal parameters are zero and estimate uncorrelated random parameters by estimating the  $Q$  diagonal elements of  $\Sigma_\theta$ . If all elements of  $\Sigma_\theta$  are assumed to be zero, the model collapses to the fixed KT structures.

### A note on Bayesian versus classical maximum likelihood estimation

The KT model without unobserved heterogeneity can be estimated using Bayesian or classical maximum likelihood techniques. The LC-KT model can only be estimated using classical maximum likelihood techniques as Bayesian approaches are challenged by the 'label switching' problem (Jasra et al., 2005). The RP-KT models can only be estimated using Bayesian techniques as random parameter models require simulated maximum likelihood estimators and these are not implemented in **rmdecv** at this time.

While there are philosophical differences between Bayesian and classical maximum likelihood techniques to estimating models, the Bernstein-von Mises theorem suggests that the Bayesian posterior distribution are asymptotically equivalent to maximum likelihood estimates if the data generating process has been correctly specified (Train, 2009).

## The rmdecv package

### Data format

The **rmdecv** uses `mdcev.data` function for handling multiple discrete-continuous data while ensuring the data is in the correct format and is suitable for estimation. The **rmdecv** package accepts data in "long" format (i.e. one row per available non-numeraire alternative for each individual). There is no row for the numeraire (i.e. outside) good. If there are  $I$  individuals and  $J$  non-numeraire alternatives, then the data frame should have  $I \times J$  rows.

To illustrate the suitable form of the data, we can load the recreation data included with the **rmdecv** package. This data is from the Canadian Nature Survey and includes choices for number of days spent recreating in 17 different outdoor activities for 2,000 people (Federal, Provincial, and Territorial Governments of Canada, 2014).

```
data(data_rec, package = "rmdcev")
```

Each recreation activity is characterized by the daily costs of participation for each individual. In addition to the recreation behaviour and prices, the data includes information on three individual characteristics: university (a dummy variable if the person has completed a university degree), ageindex (a person's age divided by the average age in sample), and urban (a dummy variable if a person lives in an urban area). Additional details on the data and price construction are provided in ?. We can summarize the average consumption and price levels for each alternative as:

```
aggregate(cbind(quant, price) ~ alt, data = data_rec, FUN = mean )
```

```
#>      alt    quant    price
#> 1    beach  6.5375  53.18359
#> 2   birding 14.3835  44.01734
#> 3   camping  2.5125  61.38326
#> 4   cycling  9.4700  45.99470
#> 5     fish   3.3435  86.22383
#> 6   garden 21.5710  38.28073
#> 7     golf   4.0260 134.10374
#> 8    hiking 41.4150  37.53204
#> 9  hunt_birds  0.4855 111.00176
#> 10 hunt_large  0.9480 184.46812
#> 11  hunt_trap  0.6290  95.33228
#> 12 hunt_waterfowl 0.2085 159.66605
#> 13   motor_land  3.7040 123.10169
#> 14  motor_water  2.8390 139.63845
#> 15     photo   8.6415  67.13733
#> 16   ski_cross  2.6450  32.65243
#> 17   ski_down  1.2065 151.01398
```

The data can be transformed into the structure for MDCEV estimation using the `mdcev.data` function:

```
data_mdcev <- mdcev.data(data_rec,
                        id.var = "id",
                        alt.var = "alt",
                        choice = "quant")
```

```
#> Sorting data by id.var then alt...
```

```
#> Checking data...
```

```
#> Data is good
```

The `id.var` argument indicates what variable uniquely identifies individuals in the data set, `alt.var` indicates the variable that identifies the non-numeraire alternatives, and `choice` indicates the level of consumption made by the individuals. Two other optional arguments of `mdcev.data` are `price` and `income` indicating the individual-specific price levels for each alternative, and the income level for each individual. These two arguments only need to be explicitly specified if they are not labeled `price` and `income`. Alternative-specific attributes and individual-specific characteristics can be included as additional columns and do not need to be specified in `mdcev.data`.

The `mdcev.data` function also checks to ensure the data has the necessary variables, and that all individuals spend positive amounts on the numeraire good. If an individual does not have positive expenditures on the numeraire good, an error message is given.

## KT estimation

### A general overview of `mdcev`

#### The `rmdcev`

All the various KT model specifications are estimated using the `mdcev` function.

```
args(mdcev)
```

```
#> function (formula = NULL, data, weights = NULL, model = c("alpha",
#>   "gamma", "hybrid", "hybrid0", "kt_ee"), n_classes = 1, fixed_scale1 = 0,
#>   trunc_data = 0, psi_ascs = NULL, gamma_ascs = 1, seed = "123",
#>   max_iterations = 2000, jacobian_analytical_grad = 1, initial.parameters = NULL,
#>   hessian = TRUE, algorithm = c("MLE", "Bayes"), flat_priors = NULL,
#>   print_iterations = TRUE, prior_psi_sd = 10, prior_gamma_sd = 10,
#>   prior_phi_sd = 10, prior_alpha_shape = 1, prior_scale_sd = 1,
#>   prior_delta_sd = 10, gamma_nonrandom = 0, alpha_nonrandom = 0,
#>   std_errors = "deltamethod", n_draws = 50, keep_loglik = 0,
#>   random_parameters = "fixed", show_stan_warnings = TRUE, n_iterations = 200,
#>   n_chains = 4, n_cores = 4, max_tree_depth = 10, adapt_delta = 0.8,
#>   lkj_shape_prior = 4, ...)
#> NULL
```

The main arguments are briefly explained below:

- **formula**: Formula for the model to be estimated as described in Section {#formula}.
- **data** The  $(I, x)$  data to be used in estimation as described above.
- **weights** An optional vector of sampling or frequency weights.
- **model** A string indicating which model specification to estimate. The four options are presented below:
  - “alpha”:  $\alpha$ -profile with all  $\gamma_k$  parameters fixed equal to 1 (Equation (6)).
  - “gamma”:  $\gamma$ -profile with one estimated  $\alpha_1$  and all non-numeraire  $\alpha_k$  parameters equal to 0 (Equation (8)).
  - “hybrid”: hybrid-profile with a single estimated  $\alpha$  parameter (i.e.  $\alpha_1 = \alpha_k = \alpha$ ) (Equation (9)).
  - “hybrid0”: hybrid-profile with all  $\alpha$  parameters fixed equal to 1e-3 (Equation (9)).
  - “kt\_ee”: Environmental economics version of KT model (Equation (11)).
- **n\_classes** The number of latent classes. Note that the LC model is automatically estimated as long as the prespecified number of classes is set greater than 1.
- **gamma\_ascs** Indicator to include alternative-specific gammas parameters.
- **psi\_ascs** Whether to include alternative-specific psi parameters. The first alternative is used as the reference category. Only specify to 1 for MDCEV models.
- **fixed\_scale1** Whether to fix the scale parameter at 1.
- **trunc\_data** Whether the estimation should be adjusted for truncation of non-numeraire alternatives. This option is useful if the data only includes individuals with positive non-numeraire consumption levels such as recreation data collected on-site. To account for the truncation of consumption, the likelihood is normalized by one minus the likelihood of observing zero consumption for all non-numeraire alternatives (i.e. likelihood of positive consumption) following Englin, Boxall and Watson (1998) and von Haefen (2003).
- **seed** Random seed.
- **algorithm** Either “Bayes” for Bayesian estimation or “MLE” for maximum likelihood estimation. The MLE algorithm uses the Limited-memory BFGS which approximates the Broyden–Fletcher–Goldfarb–Shanno (BFGS) algorithm but uses less computer memory.
- **flat\_priors** indicator if completely uninformative priors should be specified. Defaults to 1 if MLE used and 0 if Bayes used. If using MLE and set flat\_priors = 0, penalized MLE is used and the optimizing objective is augmented with the priors.
- **print\_iterations** Whether to print intermediate iteration information or not.
- **std\_errors** Compute standard errors using the delta method (“deltamethod”) or multivariate normal draws (“mvn”). The default is “deltamethod”. Note that mvn parameter draws should be used to incorporate parameter uncertainty for demand and welfare simulation. For maximum likelihood estimation only.
- **n\_draws** The number of multivariate normal draws for standard error calculations if “mvn” is specified.
- **initial.parameters** The default for fixed and random parameter specifications is to use random starting values. For LC models, the default is to use slightly adjusted MLE point estimates from the single class model. Initial parameter values should be included in a named list. For example, the LC “hybrid” specification initial parameters can be specified as: initial.parameters



= list(psi = array(0, dim = c(K, num\_psi)), gamma = array(1, dim = c(K, num\_alt)), alpha = array(0.5, dim = c(K, 1)), scale = array(1, dim = c(K))) where K is the number of classes (i.e. K = 1 is used for single class models), num\_psi is number of psi parameters, and num\_alt is number of non-numeraire alternatives.

### Formula format

The formula is used to incorporate alternative-specific variables and individual-specific characteristics into the  $\psi_k$  parameters, the membership equation of the LC-KT models, and  $\phi_k$  parameters for the KT-EE specification. By default, alternative-specific constants (ASCs) for all non-numeraire alternatives are included in the  $\psi_k$  and  $\gamma_k$  parameters. For the  $\psi_k$ , the first ASC is fixed at 0 due to identification concerns. They can be omitted using the `psi_ascs = 0` and `gamma_ascs = 0` arguments. Furthermore, the  $\gamma_k$ ,  $\alpha_k$ , and  $\sigma$  parameters cannot include alternative- or individual specific variables besides ASCs.

The formula is divided in three parts, separated by the symbol `|` and is based on the R package **Formula** (Zeileis and Croissant, 2010). The first part is reserved for the  $z_k$  variables in  $\psi_k$  as in Equation (5), excluding ASCs. These can include alternative-specific and individual-specific variables. Interaction terms between variables can be included using the normal **Formula** syntax of `z1 : z2`. This is particularly useful for creating interaction terms to incorporate observed preference heterogeneity for alternative-specific variables and individual-specific characteristics.

For a model with only ASCs in  $\psi_k$ , the formula can be specified as

```
f1 = ~ 0
```

We can add individual-specific variables to the  $\psi_k$  parameters as follows

```
f2 = ~ university + ageindex
```

Alternative-specific variables such as `z1` and `z2` can be included in the same way such as

```
f2 = ~ z1 + z2
```

The second part corresponds to individual-specific characteristics that enter in the probability assignment in models with latent classes. The formula will automatically include a constant in the membership equation but this can be omitted if `-1` is used in the formula. For example, a LC model with no alternative-specific variables in the  $\psi_k$  parameters and `university`, `ageindex` and a constant determine the class membership can be specified as

```
f3 = ~ 0 | university + ageindex
```

The third part is reserved for the  $q_k$  variables included in the  $\phi_k$  parameters in the KT-EE model specification (Equation 11) as in Equation (??). For example, if there was an alternative-specific variable named `'q1'`, it can be included as below

```
f4 = ~ 0 | 0 | q1
```

### Estimating KT using maximum likelihood techniques

We estimate a KT model by first calling `mdcev.data` on the **Recreation** data. For these examples we are going to use a subset of 200 individuals from the data.

```
data_model <- mdcev.data(data_rec, subset = id <= 200,
                        id.var = "id",
                        alt.var = "alt",
                        choice = "quant")
```

```
#> Sorting data by id.var then alt...
```

```
#> Checking data...
```

```
#> Data is good
```

We might think that older people prefer gardening to other activities and so we can include an interaction term between the activity `garden` and the variable `ageindex`. There are no alternative-specific variables besides constant terms to include in  $\psi$  and therefore the formula can be specified as



```
data_model$age_garden = ifelse(data_model$salt == "garden",
                              data_model$ageindex,0)

f5 = ~ age_garden
```

We specify the  $\gamma$ -profile of the MDCEV model specification where a single  $\alpha_1$  is estimated for the numeraire alternative and all non-numeraire alternatives are fixed at zero by setting `model = "gamma"`. We use maximum likelihood estimation by setting `algorithm = "MLE"`.

The syntax for the model is the following:

```
mdcev_mle <- mdcev(~ age_garden,
                  data = data_model,
                  model = "gamma",
                  algorithm = "MLE",
                  print_iterations = FALSE)
```

```
#> Using MLE to estimate KT model
```

Setting `print_iterations = TRUE` will print out intermediate iteration results as the model converges.

The output of the function can be accessed by calling `summary`.

```
summary(mdcev_mle)

#> Model run using rmdcev for R, version 1.2.0
#> Estimation method      : MLE
#> Model type             : gamma specification
#> Number of classes      : 1
#> Number of individuals  : 200
#> Number of non-numeraire alts : 17
#> Estimated parameters    : 36
#> LL                     : -5119.11
#> AIC                    : 10310.21
#> BIC                    : 10428.95
#> Standard errors calculated using : Delta method
#> Exit of MLE            : successful convergence
#> Time taken (hh:mm:ss)   : 00:00:0.53
#>
#> Average consumption of non-numeraire alternatives:
#>      beach      birding      camping      cycling      fish
#>      6.70      12.75      2.60      7.89      4.00
#>      garden      golf      hiking      hunt_birds      hunt_large
#>      23.18      5.42      41.62      0.58      1.03
#>      hunt_trap hunt_waterfowl      motor_land      motor_water      photo
#>      0.80      0.24      5.92      3.53      11.00
#>      ski_cross      ski_down
#>      3.12      1.85
#>
#> Parameter estimates -----
#>      Estimate Std.err z.stat
#> psi_birding      -0.762  0.113 -6.75
#> psi_camping      -0.534  0.115 -4.64
#> psi_cycling      -0.455  0.110 -4.13
#> psi_fish         -0.162  0.116 -1.39
#> psi_garden       -0.537  0.176 -3.05
#> psi_golf         0.553  0.112  4.94
#> psi_hiking       -0.038  0.107 -0.36
#> psi_hunt_birds   -1.034  0.194 -5.33
#> psi_hunt_large   -0.234  0.160 -1.46
#> psi_hunt_trap    -1.280  0.208 -6.16
#> psi_hunt_waterfowl -0.886  0.254 -3.49
#> psi_motor_land    0.119  0.126  0.94
#> psi_motor_water   0.458  0.115  3.98
#> psi_photo        0.011  0.105  0.11
#> psi_ski_cross    -1.164  0.122 -9.54
#> psi_ski_down     0.229  0.134  1.71
```

```
#> psi_age_garden      0.513    0.155    3.31
#> gamma_beach         8.665    1.457    5.95
#> gamma_birding       22.363    4.944    4.52
#> gamma_camping        7.546    1.482    5.09
#> gamma_cycling       16.180    3.115    5.19
#> gamma_fish          11.830    2.276    5.20
#> gamma_garden        17.762    2.711    6.55
#> gamma_golf          11.080    2.393    4.63
#> gamma_hiking        17.470    2.873    6.08
#> gamma_hunt_birds     9.667    3.686    2.62
#> gamma_hunt_large    12.563    3.590    3.50
#> gamma_hunt_trap     12.727    5.663    2.25
#> gamma_hunt_waterfowl 7.735    4.165    1.86
#> gamma_motor_land    16.273    4.008    4.06
#> gamma_motor_water   11.245    2.351    4.78
#> gamma_photo         14.475    2.634    5.50
#> gamma_ski_cross     10.362    2.387    4.34
#> gamma_ski_down       9.056    2.405    3.77
#> alpha_num           0.667    0.008   83.43
#> scale               0.607    0.027   22.47
#> Note: All non-numeraire alpha's fixed to 0.
```

The summary includes overall model and estimation information and the parameter estimates. All parameters have been transformed to their original form.<sup>5</sup> Interpreting the parameter estimates of KT models directly is challenging due to the non-linearities implied by the utility function and the partial confounding of  $\alpha_k$  and  $\gamma_k$  parameters (see Bhat (2008) for a in-depth discussion). Examining the  $\psi_k$  parameters first which represent the marginal utility when consumption is zero, we can see that relative to the beach recreation activity (i.e. the omitted reference category), hunting and trapping and cross country skiing have the largest negative ASCs suggesting these activities are less preferred starting from zero consumption levels. The interaction parameter between age and gardening is positive and significant suggesting that older people gain a higher utility from gardening compared to younger people. Because all non-numeraire  $\alpha$  parameters are fixed at zero, the  $\gamma_k$  parameters can be interpreted as capturing satiation and these satiation effects are lowest for the activities with the highest  $\gamma_k$  parameter values such as birding, cycling, and motorized land vehicles. The  $\alpha_1$  is estimated to be less than 1 which also implies satiation in the numeraire good. Bhat (2008); Lloyd-Smith et al. (2019) provide empirical applications of this model.

In the next example, we estimate the  $\alpha$ -profile of the MDCEV utility function by changing the model argument to "alpha".

```
mdcev_mle <- mdcev(~ age_garden,
                  data = data_model,
                  model = "alpha",
                  algorithm = "MLE",
                  print_iterations = FALSE)

summary(mdcev_mle)

#> Model run using rmdcev for R, version 1.2.0
#> Estimation method      : MLE
#> Model type              : alpha specification
#> Number of classes      : 1
#> Number of individuals  : 200
#> Number of non-numeraire alts : 17
#> Estimated parameters    : 36
#> LL                      : -5354.33
#> AIC                     : 10780.67
#> BIC                     : 10899.41
#> Standard errors calculated using : Delta method
#> Exit of MLE             : successful convergence
#> Time taken (hh:mm:ss)   : 00:00:0.59
#>
#> Average consumption of non-numeraire alternatives:
```

<sup>5</sup> $\gamma_k = \exp(\gamma_k^*)$ ,  $\alpha_1 = \exp(\alpha_1^*) / (1 + \exp(\alpha_1^*))$ , and  $\sigma = \exp(\sigma^*)$ , where  $\gamma_k^*$ ,  $\alpha_1^*$ , and  $\sigma^*$  are estimated but the transformed parameters are returned to users.

```
#>      beach      birding      camping      cycling      fish
#>      6.70      12.75      2.60      7.89      4.00
#>      garden      golf      hiking      hunt_birds      hunt_large
#>      23.18      5.42      41.62      0.58      1.03
#>      hunt_trap hunt_waterfowl      motor_land      motor_water      photo
#>      0.80      0.24      5.92      3.53      11.00
#>      ski_cross      ski_down
#>      3.12      1.85
```

```
#> Parameter estimates -----
```

```
#>      Estimate Std.err z.stat
#> psi_birding      -0.820  0.115 -7.13
#> psi_camping      -0.582  0.117 -4.97
#> psi_cycling      -0.500  0.111 -4.51
#> psi_fish         -0.208  0.117 -1.77
#> psi_garden       -0.480  0.176 -2.73
#> psi_golf          0.492  0.114  4.32
#> psi_hiking        0.127  0.109  1.17
#> psi_hunt_birds    -1.121  0.199 -5.63
#> psi_hunt_large    -0.309  0.164 -1.88
#> psi_hunt_trap     -1.359  0.213 -6.38
#> psi_hunt_waterfowl -0.976  0.261 -3.74
#> psi_motor_land     0.040  0.129  0.31
#> psi_motor_water    0.396  0.117  3.39
#> psi_photo         -0.030  0.105 -0.29
#> psi_ski_cross     -1.229  0.125 -9.83
#> psi_ski_down       0.158  0.138  1.14
#> psi_age_garden     0.494  0.156  3.17
#> alpha_num         0.658  0.008 82.21
#> alpha_beach        0.593  0.040 14.83
#> alpha_birding      0.720  0.038 18.94
#> alpha_camping      0.596  0.049 12.16
#> alpha_cycling      0.700  0.039 17.94
#> alpha_fish         0.660  0.043 15.34
#> alpha_garden       0.647  0.030 21.55
#> alpha_golf         0.669  0.045 14.87
#> alpha_hiking       0.595  0.030 19.82
#> alpha_hunt_birds    0.665  0.090  7.39
#> alpha_hunt_large    0.701  0.068 10.31
#> alpha_hunt_trap     0.710  0.094  7.55
#> alpha_hunt_waterfowl 0.652  0.132  4.94
#> alpha_motor_land    0.721  0.048 15.02
#> alpha_motor_water    0.663  0.047 14.12
#> alpha_photo        0.680  0.037 18.37
#> alpha_ski_cross     0.661  0.051 12.97
#> alpha_ski_down     0.658  0.060 10.96
#> scale              0.602  0.034 17.71
```

```
#> Note: All non-numeraire gamma's fixed to 1.
```

Estimating alternative-specific  $\alpha_k$  parameters and fixing all the non-numeraire  $\gamma$  parameters at 1, allows us to see the heterogeneity in  $\alpha_k$  parameters across recreation activities.

The hybrid model specification of the MDCEV model where a single  $\alpha$  is estimated for the numeraire and non-numeraire alternatives can be estimated by setting `model = "hybrid"` as the next example demonstrates.

```
mdcev_mle <- mdcev(~ age_garden,
  data = data_model,
  model = "hybrid",
  algorithm = "MLE",
  print_iterations = FALSE)
```

```
#> Using MLE to estimate KT model
```

```
summary(mdcev_mle)
```

```

#> Model run using rmdcev for R, version 1.2.0
#> Estimation method      : MLE
#> Model type              : hybrid specification
#> Number of classes      : 1
#> Number of individuals   : 200
#> Number of non-numeraire alts : 17
#> Estimated parameters    : 36
#> LL                      : -5230.91
#> AIC                     : 10533.81
#> BIC                     : 10652.55
#> Standard errors calculated using : Delta method
#> Exit of MLE             : successful convergence
#> Time taken (hh:mm:ss)   : 00:00:0.64
#>
#> Average consumption of non-numeraire alternatives:
#>      beach      birding      camping      cycling      fish
#>      6.70      12.75      2.60      7.89      4.00
#>      garden      golf      hiking      hunt_birds      hunt_large
#>      23.18      5.42      41.62      0.58      1.03
#>      hunt_trap hunt_waterfowl      motor_land      motor_water      photo
#>      0.80      0.24      5.92      3.53      11.00
#>      ski_cross      ski_down
#>      3.12      1.85
#>
#> Parameter estimates -----
#>      Estimate Std.err z.stat
#> psi_birding      -0.783  0.081 -9.67
#> psi_camping      -0.570  0.082 -6.95
#> psi_cycling      -0.487  0.078 -6.25
#> psi_fish         -0.206  0.083 -2.48
#> psi_garden       -0.580  0.128 -4.53
#> psi_golf         0.565  0.080  7.06
#> psi_hiking       -0.285  0.076 -3.74
#> psi_hunt_birds   -0.832  0.137 -6.07
#> psi_hunt_large   -0.095  0.113 -0.84
#> psi_hunt_trap    -1.029  0.146 -7.05
#> psi_hunt_waterfowl -0.524  0.178 -2.94
#> psi_motor_land   0.172  0.090  1.91
#> psi_motor_water  0.449  0.082  5.48
#> psi_photo       -0.102  0.074 -1.38
#> psi_ski_cross    -1.112  0.087 -12.78
#> psi_ski_down     0.346  0.095  3.64
#> psi_age_garden   0.312  0.112  2.79
#> gamma_beach     2.197  0.445  4.94
#> gamma_birding    5.721  1.484  3.85
#> gamma_camping    2.668  0.649  4.11
#> gamma_cycling    5.742  1.306  4.40
#> gamma_fish       4.162  1.008  4.13
#> gamma_garden     4.780  0.910  5.25
#> gamma_golf       3.446  0.873  3.95
#> gamma_hiking     3.313  0.719  4.61
#> gamma_hunt_birds 3.701  1.696  2.18
#> gamma_hunt_large 5.533  1.922  2.88
#> gamma_hunt_trap  4.583  2.434  1.88
#> gamma_hunt_waterfowl 3.268  2.054  1.59
#> gamma_motor_land 5.691  1.642  3.47
#> gamma_motor_water 3.940  1.011  3.90
#> gamma_photo      4.725  1.012  4.67
#> gamma_ski_cross  3.593  0.994  3.61
#> gamma_ski_down   3.264  1.026  3.18
#> alpha           0.648  0.005 129.53
#> scale           0.431  0.014 30.78
#> Note: Alpha parameter is equal for all alternatives.

```

The same number of parameters are estimated in all three models and the log-likelihood is highest for the  $\gamma$ -profile specification. The ease of estimating different MDCEV model specifications can be used to compare models quickly and help the analyst pick their preferred specification for each empirical application.

We can also estimate the KT-EE specification by changing the formula call and the model call to "kt\_ee".

```
kt_mle <- mdcev(~ age_garden | 0 | 0,
               data = data_model,
               model = "kt_ee",
               algorithm = "MLE",
               print_iterations = FALSE)

summary(kt_mle)

#> Model run using rmdcev for R, version 1.2.0
#> Estimation method           : MLE
#> Model type                   : kt_ee specification
#> Number of classes           : 1
#> Number of individuals       : 200
#> Number of non-numeraire alts : 17
#> Estimated parameters        : 20
#> LL                          : -5360.46
#> AIC                         : 10760.93
#> BIC                         : 10826.89
#> Standard errors calculated using : Delta method
#> Exit of MLE                  : successful convergence
#> Time taken (hh:mm:ss)       : 00:00:0.3
#>
#> Average consumption of non-numeraire alternatives:
#>      beach      birding      camping      cycling      fish
#>      6.70      12.75      2.60      7.89      4.00
#>      garden      golf      hiking      hunt_birds      hunt_large
#>      23.18      5.42      41.62      0.58      1.03
#>      hunt_trap hunt_waterfowl      motor_land      motor_water      photo
#>      0.80      0.24      5.92      3.53      11.00
#>      ski_cross      ski_down
#>      3.12      1.85
#>
#> Parameter estimates -----
#>
#>      Estimate Std.err z.stat
#> psi_age_garden      0.396  0.110  3.60
#> gamma_beach      10.546  1.083  9.74
#> gamma_birding     22.269  2.484  8.96
#> gamma_camping     16.201  1.778  9.11
#> gamma_cycling     16.238  1.743  9.32
#> gamma_fish        12.238  1.359  9.01
#> gamma_garden      16.644  2.166  7.68
#> gamma_golf         6.236  0.699  8.92
#> gamma_hiking      11.910  1.322  9.01
#> gamma_hunt_birds   25.821  4.428  5.83
#> gamma_hunt_large   13.798  2.019  6.83
#> gamma_hunt_trap    32.817  6.094  5.39
#> gamma_hunt_waterfowl 24.675  5.567  4.43
#> gamma_motor_land   10.400  1.281  8.12
#> gamma_motor_water   7.117  0.812  8.76
#> gamma_photo        11.153  1.183  9.43
#> gamma_ski_cross    28.684  3.201  8.96
#> gamma_ski_down      8.403  1.065  7.89
#> alpha_num          0.475  0.007 67.93
#> scale              0.713  0.025 28.54
```

This model does not include ASCs in the  $\psi_{ik}$  parameters due to concerns about weak complementarity.

## Estimating KT using Bayesian techniques

The exact same models can be fit using Bayesian estimation by changing the algorithm call to "Bayes". Bayesian estimation is implemented using the Stan programming language (Carpenter et al., 2017). The Bayesian framework requires careful choice of priors for the parameters, especially in data sparse contexts. The specific prior distributions for the fixed parameter specifications is presented below. The user has the ability to change the standard deviation and shape of these priors through these options in the `mdcev` function:

- `prior_psi_sd` standard deviation for normal prior with mean 0.
- `prior_phi_sd` standard deviation for normal prior with mean 0.
- `prior_gamma_sd` standard deviation for half-normal prior with mean 1.
- `prior_alpha_shape` shape parameter for beta distribution.
- `prior_scale_sd` standard deviation for half-normal prior with mean 0.

For the random parameter model specifications, the priors for the means of all random parameters follow a normal distribution with mean 0 on the unconstrained space.

There are also a number of further options for Bayesian estimation. For example, the number of iterations (`n_iterations`), number of chains (`n_chains`), and number of cores (`n_cores`) for parallel implementation of the chains can also be chosen. The full set of options for Bayesian estimation are presented below.

- `random_parameters` The form of the covariance matrix for the parameters. Options are
  - 'fixed' for no random parameters,
  - 'uncorr' for uncorrelated random parameters, or
  - 'corr' for correlated random parameters.
- `n_iterations` The number of iterations to use in Bayesian estimation. The default is for the number of iterations to be split evenly between warmup and posterior draws. The number of warmup draws can be directly controlled using the `warmup` argument (see `rstan::sampling`)
- `n_chains` The number of independent Markov chains in Bayesian estimation.
- `n_cores` The number of cores used to execute the Markov chains in parallel in Bayesian estimation. Can set using `options(mc.cores = parallel::detectCores())`.
- `max_tree_depth` <http://mc-stan.org/misc/warnings.html#maximum-treedepth-exceeded>
- `adapt_delta` <http://mc-stan.org/misc/warnings.html#divergent-transitions-after-warmup>
- `lkj_shape_prior` Prior for Cholesky matrix for correlated random parameters.

In this example, we estimate the  $\gamma$ -profile of the MDCEV specification using Bayesian techniques. We set the number of iterations to 200 and use 4 independent chains across 4 cores.

```
mdcev_bayes <- mdcev(~ age_garden,
  data = data_model,
  model = "gamma",
  algorithm = "Bayes",
  n_iterations = 200,
  n_chains = 4,
  n_cores = 4,
  print_iterations = FALSE)
```

The output of the function can be accessed by calling `summary`.

```
summary(mdcev_bayes)

#> Model run using rmdcev for R, version 1.2.0
#> Estimation method           : Bayes
#> Model type                   : gamma specification
#> Number of classes           : 1
#> Number of individuals       : 200
#> Number of non-numeraire alts : 17
#> Estimated parameters        : 36
#> LL                          : -5137.69
#> Number of chains            : 4
#> Number of warmup draws per chain : 100
#> Total post-warmup sample     : 400
```

```

#> Time taken (hh:mm:ss)          : 00:00:40.03
#>
#> Average consumption of non-numeraire alternatives:
#>      beach      birding      camping      cycling      fish
#>      6.70      12.75      2.60      7.89      4.00
#>      garden      golf      hiking      hunt_birds      hunt_large
#>      23.18      5.42      41.62      0.58      1.03
#>      hunt_trap hunt_waterfowl      motor_land      motor_water      photo
#>      0.80      0.24      5.92      3.53      11.00
#>      ski_cross      ski_down
#>      3.12      1.85
#>
#> Parameter estimates -----
#>      Estimate Std.dev z.stat n_eff Rhat
#> psi_birding      -0.786  0.118 -6.69  366 1.00
#> psi_camping      -0.572  0.125 -4.57  272 1.00
#> psi_cycling      -0.489  0.118 -4.13  446 1.00
#> psi_fish         -0.196  0.135 -1.45  467 1.00
#> psi_garden       -0.553  0.180 -3.08  327 1.00
#> psi_golf         0.515  0.117  4.41  397 0.99
#> psi_hiking       0.017  0.108  0.16  531 1.00
#> psi_hunt_birds   -1.149  0.197 -5.82  478 0.99
#> psi_hunt_large   -0.322  0.168 -1.92  538 0.99
#> psi_hunt_trap    -1.413  0.223 -6.35  430 1.00
#> psi_hunt_waterfowl -1.039  0.273 -3.80  305 1.00
#> psi_motor_land   0.059  0.130  0.45  284 1.00
#> psi_motor_water  0.415  0.122  3.40  296 1.00
#> psi_photo        0.000  0.104  0.00  372 1.00
#> psi_ski_cross    -1.225  0.130 -9.39  427 1.00
#> psi_ski_down     0.157  0.147  1.06  495 1.00
#> psi_age_garden   0.551  0.157  3.50  407 1.00
#> gamma_beach      7.959  1.343  5.93  396 1.00
#> gamma_birding    18.287  3.367  5.43  428 1.00
#> gamma_camping     7.249  1.450  5.00  568 0.99
#> gamma_cycling    14.410  2.731  5.28  681 0.99
#> gamma_fish       11.064  2.172  5.09  414 1.00
#> gamma_garden     15.744  2.268  6.94  503 1.00
#> gamma_golf       10.126  2.093  4.84  370 1.00
#> gamma_hiking     15.319  2.387  6.42  566 0.99
#> gamma_hunt_birds  9.488  3.404  2.79  214 1.00
#> gamma_hunt_large  11.626  3.189  3.65  254 1.01
#> gamma_hunt_trap   11.425  4.214  2.71  230 1.02
#> gamma_hunt_waterfowl 8.754  3.991  2.19  200 1.02
#> gamma_motor_land  14.303  3.295  4.34  657 1.00
#> gamma_motor_water 10.573  2.139  4.94  430 1.00
#> gamma_photo      13.169  2.239  5.88  514 1.00
#> gamma_ski_cross   9.894  2.311  4.28  410 1.00
#> gamma_ski_down    8.767  2.335  3.75  433 1.00
#> alpha_num        0.668  0.008 88.60  397 1.00
#> scale            0.651  0.029 22.72  280 1.00
#> Note: All non-numeraire alpha's fixed to 0.
#> Note from Rstan: 'For each parameter, n_eff is a crude measure of effective sample size, and Rhat is the p

```

Comparing these parameter values to the maximum likelihood estimates of the  $\gamma$ -profile MDCEV specification, the values are quite similar. As the data set is rather small with only 200 individuals, the priors play a role in reducing the estimates closer to 1 for the  $\gamma_k$ , but this role will lessen in larger data applications.

One benefit of using the Bayesian approach is that one can take advantage of the postestimation commands, interactive diagnostics, and posterior analysis in **rstan**, **bayesplot** (Gabry et al., 2019), and **shinystan** (Muth et al., 2018). For example, the effective sample size reports the estimated number of independent draws from the posterior distribution for each parameter (Stan Development Team, 2019). The interested reader is referred to these packages for additional details.



## Estimating LC-KT models

In this example, we estimate a LC-KT model using the **Recreation** data. We set the number of classes equal to 2 and we use data on 500 individuals. We would like to include the university, ageindex, and urban in the membership equation and we include them in the formula interface. The constant for the membership equation is included automatically. The LC model is automatically estimated as long as the prespecified number of classes (n\_classes) is set greater than 1. The scale parameters are fixed at 1 using fixed\_scale1 = 1.

```
data_model <- mdcev.data(data_rec, subset = id <= 500,
                        id.var = "id",
                        alt.var = "alt",
                        choice = "quant")

mdcev_lc <- mdcev(~ 0 | university + ageindex + urban,
                 data = data_model,
                 n_classes = 2,
                 model = "gamma",
                 fixed_scale1 = 1,
                 algorithm = "MLE",
                 print_iterations = FALSE)

#> Error in chol.default(-H) :
#>   the leading minor of order 43 is not positive definite

summary(mdcev_lc)

#> Warning in sqrt(diag(cov_mat)): NaNs produced

#> Model run using rmdcev for R, version 1.2.0
#> Estimation method           : MLE
#> Model type                   : gamma specification
#> Number of classes           : 2
#> Number of individuals       : 500
#> Number of non-numeraire alts : 17
#> Estimated parameters        : 72
#> LL                           : -12073.55
#> AIC                         : 24291.1
#> BIC                         : 24594.55
#> Standard errors calculated using : Delta method
#> Exit of MLE                  : successful convergence
#> Time taken (hh:mm:ss)       : 00:00:7.5
#>
#> Average consumption of non-numeraire alternatives:
#>      beach      birding      camping      cycling      fish
#>      6.47      12.27      2.33      9.03      3.84
#>      garden      golf      hiking      hunt_birds      hunt_large
#>      21.97      4.44      39.52      0.64      1.14
#>      hunt_trap hunt_waterfowl      motor_land      motor_water      photo
#>      0.84      0.27      4.48      3.03      9.20
#>      ski_cross      ski_down
#>      2.36      1.49
#>
#>
#> Class average probabilities:
#> class1 class2
#>  0.85  0.15
#> Parameter estimates -----
#>
#>      Estimate Std.err z.stat
#> class1.psi_birding      -1.290  0.133 -9.70
#> class2.psi_birding      -1.154  0.133 -8.68
#> class1.psi_camping      -0.977  0.125 -7.81
#> class2.psi_camping      -0.580  0.161 -3.60
#> class1.psi_cycling      -0.777  0.113 -6.88
#> class2.psi_cycling      -0.949  0.139 -6.83
#> class1.psi_fish        -1.117  0.121 -9.23
```

```

#> class2.psi_fish          1.446  0.585  2.47
#> class1.psi_garden        0.025  1.110  0.02
#> class2.psi_garden       -0.073 10.998 -0.01
#> class1.psi_golf         -0.113  0.617 -0.18
#> class2.psi_golf          0.689  0.159  4.33
#> class1.psi_hiking        0.473  0.154  3.07
#> class2.psi_hiking        0.186  0.119  1.57
#> class1.psi_hunt_birds   -4.017  0.143 -28.09
#> class2.psi_hunt_birds    0.591  0.151  3.91
#> class1.psi_hunt_large   -4.259  0.343 -12.42
#> class2.psi_hunt_large    1.609  0.313  5.14
#> class1.psi_hunt_trap   -10.286  0.332 -30.98
#> class2.psi_hunt_trap    -0.247  0.301 -0.82
#> class1.psi_hunt_waterfowl -3.723  0.281 -13.25
#> class2.psi_hunt_waterfowl -0.019  0.295 -0.06
#> class1.psi_motor_land   -0.740  0.271 -2.73
#> class2.psi_motor_land    1.210  0.305  3.97
#> class1.psi_motor_water  -0.456  0.325 -1.40
#> class2.psi_motor_water    1.624  0.323  5.03
#> class1.psi_photo        -0.130  0.336 -0.39
#> class2.psi_photo        -0.724  0.288 -2.51
#> class1.psi_ski_cross    -1.783  0.278 -6.41
#> class2.psi_ski_cross    -1.438  0.344 -4.18
#> class1.psi_ski_down     -0.417  0.348 -1.20
#> class2.psi_ski_down     -0.107  0.356 -0.30
#> class1.gamma_beach      3.814  0.479  7.96
#> class2.gamma_beach      7.426  1.366  5.44
#> class1.gamma_birding    11.328  1.781  6.36
#> class2.gamma_birding     6.982  1.060  6.59
#> class1.gamma_camping     3.558  0.720  4.94
#> class2.gamma_camping     8.956  1.062  8.43
#> class1.gamma_cycling     9.239  1.649  5.60
#> class2.gamma_cycling    13.108  1.671  7.84
#> class1.gamma_fish        5.386  4.401  1.22
#> class2.gamma_fish        3.472  3.955  0.88
#> class1.gamma_garden      8.718    NaN    NaN
#> class2.gamma_garden      8.816  8.146  1.08
#> class1.gamma_golf        7.088  1.533  4.62
#> class2.gamma_golf        3.805  0.778  4.89
#> class1.gamma_hiking      6.160  0.829  7.43
#> class2.gamma_hiking     12.122  2.280  5.32
#> class1.gamma_hunt_birds   3.228  0.651  4.96
#> class2.gamma_hunt_birds   2.553  0.811  3.15
#> class1.gamma_hunt_large  11.305  5.693  1.99
#> class2.gamma_hunt_large   3.133  1.152  2.72
#> class1.gamma_hunt_trap    0.803  0.337  2.38
#> class2.gamma_hunt_trap    5.184  1.465  3.54
#> class1.gamma_hunt_waterfowl 4.955  1.522  3.26
#> class2.gamma_hunt_waterfowl 3.796  1.348  2.82
#> class1.gamma_motor_land   5.093  1.436  3.55
#> class2.gamma_motor_land   8.906  2.587  3.44
#> class1.gamma_motor_water  3.010  0.851  3.54
#> class2.gamma_motor_water  5.688  2.028  2.80
#> class1.gamma_photo        6.103  2.394  2.55
#> class2.gamma_photo     10.199  2.794  3.65
#> class1.gamma_ski_cross    4.948  1.339  3.70
#> class2.gamma_ski_cross    7.200  3.109  2.32
#> class1.gamma_ski_down     3.997  1.669  2.39
#> class2.gamma_ski_down     4.691  2.170  2.16
#> class1.alpha_num         0.672  0.009 74.68
#> class2.alpha_num         0.686  0.022 31.17
#> class2.(Intercept)       -1.100  0.479 -2.30
#> class2.university        -0.406  0.338 -1.20
#> class2.ageindex          0.186  0.364  0.51

```

```
#> class2.urban          -0.837   0.350  -2.39
#> Note: Scale parameter fixed to 1.
#> Note: All non-numeraire alpha's fixed to 0.
#> Note: The membership equation parameters for class 1 are normalized to 0.
```

In this LC example, we assume that there are two types of people that have different preferences for recreation. The probability of class assignment depends on unobserved factors and the three sociodemographic factors included in the membership equation with only urban having a statistically significant effect on class probability. People living in urban areas are less likely to be in class 2. The summary output reports the average class probabilities as being 32% for class 1 and 68% for class 2. The  $\psi$  parameters across classes are similar although there are some noticeable differences such as the hunting and trapping preferences. The  $\gamma$  parameters, on the other hand, show that satiation between classes is quite different. [Sobhani et al. \(2013\)](#); [Kuriyama et al. \(2010\)](#) provide empirical applications of these models.

If initial parameter are not provided, the default is to use slightly adjusted parameter estimates of the MDCEV model as starting values when estimating the LC-MDCEV model to assist speed and convergence issues.<sup>6</sup> The MDCEV model output can be accessed from `mdcev_lc[["mdcev_fit"]]` object for comparison.

### Estimating RP-KT models

Random parameter models require defining and parameterizing the variance covariance matrix. For uncorrelated random parameters, the diagonal elements of the variance covariance matrix are estimated and the off-diagonal elements are assumed to be zero. For correlated random parameters, the variance covariance matrix is fully estimated and can be parameterized in many ways. The **rmdcev** package defines the variance covariance matrix in terms of Cholesky factors of the correlation matrix and a vector of standard deviations for numerical stability. Thus the variance covariance matrix is specified as

$$\Sigma = \text{diag}(\tau) \times LL^T \times \text{diag}(\tau), \quad (17)$$

where  $\tau$  is a vector of standard deviations, and  $L$  is the cholesky factors of the correlation matrix.

In this example, we estimate an uncorrelated random parameters  $\gamma$ -specification of the MDCEV model without any  $\psi_k$  parameters. We set the argument `random_parameters = "uncorr"` to indicate that uncorrelated random parameters will be estimated. As noted earlier, all random parameters follow a normal distribution. We change the `psi_ascs = 0` to omit the ASCs in the  $\psi_k$  parameters.

```
data_model <- mdcev.data(data_rec, subset = id <= 200,
                        id.var = "id",
                        alt.var = "alt",
                        choice = "quant")

mdcev_rp <- mdcev(~ 0,
                 data = data_model,
                 model = "gamma",
                 algorithm = "Bayes",
                 n_chains = 4,
                 psi_ascs = 0,
                 fixed_scale1 = 1,
                 n_iterations = 200,
                 random_parameters = "uncorr",
                 print_iterations = FALSE)

summary(mdcev_rp)

#> Model run using rmdcev for R, version 1.2.0
#> Estimation method          : Bayes
#> Model type                  : gamma specification
#> Number of classes           : 1
#> Number of individuals       : 200
#> Number of non-numeraire alts : 17
#> Estimated parameters        : 36
```

<sup>6</sup>In particular, the estimated  $\psi_k$  and  $\gamma_k$  parameters from the MDCEV model are randomly adjusted by 0.02.

```

#> LL : -5362.51
#> Random parameters : uncorrelated random parameters
#> Number of chains : 4
#> Number of warmup draws per chain : 100
#> Total post-warmup sample : 400
#> Time taken (hh:mm:ss) : 00:01:57.22
#>
#> Average consumption of non-numeraire alternatives:
#>      beach      birding      camping      cycling      fish
#>      6.70      12.75      2.60      7.89      4.00
#>      garden      golf      hiking      hunt_birds      hunt_large
#>      23.18      5.42      41.62      0.58      1.03
#>      hunt_trap hunt_waterfowl      motor_land      motor_water      photo
#>      0.80      0.24      5.92      3.53      11.00
#>      ski_cross      ski_down
#>      3.12      1.85
#>
#> Parameter estimates -----
#>
#>      Estimate Std.dev z.stat n_eff Rhat
#> gamma_beach      5.661 1.075 5.27 427 1.00
#> gamma_birding     8.249 2.272 3.63 479 1.00
#> gamma_camping     3.795 0.732 5.18 404 0.99
#> gamma_cycling     7.980 1.480 5.39 463 1.00
#> gamma_fish        7.020 1.649 4.26 400 1.00
#> gamma_garden     12.635 2.048 6.17 407 1.00
#> gamma_golf        6.506 1.605 4.05 362 1.00
#> gamma_hiking     15.176 2.188 6.94 455 1.00
#> gamma_hunt_birds  4.135 1.962 2.11 534 0.99
#> gamma_hunt_large  7.206 2.399 3.00 351 1.00
#> gamma_hunt_trap   5.689 3.753 1.52 595 1.00
#> gamma_hunt_waterfowl 4.052 2.916 1.39 706 0.99
#> gamma_motor_land  8.446 2.323 3.64 344 1.01
#> gamma_motor_water 6.649 1.644 4.05 358 1.00
#> gamma_photo       8.786 1.588 5.53 357 1.00
#> gamma_ski_cross   3.343 0.796 4.20 334 1.01
#> gamma_ski_down    4.933 1.400 3.52 537 1.00
#> alpha_num         0.725 0.007 98.82 396 1.00
#> sd.gamma_beach    1.263 0.223 5.65 365 1.00
#> sd.gamma_birding  1.800 0.609 2.95 242 1.00
#> sd.gamma_camping  1.283 0.248 5.18 585 0.99
#> sd.gamma_cycling  1.307 0.256 5.12 331 1.00
#> sd.gamma_fish     1.262 0.236 5.36 531 1.00
#> sd.gamma_garden   1.297 0.242 5.36 361 0.99
#> sd.gamma_golf     1.565 0.475 3.30 176 1.00
#> sd.gamma_hiking   1.334 0.265 5.04 300 0.99
#> sd.gamma_hunt_birds 1.767 1.185 1.49 556 0.99
#> sd.gamma_hunt_large 1.415 0.422 3.35 776 0.99
#> sd.gamma_hunt_trap 2.472 2.752 0.90 260 1.02
#> sd.gamma_hunt_waterfowl 2.622 2.382 1.10 383 1.00
#> sd.gamma_motor_land 1.544 0.522 2.96 395 1.00
#> sd.gamma_motor_water 1.407 0.321 4.38 422 1.00
#> sd.gamma_photo    1.275 0.253 5.05 248 1.00
#> sd.gamma_ski_cross 1.509 0.472 3.20 285 1.00
#> sd.gamma_ski_down 1.500 0.448 3.35 381 1.01
#> sd.alpha_num      0.515 0.010 50.79 365 1.00
#> Note: Scale parameter fixed to 1.
#> Note: All non-numeraire alpha's fixed to 0.
#> Note from Rstan: 'For each parameter, n_eff is a crude measure of effective sample size, and Rhat is the p

```

The results show the means of the random parameters followed by the estimated standard deviations. The standard deviations that are estimated to be different from zero suggest there is heterogeneity in preference parameters. The correlated random parameters specification can be estimated by setting `random_parameters = "corr"`. [Bhat and Sen \(2006\)](#) provide an empirical application of this type of model.

## Computational and estimation issues

KT models are notoriously tricky to estimate relative to standard discrete choice models. This section provides some guidance for estimating these models and common convergence issues:

- **Starting values:** Model parameter estimates can be sensitive to starting values, especially the more complex LC-KT specification. Users should use several different initial parameter values for model estimation to ensure robust results and a global maxima is found rather than a local maxima. The default behaviour for LC-KT models is to use KT parameters as starting values. In practice the author has found this to be quite effective at finding global maxima. However, users are encouraged to use random starting values as a robustness check.
- **Identification issues:** Depending on the model specification and included variables the model may not be properly identified. If you receive an error such as `Error in chol.default(-H) : the leading minor of order 9 is not positive definite`, this usually suggests an identification issue. Users should double check all variables included in the model are appropriate. One solution is to start with a simpler model first and then slowly add variables to help locate any problematic variables.
- **Parameter estimates near boundaries:** Interpret models with parameter estimates that are near the boundaries (e.g.  $\alpha$  close to 1) with caution. Users are recommended to re-estimate the model with starting values far from this boundary.
- **Bayesian estimation:** For models estimated using Bayesian estimation, users should consult the **rstan** User Guide for additional guidance on model estimation options and postestimation checks (Stan Development Team, 2019). Additional information is available by typing `help(rstan)`.

## Simulating KT demand and welfare scenarios

The **rmdecv** package includes simulation functions for calculating welfare measures and forecasting demand under alternative policy scenarios. The overall approach used for simulation is first introduced and then code examples are given.

### Overview of simulation steps

Once the model parameters are estimated, there are two steps to simulation in KT models. In the first step we draw simulated values for the unobserved heterogeneity term ( $\epsilon$ ) using Monte Carlo techniques. The second step uses these error draws, the previously estimated model parameters, and the underlying data to calculate Marshallian demands for forecasting or Hicksian demands for welfare analysis. These two steps are described below.

#### Step 1: simulating unobserved heterogeneity

Monte Carlo simulation techniques can be employed to draw simulated values of the unobserved heterogeneity ( $\epsilon$ ) using either unconditional or conditional draws.

1. Unconditional error draws: draw from the entire distribution of unobserved heterogeneity using the following formula

$$\epsilon_k = -\log(-\log(\text{draw}(0,1))) * \sigma, \quad (18)$$

where  $\text{draw}(0,1)$  is a draw between 0 and 1 and  $\sigma$  is the scale parameter. **rmdecv** allows errors to be drawn using uniform draws or the Modified Latin Hypercube Sampling algorithm (Hess et al., 2006).

2. Conditional error draws: draw errors terms to reflect behaviour and dependent on whether alternative is consumed or not (von Haefen, 2003; von Haefen et al., 2004):
  - If  $x_k > 0$ , set  $\epsilon_k = (V_1 - V_k)/\sigma$  for the MDCEV specifications where  $V_1$  and  $V_k$  depend on the model specification as detailed above. If using the environmental economics KT model specification ("kt\_ee"), set  $\epsilon_k = g_k(\cdot)$  from Equation (12).
  - If  $x_k = 0$ ,  $\epsilon_k < (V_1 - V_k)/\sigma$  and simulate  $\epsilon_k$  from the truncated type I extreme value distribution such that

$$\epsilon_k = -\log(-\log(\text{draw}(0,1) * \exp(-\exp(\frac{V_1 - V_k}{\sigma})))) * \sigma \text{ for the MDCEV specifications, or } \epsilon_k = -\log(-\log(\text{draw}(0,1) * \exp(-\exp(\frac{V_1 - V_k}{\sigma})))) * \sigma \text{ for the EE specifications} \quad (19)$$

In the conditional error draw approach, we normalize  $\varepsilon_1 = 0$ .

The main differences between these two error draw approaches is that in the conditional approach, errors are drawn such that the model perfectly predicts the observed consumption patterns in the baseline state (von Haefen and Phaneuf, 2005). The conditional approach uses observed behaviour by individuals to characterize unobserved heterogeneity and can be useful for scenario simulation as the baseline matches observed behavior. This is especially true if poor in-sample behavioral predictions is found using the unconditional approach (von Haefen, 2003). The unconditional approach draws all errors based on distributional assumptions and is necessary for out-of-sample forecasting. If the model correctly specifies the data generating process, the sample means of the conditional and unconditional approaches should converge in expectation. Another difference between the two approaches is that the unconditional approach uses more computation time as there is a need to calculate consumption patterns in the baseline state as well as simulate the entire distribution of unobserved heterogeneity.

(von Haefen, 2003)

### Step 2: Calculating welfare measures and demand forecasts

With the error draws in hand, the second step is to simulate demand or welfare changes. Compared to welfare measures in discrete choice models, welfare calculation in KT models is more challenging because of the two KT conditions in Equation (2). For a given policy scenario, a priori, we do not know which alternatives have a positive or zero consumption level. **rmdcev** implements the Pinjari and Bhat (2011) efficient demand forecasting routine for simulating demand behaviour for MDCEV models which relies on calculating Marshallian demands. For welfare calculations, we need to calculate the expenditure function in Equation (3) which relies on Hicksian demands. These are calculated using the approach described by Lloyd-Smith (2018) and the **rmdcev** extends these approaches to the environmental economics KT model specifications. The demand and welfare simulation approaches share a lot of commonalities and thus only the approach used for welfare calculations are fully described in the appendix. The specific steps for demand simulation is explained in-depth in Pinjari and Bhat (2011) and the interested reader is encouraged to read Section 4 of the paper for the exact details.

## Welfare analysis

In **rmdcev**, the functions for welfare and demand simulation have been divided into 3 steps to allow users to parallelize operations as necessary.

We first estimate the model using **mdcev** and we set `std_errors = "mvn"` to generate multivariate normal draws.

```
mdcev_mle <- mdcev(~0,
  data = data_model,
  model = "hybrid",
  algorithm = "MLE",
  std_errors = "mvn",
  print_iterations = FALSE)

#> Using MLE to estimate KT model
```

**1. Define policy scenarios** In the first step, we define the number of alternative policy scenarios to use in simulation and then specify changes to the  $\psi$  variables and prices of alternatives. The `CreateBlankPolicies` function has been created to easily set-up the required lists for the simulation. These policies can then be manually edited according to the specific policy scenario. For prices, **rmdcev** is set up to accept additive changes in prices that impact all individuals the same. For the  $\psi$  and  $\phi$  variable changes, the package is set up to accept any new values for these variables. Depending on the number of individuals and number of policies, the generated policies list can be quite large. If the user is only interested in assessing price changes, then you can use `price_change_only = TRUE` which ensures duplicate  $\psi$  and  $\phi$  data is not created.

In this example, we are interested in two separate policies. The first policy increases the costs of all recreation activities by \$1 and the second policy increases the cost of all four hunting activities by \$10. The policy set-up for these two scenarios is demonstrated below.

```
nalts <- mdcev_mle$stan_data[["J"]]
npols <- 2

policies<- CreateBlankPolicies(npols = npols,
  model = mdcev_mle,
```

```

      price_change_only = TRUE)

policies$price_p[[1]] <- c(0, rep(1, nalts))
policies$price_p[[2]][10:13] <- rep(10, 4)

```

For policy scenarios that involve changes in the  $\psi$  or  $\phi$  variables, the user can change the `dat_psi` or `dat_phi` list of the `policies` object. For example, the following code will increase the value of the third  $\psi$  variable by 20% in policy scenario 1.

```
policies$dat_psi[[1]][3] <- policies$dat_psi[[1]][3]*1.2
```

**2. Prepare simulation data** The second step is to combine the parameter estimates, data, and policy scenarios into a data format for simulation. The `PrepareSimulationData` function uses the model fit and the user defined policy scenarios to create this specific data format. This function separates the output into individual-specific data (`df_indiv`), data common to all individuals (`df_common`), and simulation options (`sim_options`).

```
df_sim <- PrepareSimulationData(mdcev_mle, policies)
```

**3. Simulate MDCEV model** The third step is to simulate the policy scenario using the formatted data and the `mdcev.sim` function. The specific steps for the simulation algorithms are described in Appendix A. The user chooses the type of error draws (unconditional or conditional as described above), the number of error draws, and whether to simulate the demand or welfare changes.

```

welfare <- mdcev.sim(df_sim$df_indiv,
                    df_common = df_sim$df_common,
                    sim_options = df_sim$sim_options,
                    cond_err = 1,
                    nerrs = 25,
                    sim_type = "welfare")

#> Using hybrid approach in simulation...

#>
#> 3.00e+05simulations finished in0.07minutes.(74074per second)

summary(welfare)

#> # A tibble: 2 x 5
#>   policy   mean std.dev `ci_lo2.5%` `ci_hi97.5%`
#>   <chr>   <dbl>   <dbl>      <dbl>      <dbl>
#> 1 policy1 -125.    0.116    -125.      -125.
#> 2 policy2 -20.2    0.418    -20.9      -19.4

```

The output of the `mdcev.sim` for welfare analysis is an object of class `mdcev.sim` which contains a list of matrices where each element of the list is for an individual and the matrix consists of rows for each policy scenario and columns for each parameter simulation.

The `summary` function computes summary statistics across all individuals. For example, the average welfare change for a \$1 daily increase in all recreation costs is -\$125.

The reason these last two steps are separate is to allow users to parallelize the simulation step as the last step can be computationally intensive. The number of simulations is a multiplicative function of the number of individuals, number of policies, number of parameter estimate simulations, and the number of error draws ( $I \times npols \times nsims \times nerrs$ ). Even for modestly sized data, the total number of simulations can easily reach well into the millions or billions. All simulations are conducted at the individual level which allows the user to easily parallelize the `mdcev.sim` function using the **parallel** package or similar packages.

## Demand forecasting

This section demonstrates the demand forecasting capabilities of **rmdcev**. Please refer to the previous section for an overview of the three steps to simulation.

```
policies <- CreateBlankPolicies(npols = 2, model = mdcev_mle)
```



```

policies$price_p[[1]] <- c(0, rep(1, nalts))
policies$price_p[[2]][10:13] <- rep(10, 4)

df_sim <- PrepareSimulationData(mdcev_mle, policies)

demand <- mdcev.sim(df_sim$df_indiv,
  df_common = df_sim$df_common,
  sim_options = df_sim$sim_options,
  cond_err = 1,
  nerrs = 25,
  sim_type = "demand")

#> Using hybrid approach in simulation...

#>
#> 5.40e+06simulations finished in0.07minutes.(1291866per second)

summary(demand)

#> # A tibble: 36 x 6
#> # Groups:   policy [2]
#>   policy alt      mean std.dev `ci_lo2.5%` `ci_hi97.5%`
#>   <chr> <int>    <dbl>    <dbl>    <dbl>    <dbl>
#> 1 policy1      0 69085.    6.07    69073.    69095.
#> 2 policy1      1   6.22    0.02     6.18     6.24
#> 3 policy1      2  11.0    0.05    11.0    11.1
#> 4 policy1      3   2.32    0.02     2.29     2.34
#> 5 policy1      4   7.19    0.03     7.13     7.24
#> 6 policy1      5   3.75    0.01     3.72     3.77
#> 7 policy1      6  20.5    0.07    20.4    20.6
#> 8 policy1      7   5.29    0.01     5.28     5.3
#> 9 policy1      8  35.6    0.11    35.4    35.7
#> 10 policy1     9   0.54     0     0.54     0.55
#> # ... with 26 more rows

```

The output of the demand simulation a `mdcev.sim` object with a list of  $I$  elements, one for each individual. Within each element there are `nsim` lists each containing a matrix of demands. The rows of the matrix are for each policy scenario and the columns represent each alternative. The summary function computes summary statistics.

## Generating simulated data

The `rmdcev` package has the capability to simulate KT data. Simulated KT data can be easily created for model assessment and Monte Carlo analysis using the `GenerateMDCEVData` function. The following example will generate a simulated data set with 1,000 individuals, 10 non-numeraire alternatives, and particular parameter values.

```

model = "gamma"
nobs = 1000
nalts = 10
sim.data <- GenerateMDCEVData(model = model,
  nobs = nobs,
  nalts = nalts,
  psi_j_parms = c(-5, 0.5, 2), # alternative-specific variables
  psi_i_parms = c(-1.5, 3, -2, 1, 2), # individual-specific variables
  gamma_parms = stats::runif(nalts, 1, 10),
  alpha_parms = 0.5,
  scale_parms = 1)

#> Sorting data by id.var then alt...

#> Checking data...

#> Data is good

```

Next, we can estimate the model using maximum likelihood techniques to recover the parameter estimates.

```
mdcev_mle <- mdcev(formula = ~ b1 + b2 + b3 + b4 + b5 + b6 + b7 + b8,
  data = sim.data$data,
  model = model,
  psi_ascs = 0,
  algorithm = "MLE",
  print_iterations = FALSE)
```

## Conclusions

The **mdcev** package implements several Kuhn-Tucker model specifications including MDCEV with heterogeneity that can be continuous (i.e. random parameters) or discrete (i.e. latent classes). Models can be estimated using maximum likelihood or Bayesian techniques. This paper demonstrates the use of the package to estimate several model specifications and to derive demand forecasts and welfare implications of policy scenarios. To my knowledge, there is no other available statistical package that can estimate welfare implications of policy scenarios using MDCEV models. I hope that the publication of **mdcev** will make KT modeling available to a wider audience.

## Appendix A: Specific steps for simulating KT models

Welfare and demand simulation follow similar approaches and this section details the welfare simulation approach. There are two algorithms that differ depending on the model specification. If a single  $\alpha$  parameter is estimated (i.e. model = "hybrid" or "hybrid0"), then we can use the hybrid approach to welfare simulation. If there are heterogeneous  $\alpha$  parameters (i.e. model = "gamma", "alpha", or "kt\_ee"), then we can use the general approach to welfare simulation. The hybrid approach is less computationally intensive and provides an exact analytical solution but the general approach can be used with all utility specifications. The specific steps for both algorithms are described below. Additional details are provided in [Lloyd-Smith \(2018\)](#).

### Steps in algorithm for hybrid-profile MDCEV utility specifications

**Step 0:** Assume that only the numeraire alternative is chosen and let the number of chosen alternatives equal one ( $M=1$ ).

**Step 1:** Using the data, model parameters, and either conditional or unconditional simulated error term draws, calculate the price-normalized baseline utility values ( $\psi_k/p_k$ ) for all alternatives. Sort the  $K$  alternatives in the descending order of their price-normalized baseline utility values. Note that the numeraire alternative is in the first place. Go to step 2.

**Step 2:** Compute the value of  $\lambda^E$  using the following equation:

$$\frac{1}{\lambda^E} = \left[ \frac{\alpha \bar{U} + \sum_{m=2}^M \gamma_m \psi_m}{\sum_{m=2}^M \gamma_m \psi_m \left( \frac{p_m}{\psi_m} \right)^{\frac{\alpha}{\alpha-1}} + \psi_1 \left( \frac{p_1}{\psi_1} \right)^{\frac{\alpha}{\alpha-1}}} \right]^{\frac{\alpha-1}{\alpha}}. \quad (20)$$

Go to step 3.

**Step 3:** If  $\frac{1}{\lambda^E} > \frac{\psi_{M+1}}{p_{M+1}}$ , go to step 4. Else if  $\frac{1}{\lambda^E} < \frac{\psi_{M+1}}{p_{M+1}}$ , set  $M = M + 1$ . If  $M < K$ , go back to step 2. If  $M = K$ , go to step 4.

**Step 4:** Compute the optimal Hicksian consumption levels for the first  $I$  alternatives in the above descending order using the following equations

$$x_1 = \left( \frac{p_1}{\lambda^E \psi_1} \right)^{\frac{1}{\alpha_1-1}}, \text{ and} \quad (21)$$

$$x_m = \left[ \left( \frac{p_m}{\lambda^E \psi_m} \right)^{\frac{1}{\alpha_m-1}} - 1 \right] \gamma_m, \text{ if } x_m > 0. \quad (22)$$

Set the remaining alternative consumption levels to zero and stop.

### Steps in algorithm for general utility specifications

In this context, there is no closed-form expressions for  $\lambda^E$  and we need to conduct a numerical bisection routine. The following routine describes the approach for the MDCEV utility specifications. The approach used for the KT-EE specification is omitted due to space, but the overall strategy is the same with the only differences being the definitions for utility functions and optimal demands. Let  $\hat{\lambda}^E$

and  $\hat{U}$  be estimates of  $\lambda^E$  and  $U$  and let  $tol_\lambda$  and  $tol_U$  be the tolerance levels for estimating  $\lambda^E$  and  $U$  that can be arbitrarily small. The algorithm works as follows:

**Step 0:** Assume that only the numeraire is chosen and let the number of chosen alternatives equal one ( $M=1$ ).

**Step 1:** Using the data, model parameters, and either conditional or unconditional simulated error term draws, calculate the price-normalized baseline utility values ( $\psi_k/p_k$ ) for all alternatives. Sort the  $K$  alternatives in the descending order of their price-normalized baseline utility values. Note that the numeraire is in the first place. Go to step 2.

**Step 2:** Let  $\frac{1}{\lambda^E} = \frac{\psi_{M+1}}{p_{M+1}}$  and substitute  $\hat{\lambda}^E$  into the following equation to obtain an estimate of  $\hat{U}$ .

$$\bar{U} = \sum_{m=2}^M \frac{\gamma_m}{\alpha_m} \psi_m \left[ \left( \frac{p_m}{\lambda^E \psi_m} \right)^{\frac{\alpha_m}{\alpha_m-1}} - 1 \right] + \frac{\psi_1}{\alpha_1} \left( \frac{p_1}{\lambda^E \psi_1} \right)^{\frac{\alpha_1}{\alpha_1-1}}. \quad (23)$$

**Step 3:** If  $\hat{U} < \bar{U}$ , go to step 4. Else, if  $\hat{U} \geq \bar{U}$ , set  $\frac{1}{\lambda_l^E} = \frac{\psi_{M+1}}{p_{M+1}}$  and  $\frac{1}{\lambda_u^E} = \frac{\psi_M}{p_M}$ . Go to step 5.

**Step 4:** Set  $M = M + 1$ . If  $M < K$ , go to step 2. Else if  $M = K$ , set  $\frac{1}{\lambda_l^E} = 0$  and  $\frac{1}{\lambda_u^E} = \frac{\psi_K}{p_K}$ . Go to step 5.

**Step 5:** Let  $\hat{\lambda}^E = (\lambda_l^E + \lambda_u^E)/2$  and substitute  $\hat{\lambda}^E$  into the equation of step 2 to obtain an estimate of  $\hat{U}$ . Go to step 6.

**Step 6:** If  $|\lambda_l^E - \lambda_u^E| \leq tol_\lambda$  or  $|\hat{U} - \bar{U}| \leq tol_U$ , go to step 7. Else if  $\hat{U} < \bar{U}$ , update  $\lambda_u^E = (\lambda_l^E + \lambda_u^E)/2$  and go to step 5. Else if  $\hat{U} > \bar{U}$ , update  $\lambda_l^E = (\lambda_l^E + \lambda_u^E)/2$  and go to step 5.

**Step 7:** Compute the optimal Hicksian consumption levels for the first  $M$  alternatives in the above descending order using Equation (21). Set the remaining alternative consumption levels to zero and stop.

## Bibliography

- C. Bhat and A. Pinjari. Multiple discrete-continuous choice models: A reflective analysis and a prospective view. In S. Hess and A. Daly, editors, *Handbook of Choice Modelling*, chapter 19, pages 427–454. Edward Elgar Publishing, 2014. [p1]
- C. R. Bhat. The multiple discrete-continuous extreme value (MDCEV) model: Role of utility function parameters, identification considerations, and model extensions. *Transportation Research Part B: Methodological*, 42(3):274–303, 2008. ISSN 0191-2615. doi: 10.1016/j.trb.2007.06.002. [p1, 3, 4, 5, 10]
- C. R. Bhat and S. Sen. Household vehicle type holdings and usage: an application of the multiple discrete-continuous extreme value (MDCEV) model. *Transportation Research Part B: Methodological*, 40(1):35–53, Jan. 2006. ISSN 0191-2615. doi: 10.1016/j.trb.2005.01.003. URL <http://www.sciencedirect.com/science/article/pii/S0191261505000093>. [p19]
- B. Carpenter, A. Gelman, M. Hoffman, D. Lee, B. Goodrich, M. Betancourt, M. Brubaker, J. Guo, P. Li, and A. Riddell. Stan: A Probabilistic Programming Language. *Journal of Statistical Software, Articles*, 76(1):1–32, 2017. ISSN 1548-7660. doi: 10.18637/jss.v076.i01. [p2, 14]
- Y. Croissant. *mlogit: Multinomial Logit Models*, 2019. URL <https://CRAN.R-project.org/package=mlogit>. R package version 1.0-1. [p1]
- Federal, Provincial, and Territorial Governments of Canada. 2012 Canadian Nature Survey: Awareness, participation, and expenditures in nature-based recreation, conservation, and subsistence activities. Technical report, Canadian Councils of Resource Ministers, Ottawa, ON, 2014. [p5]
- J. Gabry, D. Simpson, A. Vehtari, M. Betancourt, and A. Gelman. Visualization in bayesian workflow. *J. R. Stat. Soc. A*, 182:389–402, 2019. doi: 10.1111/rssa.12378. [p15]
- J. A. Herriges, C. L. Kling, and D. J. Phaneuf. What’s the use? welfare estimates from revealed preference models when weak complementarity does not hold. *Journal of Environmental Economics and Management*, 47(1):55–70, 2004. Publisher: Elsevier. [p4]
- S. Hess and D. Palma. Apollo: A flexible, powerful and customisable freeware package for choice model estimation and application. *Journal of Choice Modelling*, page 100170, June 2019. ISSN 1755-5345. doi: 10.1016/j.jocm.2019.100170. [p1, 2]

- S. Hess, K. E. Train, and J. W. Polak. On the use of a Modified Latin Hypercube Sampling (MLHS) method in the estimation of a Mixed Logit Model for vehicle choice. *Transportation Research Part B: Methodological*, 40(2):147–163, Feb. 2006. ISSN 0191-2615. doi: 10.1016/j.trb.2004.10.005. [p20]
- A. Jasra, C. C. Holmes, and D. A. Stephens. Markov chain monte carlo methods and the label switching problem in bayesian mixture modeling. *Statist. Sci.*, 20(1):50–67, 02 2005. doi: 10.1214/088342305000000016. URL <https://doi.org/10.1214/088342305000000016>. [p5]
- K. Kuriyama, M. Hanemann, and J. Hilger. A latent segmentation approach to a Kuhn-Tucker model: An application to recreation demand. *Journal of Environmental Economics and Management*, 60(3): 209–220, Nov. 2010. ISSN 0095-0696. doi: 10.1016/j.jeem.2010.05.005. [p5, 18]
- P. Lloyd-Smith. A new approach to calculating welfare measures in kuhn-tucker demand models. *Journal of Choice Modelling*, 26:19 – 27, 2018. ISSN 1755-5345. doi: <https://doi.org/10.1016/j.jocm.2017.12.002>. [p21, 24]
- P. Lloyd-Smith, J. K. Abbott, W. Adamowicz, and D. Willard. Decoupling the Value of Leisure Time from Labor Market Returns in Travel Cost Models. *Journal of the Association of Environmental and Resource Economists*, 6(2):215–242, Jan. 2019. ISSN 2333-5955. doi: 10.1086/701760. URL <https://www.journals.uchicago.edu/doi/abs/10.1086/701760>. [p10]
- K. Mäler. *Environmental Economics: A Theoretical Inquiry*. Johns Hopkins University Press for Resources for the Future, 1974. [p3]
- C. Muth, Z. O. Jonah, and Gabry. User-friendly bayesian regression modeling: A tutorial with rstanarm and shinystan. *The Quantitative Methods for Psychology*, 14(2):99–119, 2018. doi: 10.20982/tqmp.14.2.p099. [p15]
- A. R. Pinjari and C. R. Bhat. Computationally efficient forecasting procedures for Kuhn-Tucker consumer demand model systems: Application to residential energy consumption analysis. Technical report, 2011. [p21]
- M. Sarrias and R. Daziano. Multinomial Logit Models with Continuous and Discrete Individual Heterogeneity in R: The gmnL Package. *Journal of Statistical Software*, 79(1):1–46, July 2017. ISSN 1548-7660. doi: 10.18637/jss.v079.i02. URL <https://www.jstatsoft.org/index.php/jss/article/view/v079i02>. [p1]
- A. Sobhani, N. Eluru, and A. Faghih-Imani. A latent segmentation based multiple discrete continuous extreme value model. *Transportation Research Part B: Methodological*, 58:154–169, Dec. 2013. ISSN 0191-2615. doi: 10.1016/j.trb.2013.07.009. [p5, 18]
- Stan Development Team. RStan: the R interface to Stan. R package version 2.19, 2019. URL <https://mc-stan.org/rstan>. [p15, 20]
- K. Train. *Discrete Choice Methods with Simulation*. Cambridge University Press, 2009. [p5]
- R. H. von Haefen. Incorporating observed choice into the construction of welfare measures from random utility models. *Journal of Environmental Economics and Management*, 45(2):145–165, Mar. 2003. ISSN 0095-0696. doi: 10.1016/S0095-0696(02)00047-5. [p20, 21]
- R. H. von Haefen and D. J. Phaneuf. Kuhn-Tucker Demand System Approaches to Non-Market Valuation. In R. Scarpa and A. Alberini, editors, *Applications of Simulation Methods in Environmental and Resource Economics*, pages 135–157. Springer Netherlands, Dordrecht, 2005. ISBN 978-1-4020-3684-2. [p1, 2, 4, 5, 21]
- R. H. von Haefen, D. J. Phaneuf, and G. R. Parsons. Estimation and Welfare Analysis with Large Demand Systems. *Journal of Business & Economic Statistics*, 22(2):194–205, 2004. ISSN 0735-0015. [p1, 20]
- A. Zeileis and Y. Croissant. Extended Model Formulas in R: Multiple Parts and Multiple Responses. *Journal of Statistical Software*, 34(1):1–13, Apr. 2010. ISSN 1548-7660. doi: 10.18637/jss.v034.i01. [p8]

Patrick Lloyd-Smith

University of Saskatchewan

Department of Agricultural and Resource Economics Room 3D34, Agriculture Building 51 Campus Drive  
Saskatoon, SK S7N 5A8 Canada

[patrick.lloydsmith@usask.ca](mailto:patrick.lloydsmith@usask.ca)