

# The Journal

Volume 16/2, June 2024

---

A peer-reviewed, open-access publication of the  
R Foundation for Statistical Computing

## Contents

Editorial . . . . .	3
---------------------	---

### Contributed Research Articles

LUCIDus: An R Package For Implementing Latent Unknown Clustering By Integrating Multi-omics Data (LUCID) With Phenotypic Traits . . . . .	4
Current State and Prospects of R-Packages for the Design of Experiments . . . . .	27
BCClong: An R Package for Bayesian Consensus Clustering for Multiple Longitudinal Features. . . . .	41
reslr: An R Package for Relative Sea Level Modelling . . . . .	61
Rfssa: An R Package for Functional Singular Spectrum Analysis . . . . .	82
mcmcSupply: An R Package for Estimating Contraceptive Method Market Supply Shares . . . . .	99
Pomdp: A Computational Infrastructure for Partially Observable Markov Decision Processes . . . . .	116
pencal: an R Package for the Dynamic Prediction of Survival with Many Longitudinal Predictors . . . . .	134
clarify: Simulation-Based Inference for Regression Models. . . . .	154
FuzzySimRes: Epistemic Bootstrap – an Efficient Tool for Statistical Inference Based on Imprecise Data . . . . .	175
Fast and Flexible Search for Homologous Biological Sequences with DECIPHER v3. .	191

### News and Notes

Bioconductor Notes, June 2024 . . . . .	201
Changes on CRAN . . . . .	203
Changes in R . . . . .	205
R Foundation News . . . . .	211

---

The R Journal is a peer-reviewed publication of the R Foundation for Statistical Computing. Communications regarding this publication should be addressed to the editors. All articles are licensed under the Creative Commons Attribution 4.0 International license (CC BY 4.0, <http://creativecommons.org/licenses/by/4.0/>).

Prospective authors will find detailed and up-to-date submission instructions on the Journal's homepage.

**Editor-in-Chief:**

Mark van der Loo, Statistics Netherlands and Leiden University, Netherlands

**Executive editors:**

Simon Urbanek, University of Auckland, New Zealand

Catherine Hurley, Maynooth University, Ireland

Rob Hyndman, Monash University, Australia

Emi Tanaka, Australian National University, Australia

**Technical editors:**

Mitchell O'Hara-Wild, Monash University, Australia

**R Journal Homepage:**

<http://journal.r-project.org/>

**Email of editors and editorial board:**

[r-journal@R-project.org](mailto:r-journal@R-project.org)

The R Journal is indexed/abstracted by EBSCO, DOAJ, Thomson Reuters.

# Editorial

by *Mark P.J. van der Loo*

On behalf of the editorial board, I am pleased to present Volume 16 Issue 1 of the R Journal.

## In this issue

News from CRAN, the R Foundation, the R Core team, and Bioconductor are included in this issue.

This issue features 10 contributed research articles the majority of which relate to R packages on a diverse range of topics. All packages are available on CRAN. Supplementary material with fully reproducible code is available for download from the Journal website. Topics covered in this issue are the following.

### Reviews

- Current State and Prospects of R-Packages for the Design of Experiments

### Modeling and Inference

- BCCLong: An R Package for Bayesian Consensus Clustering for Multiple Longitudinal Features
- pencal: an R Package for the Dynamic Prediction of Survival with Many Longitudinal Predictors
- Rfssa: An R Package for Functional Singular Spectrum Analysis

### Computational Methods

- clarify: Simulation-Based Inference for Regression Models
- FuzzySimRes: Epistemic Bootstrap –an Efficient Tool for Statistical Inference Based on Imprecise Data
- Pomdp: A Computational Infrastructure for Partially Observable Markov Decision Processes

### Applications

- Fast and Flexible Search for Homologous Biological Sequences with DECIPHER v3
- reslr: An R Package for Relative Sea Level Modelling
- mcmSupply: An R Package for Estimating Contraceptive Method Market Supply Shares
- LUCIDus: An R Package For Implementing Latent Unknown Clustering By Integrating Multi-omics Data (LUCID) With Phenotypic Traits

*Mark P.J. van der Loo  
Statistics Netherlands and Leiden University*

<https://journal.r-project.org>  
[r-journal@r-project.org](mailto:r-journal@r-project.org)

# LUCIDus: An R Package For Implementing Latent Unknown Clustering By Integrating Multi-omics Data (LUCID) With Phenotypic Traits

by Yinqi Zhao, Qiran Jia, Jesse A. Goodrich, David V. Conti

**Abstract** Many studies are leveraging current technologies to obtain multiple omics measurements on the same individuals. These measurements are usually cross-sectional, and methods developed and commonly used focus on omic integration at a single time point. More unique, and a growing area of interest, are studies that leverage biology or the temporal sequence of measurements to relate long-term exposures or germline genetics to intermediate measures capturing transitional processes that ultimately result in an outcome. In this context, we have previously introduced an integrative model named Latent Unknown Clustering by Integrating multi-omics Data (LUCID) aiming to distinguish unique effects of environmental exposures or germline genetics and informative omic effects while jointly estimating subgroups of individuals relevant to the outcome of interest. This multiple omics analysis consists of early integration (concatenation of omic layers to estimate common subgroups); intermediate integration (omic-layer-specific estimation of subgroups that are all related to the outcome); and late integration (omic-layer-specific estimation of subgroups that are then interrelated by a priori structures). In this article, we introduce LUCIDus version 3, an R package to implement the LUCID model. We review the statistical background of the model and introduce the workflow of LUCIDus, including model fitting, model selection, interpretation, inference, and prediction. Throughout, we use a realistic but simulated dataset based on an ongoing study, the Human Early Life Exposome Study (HELIX), to illustrate the workflow.

## 1 Introduction

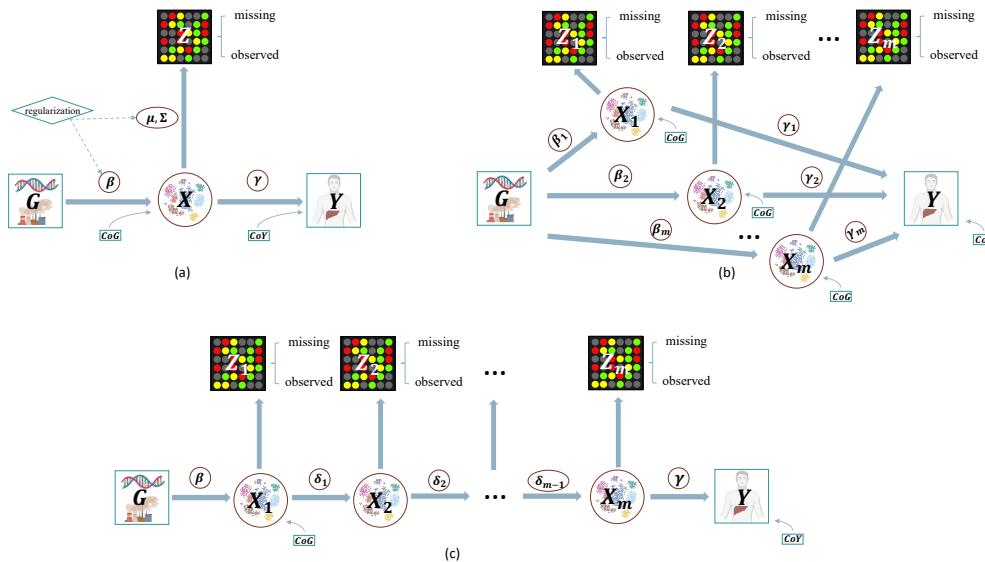
The rapid advancement in high-throughput technology has made it possible to measure multiple omics (referred to as “multi-omics”) data on the same person in many cohort studies. These datasets include DNA genome sequences (Goodwin et al., 2016), RNA expression data (Ozsolak and Milos, 2011), metabolites in biofluids (Beger, 2013), proteins in cells or tissues (Aslam et al., 2017), as well as chemical exposures in the environment (Wild, 2005). For example, the Human Early-Life Exposome project (HELIX) measured molecular omics signatures including RNA expression data, metabolites, plasma proteins, etc. from 1300 children at the age of 6-11 in six European countries (Vrijheid et al., 2014). Guided by biology or the temporal sequence of measurements, these studies often share a common structure that relates germline genetics or environmental exposures to intermediate factors capturing transitional processes that ultimately result in an outcome. While these suspected causal pathways can be measured with current omic technology, the analysis often focuses on cross-sectional relationships using cluster analysis or models focused on feature selection. This is in part due to the analytic and computational challenges resulting from the complex structure of the data collected from various sources, limited sample sizes, and the high dimensionality of multi-omics information (Tini et al., 2019). Existing methods do not consider risk factors that precede the omic measurements and link those factors with disease or trait outcomes while using multi-omics data to characterize the underlying mechanism. Thus, there is a great need to develop and easily implement approaches leveraging these mediating relationships or latent structures for the integration of multi-omics data (Subramanian et al., 2020). Conventional analysis tools range from clustering approaches to variable selection approaches (González and Cáceres, 2019). Clustering is a fundamental method for analyzing both single omic and multi-omics data (Rappoport and Shamir, 2018). Clustering aims to divide observations into several groups (called clusters) so that observations within a group are similar while samples in different groups are dissimilar. Clustering has many important applications in the field of biological science and medicine, such as defining gene sets (Hejblum et al., 2015), identifying subtypes of cancer, and performing diagnostic predictions (Curtis et al., 2012; Khan et al., 2001), defining cell types in flow cytometry and scRNA-seq experiments (Chan et al., 2008; Hejblum et al., 2019; Prabhakaran et al., 2016), and estimating protein localization (Crook et al., 2018). Conventional clustering methods applied to multi-omics data either concatenate multiple datasets (often called early integration) or analyze each dataset independently (i.e., late integration). However, both approaches fall short of adequately capturing the variation across omic levels or reflecting the heterogeneity of the integrated datasets. To overcome these limitations, several specific integrated clustering

methods for analyzing multi-omics data are available in the community. Pierre-Jean et al. (2020) conducted a thorough review comparing 13 unsupervised methods for multi-omics data integration via clustering. Methods include SGCCA (Tenenhaus et al., 2014), SNF (Wang et al., 2014), and iCluster plus (Mo et al., 2013), etc. Building upon the foundation of Principal Component Analysis (PCA), Joint and Individual Variation Explained (JIVE) conducts integrated clustering of multi-omics data via a variance decomposition framework to summarize information across multiple omic layers (Lock et al., 2013). Within the Bayesian framework, Kirk et al. (2012) proposed an unsupervised approach, Bayesian correlated clustering, that simultaneously models each omic layer via a mixture model while considering the correlation between each model. In a similar fashion, Bayesian consensus clustering can flexibly describe the dependency and the heterogeneity of the multi-omics data (Lock and Dunson, 2013). An alternative approach to clustering includes dimension reduction or variable selection to reduce noise, improve model interpretability, and avoid problems with overfitting. For example, the least absolute shrinkage and selection operator (LASSO) induces sparsity in model coefficients based on  $L_1$  norm regularization (Tibshirani, 1996). An extension of LASSO, Group LASSO, allows for variable selection on both individual and pre-specified grouped variables (Yuan and Lin, 2006).  $L_1$  norm regularization is integrated into many clustering approaches to reduce data dimensionality, such as SGCCA and iCluster plus. Extensions of variable selection approaches exist for mediation as well and include both high-dimensional and latent variable approaches. HIMA first implements pre-screening of the high-dimensional mediators and then utilizes an LASSO-type penalty to obtain a sparse solution (Zhang et al., 2016), while BAMA is a Bayesian shrinkage approach that employs continuous shrinkage priors on the key coefficients to select active mediators with large effects only (Song et al., 2020). Both Albert et al. (2016) and Derkach et al. (2019) proposed statistical models within the causal mediation framework that summarize the information from high-dimensional multi-omics data into a latent variable to facilitate interpretation. In the context of linking multi-omics to health outcomes and identifying key omic features that drive the outcomes, Singh et al. (2019) proposed Data Integration Analysis for Biomarker discovery using Latent cOmponents (DIABLO), which extends the unsupervised sGCCA to a supervised framework. DIABLO achieves summarizing common information across different omic layers and conducting variable selection while discriminating the outcome of interest. Nonetheless, DIABLO fails to take into account the impact of risk factors related to the outcome and their interplay with multi-omics data. Finally, Peng et al. (2020) proposed a model called Latent Unknown Clustering by Integrating multi-omics Data (LUCID) that incorporates both clustering and variable selection to distinguish unique effects of germline genetics or environmental exposures and informative omic effects while jointly estimating subgroups of individuals relevant to the outcome of interest. LUCID is a novel 'quasi-mediation' approach that uses latent variable analysis to estimate subgroups of individuals characterized by key omic factors and with differential associations to the outcome and exposures. They demonstrated the performance of LUCID through extensive simulation studies and real data applications to highlight the integration of genomic, exposomic, and metabolomic data. **LUCIDus**, the R package to implement the LUCID model, provides an integrated clustering framework and has numerous downloads (around 19,000 times since it was first introduced according to **dlstats** (Yu, 2022)). It has also been applied in several environmental epidemiological studies (Jin et al., 2020; Stratakis et al., 2020; Matta et al., 2022). In this paper, we introduce **LUCIDus** version 3, a major update and enhancement from the original release (Jia et al., 2024). To account for the heterogeneity of the integrated datasets based on the study design, **LUCIDus** version 3 incorporates three integration strategies into the LUCID framework: (1) Early integration (Figure 1 (a)) concatenates all multi-omics data matrices into a single matrix prior to model estimation; (2) Intermediate integration (Figure 1 (b)) estimates latent clusters for each omic layer, and the resulting clusters are then modeled in parallel in a joint model linking exposures and the outcome; and (3) Late integration (Figure 1 (c)) in which each omic layer is used to estimate omic-specific clusters that are then linked via *a priori* relationships with each other and the exposure and outcome. **LUCIDus** version 3 also includes model selection, model visualization, and inference based on bootstrap resampling. It also incorporates an integrated imputation approach to deal with missingness in multi-omics data. The paper is organized as follows: we first briefly introduce the statistical model of LUCID; we then go through the details of functions in the **LUCIDus** package; we illustrate the workflow of fitting LUCID models by using a dataset simulated based on real cases from the HELIX study (Vrijheid et al., 2014; Maitre et al., 2022).

## 2 Model and software

### 2.1 Overview of the original LUCID model (early integration)

In the LUCID model, genomic/exposomic exposures  $G$ , other multi-omics data  $Z$ , and phenotype trait  $Y$  are integrated through a latent categorical variable  $X$ . To model the complex correlation structure between each layer of the multi-omics data, we propose three types of LUCID models based on



**Figure 1:** Three DAGs represent how three different types of the LUCID model integrate genetic/environmental exposures ( $G$ ), other multi-omics data ( $Z$ ), and the phenotype trait ( $Y$ ). (a) LUCID early integration; (b) LUCID in parallel; (c) LUCID in serial. The squares represent observed data, the circles represent unobserved latent variables (clusters) and model parameters, and the diamond refers to  $L_1$  penalty terms for regularization for (a).  $CoG$  and  $CoY$  represent covariates to be adjusted in the LUCID model. Missingness is allowed in multi-omics data.  $Z$  is divided into subsets of observations with complete measurements and observations with missingness. For (b) and (c),  $Z$  is partitioned into  $m$  layers.

different integration strategies: (1) LUCID early integration; (2) LUCID in parallel or intermediate integration; and (3) LUCID in serial or late integration. Figure 1 illustrates the joint relationship between  $G$ ,  $Z$ ,  $X$ , and  $Y$  for each LUCID model. We will discuss LUCID in parallel and LUCID in serial in detail under Sections [LUCID in Parallel](#) and [LUCID in Serial](#) as two extensions, and here we focus on LUCID early integration. Because  $X$  is an unobserved categorical variable, each category of  $X$  is interpreted as a latent cluster in the data, jointly defined by  $G$ ,  $Z$ , and  $Y$ . Let  $G$  be a  $N \times P$  matrix with columns representing genetic/environmental exposures and rows representing the observations;  $Z$  be a  $N \times M$  matrix of omics data (for example, gene expression data, DNA methylation profiles, and metabolomic data, etc.) and  $Y$  be a  $N$ -length vector of phenotype traits. We further assume  $G$ ,  $Z$ , and  $Y$  are measured through a prospective sampling procedure, so we do not model the distribution of  $G$ . All three measured components ( $G$ ,  $Z$ , and  $Y$ ) are linked by a latent variable  $X$  consisting of  $K$  categories. The directed acyclic graph (DAG) in Figure 1 (a) implies that the distributions of  $X$  given  $G$ ,  $Z$  given  $X$ , and  $Y$  given  $X$  are conditionally independent of each other. Let  $f(\cdot)$  denote the probability mass functions (PMFs) for categorical random variables or the probability density functions (PDFs) for continuous random variables. The joint log-likelihood of the LUCID model is constructed as:

$$\begin{aligned} \log L(\Theta) &= \sum_{i=1}^N \log f(Z_i, Y_i | G_i; \Theta) \\ &= \sum_{i=1}^N \log \sum_{j=1}^K f(X_i = j | G_i; \Theta) f(Z_i | X_i = j; \Theta) f(Y_i | X_i = j; \Theta) \end{aligned} \quad (1)$$

where  $\Theta$  is a generic notation for all parameters in the LUCID model. Since  $X$  is a discrete variable with  $K$  categories, we assume  $X$  follows a multinomial distribution conditioning on  $G$ , denoted by the softmax function  $S(\cdot)$ . We assume omics data  $Z$  follows a multivariate Gaussian distribution conditioning on  $X$ , denoted by  $\phi(Z|X = j; \mu_j, \Sigma_j)$ , where  $\mu_j$  and  $\Sigma_j$  are cluster-specific means and variance-covariance matrices for  $Z$ . For illustrative purposes, we assume  $Y$  is a continuous outcome following a univariate Gaussian distribution, denoted by  $\phi(Y|X = j; \gamma_j, \sigma_j^2)$  ( $\gamma_j$  and  $\sigma_j^2$  are also interpreted as the cluster-specific means and variances for the effect of  $X$  on  $Y$ ). Similar development for a binary outcome is detailed in (Peng et al., 2020). Because the latent cluster  $X$  is unobserved, the maximum likelihood estimates (MLE) of the parameters associated with the LUCID model based on Equation 1 are not readily estimated. Therefore, we use the Expectation-Maximization (EM) algorithm to obtain the MLE of model parameters. We define  $I(X_i = j)$  as an indicator function representing

that observation  $i$  belongs to latent cluster  $j$ . Then, the log-likelihood function of model parameters in Equation 1 becomes:

$$\log L(\Theta) = \sum_{i=1}^N \sum_{j=1}^K I(X_i = j) (\log S(X_i = j | G_i; \Theta) + \log \phi(Z_i | X_i = j; \Theta) + \log \phi(Y_i | X_i = j; \Theta)) \quad (2)$$

We denote the observed data as  $D = \{G, Z, Y\}$ , and define the responsibility  $r$  as the inclusion probability (IP) of belonging to the latent cluster  $j$  given observed data and current estimation of model parameters at iteration  $t$ , which is:

$$\begin{aligned} r_{ij}^{(t)} &= P(X_i = j | D; \Theta^{(t)}) \\ &= \frac{S(X_i = j | G_i; \beta^{(t)}) \phi(Z_i | X_i = j; \mu_j^{(t)}, \Sigma_j^{(t)}) \phi(Y_i | X_i = j; \gamma_j^{(t)}, \sigma_j^{2(t)})}{\sum_{j=1}^K S(X_i = j | G_i; \beta^{(t)}) \phi(Z_i | X_i = j; \mu_j^{(t)}, \Sigma_j^{(t)}) \phi(Y_i | X_i = j; \gamma_j^{(t)}, \sigma_j^{2(t)})} \end{aligned} \quad (3)$$

At each iteration  $t$ , in the E-step, we compute the expectation of the complete data likelihood, which is:

$$\begin{aligned} Q(\Theta | D, \Theta^{(t)}) &= \sum_{i=1}^N \sum_{j=1}^K r_{ij} \log S(X_i = j | G_i; \beta) + \sum_{i=1}^N \sum_{j=1}^K r_{ij} \log \phi(Z_i | X_i = j; \mu_j, \Sigma_j) \\ &\quad + \sum_{i=1}^N \sum_{j=1}^K r_{ij} \log \phi(Y_i | X_i = j; \gamma_j, \sigma_j^2) \end{aligned} \quad (4)$$

In the M-step, we update parameter estimates by maximizing Equation 4 in terms of  $\Theta$ , which results in the following:

$$\beta^{(t+1)} = \arg \max_{\beta} \sum_{i=1}^N \sum_{j=1}^K r_{ij}^{(t)} \log S(X_i = j | G_i; \beta_j) \quad (5)$$

$$\mu_j^{(t+1)} = \frac{\sum_{i=1}^N r_{ij}^{(t)} Z_i}{\sum_{i=1}^N r_{ij}^{(t)}} \quad (6)$$

$$\Sigma_j^{(t+1)} = \frac{\sum_{i=1}^N r_{ij}^{(t)} (Z_i - \mu_j^{(t+1)}) (Z_i - \mu_j^{(t+1)})^T}{\sum_{i=1}^N r_{ij}^{(t)}} \quad (7)$$

$$\gamma_j^{(t+1)} = \frac{\sum_{i=1}^N r_{ij}^{(t)} Y_i}{\sum_{i=1}^N r_{ij}^{(t)}} \quad (8)$$

$$\sigma_j^{2(t+1)} = \frac{\sum_{i=1}^N r_{ij}^{(t)} (Y_i - \gamma_j^{(t+1)})^2}{\sum_{i=1}^N r_{ij}^{(t)}} \quad (9)$$

Although maximization of  $\beta_j^{(t+1)}$  in Equation 5 does not have a closed-form solution, it is equivalent to fitting a multinomial logistic regression by treating  $r_{ij}$  as the outcome and  $G_i$  as the exposure. For **LUCIDus**, we include a more flexible geometric feature of latent clusters, such as volume, shape, and orientation determined by  $\Sigma_j$ . We use the parameterization of covariance matrices by means of an eigenvalue decomposition in the form below (Banfield and Raftery, 1993):

$$\Sigma_j = \lambda_j D_j A_j D_j^T \quad (10)$$

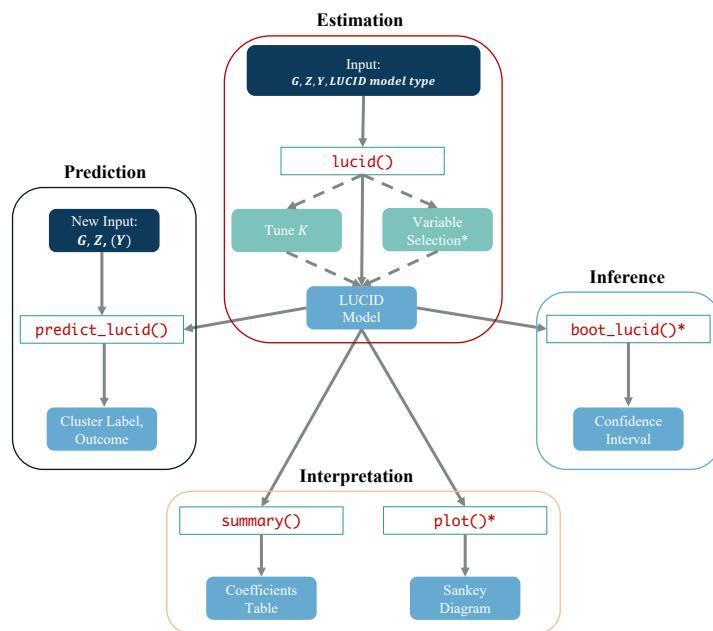
where  $\lambda_j$  is a scalar,  $D_j$  is the orthogonal matrix of eigenvectors, and  $A_j$  is a diagonal matrix whose values are proportional to eigenvalues. A detailed discussion of maximizing  $\Sigma_j$  parameterized by Equation 10 is provided by Celeux and Govaert (1995). Their algorithm is implemented in the R package **mclust** (Scrucca et al., 2016), and we leverage **mclust** to update  $\Sigma_j$  in the M-step at each iteration of the EM algorithm for LUCID.

## 2.2 General workflow of LUCIDus

The **LUCIDus** package includes five main functions and two auxiliary functions to implement the analysis framework based on LUCID. Brief descriptions of each function are listed in Table 1. The workflow of the **LUCIDus** package is shown in Figure 2. Below we describe the typical workflow

Main function	Description
<code>lucid()</code>	Main function to fit LUCID models, specified by giving integrated data, a distribution of the outcome, and the type of the LUCID model (early, parallel, serial). It also conducts model selection and variable selection for LUCID early integration and model selection for candidate models with different numbers of latent clusters for LUCID in parallel and LUCID in serial.
<code>summary()</code>	S3 method for LUCID. Create tables to summarize a LUCID model.
<code>plot()</code>	S3 method for LUCID. Visualize LUCID models through a Sankey diagram.
<code>boot_lucid()</code>	Derive confidence intervals based on bootstrap resampling.
<code>predict_lucid()</code>	Predict latent cluster assignment and outcome using integrated data.
Workhorse function	Description
<code>estimate_lucid()</code>	Fit a LUCID model to estimate latent clusters by using integrated data.
<code>tune_lucid()</code>	Fits a series of LUCID models over different combinations of tuning parameters and determines an optimal model with the minimum BIC.

**Table 1:** Functions in the **LUCIDus** package. `lucid()` calls `estimate_lucid()` and `tune_lucid()` in the backend. The two workhorse functions are not normally called directly, but they can be useful when a user wants to examine the model fitting process in more detail.



**Figure 2:** The workflow of the **LUCIDus** package. Dark blue nodes represent input data, light blue nodes represent output results. Green nodes and dashed arrows are optional steps for model estimation. Red texts correspond to 5 key functions in **LUCIDus**. Steps and functions marked with an asterisk currently work for LUCID early integration only.

of analyzing integrated data using the LUCID early integration model. The function `lucid()` is the primary function in the package, which fits a LUCID model based on an exposure matrix (argument  $G$ ), multi-omics data (argument  $Z$ ), outcome data (argument  $Y$ ), the number of latent clusters ( $K$ ; default is 2), the type of the LUCID model (argument `lucid_model`; here it is ‘early’ as we focus on the LUCID early integration model in this section), and the family of the outcome (argument `family`; default is ‘normal’). If a vector of  $K$  and/or  $L_1$  penalties are supplied, `lucid()` will automatically conduct model selection on the number of clusters  $K$ , select informative variables in  $G$  or  $Z$ , or both, and return a LUCID early integration model (an R object of class ‘`lucid_early`’) with optimal  $K$  and selected variables in  $G$  and/or  $Z$ . Several additional functions can then be applied to a fitted LUCID object. `summary()` summarizes the fitted LUCID model by producing summary tables of parameter estimation and interpretation. Visualization is performed via the `plot()` function to create a Sankey diagram showing the interplay among the three components ( $G$ ,  $Z$ , and  $Y$ ). In addition, statistical inference can be accomplished for LUCID by constructing confidence intervals (CIs) based on bootstrap resampling. This is achieved by the function `boot_lucid()`. Finally, predictions on the cluster assignment and the outcome can be obtained by calling the function `predict_lucid()`. In practice, it might not be necessary to implement the entire workflow above. For instance, if we have prior knowledge of the number of latent clusters  $K$ , model selection for the number of clusters can be skipped. If a given dataset has limited variables (for example, variables selected based on biological annotations), then variable selection may not be necessary. Note that for a LUCID in parallel model or a LUCID in serial model (the argument `lucid_model` is ‘parallel’ or ‘serial’, respectively), specification of  $K$  is different, and the features of variable selection, constructing confidence intervals (CIs) based on bootstrap resampling, and plotting are not yet available (see Sections [LUCID in Parallel](#) and [LUCID in Serial](#)). In the following sections, we show example code and demonstrate the usage of the LUCID early integration model using the simulated HELIX data.

### 3 Illustration

#### 3.1 Fitting a LUCID model by `lucid()`

To illustrate integrative clustering with the LUCID early integration model, we use simulated data based on real cases from the HELIX study based on the correlation structure. A subset of the simulated data is incorporated in the **LUCIDus** package to replicate the workflow presented here. The dataset is a list of dataframe containing data from 420 children, with 1 variable measuring maternal exposure to utero mercury (referred to as the exposome), 10 methylomes measured in children (referred to as the methylomics), 10 transcriptomes measured in children (referred to as the transcriptomics), 10 miRNA measured in children (referred to as the miRNA), childhood cytokeratin 18 level (referred to as ck18, a continuous outcome as an indicator of metabolic-dysfunction-associated fatty liver disease (MAFLD), childhood cytokeratin 18 category (a binary outcome referred to as `ck18_cat`) and 2 covariates, including child’s sex and child’s age. To organize the data into separate dataframes for each type of data and to better illustrate functionality within **LUCIDus**, we use the following code:

```
> library(LUCIDus)
> # load data
> data("simulated_HELIX_data")
> simulated_data <- simulated_HELIX_data
> exposome <- as.matrix(simulated_data[["phenotype"]]$hs_hg_m_scaled)
> colnames(exposome) <- "hs_hg_m_scaled"
> methylomics <- simulated_data$methylome
> ck18 <- as.matrix(simulated_data[["phenotype"]]$ck18_scaled)
> colnames(ck18) <- "ck18_scaled"
> ck18_cat <- ifelse(ck18 > mean(ck18), 1, 0)
> covars <- c("hs_child_age_yrs_None", "e3_sex_None")
> covs <- simulated_data[["phenotype"]][covars]
> covs$e3_sex_None <- ifelse(covs$e3_sex_None == "male", 1, 0)
```

Our goal is to conduct an integrated clustering of the exposome and methylome and to relate the estimated latent clusters to the childhood cytokeratin 18 level. The LUCID model is fitted using the function `lucid()`. The input data are specified by arguments  $G$ ,  $Z$ , and  $Y$ , corresponding to the exposome, the methylome, and the childhood cytokeratin 18 level, respectively. In practice, scaling of multi-omics data is highly recommended to obtain more stable estimates; the methylomics data used in this example are already scaled. If multi-omics data are included in the LUCID analysis as the omic intermediate  $Z$ , the user can concatenate them one by one into a single data matrix and then input the concatenated matrix as  $Z$  for early integration. For illustrative purposes, we assume that, in the example data, the optimal number of latent clusters is 2 (determining the optimal number of clusters is

a major question prior to performing clustering analysis; more details are discussed in the later Section [Model selection and variable selection](#)). For illustration, we fit two LUCID models with continuous outcome ck18 and binary outcome ck18\_cat, respectively. The parameter family specifies the outcome type. The default setting is ‘normal’, corresponding to a continuous outcome. For a binary outcome, family should be specified as ‘binary’. `lucid()` returns an object of class ‘lucid\_early’ providing the MLE of the LUCID model as well as IPs.

```
> # Fit a LUCID model with a continuous outcome
> fit1 <- lucid(G = exposome, Z = methylomics, Y = ck18, init_omic.data.model = NULL,
+                  lucid_model = "early", family = "normal", K = 2)
> # MLE of the LUCID model
> # fit1$res_Beta
> # fit1$res_Mu
> # fit1$res_Sigma
> # fit1$res_Gamma
> # IPs of the sample
> # fit1$inclusion.p
> # Fit LUCID model with a binary outcome
> fit2 <- lucid(G = exposome, Z = methylomics, Y = ck18_cat, init_omic.data.model = NULL,
+                  lucid_model = "early", family = "binary", K = 2)
```

### Initializing the parameters of LUCID

$\beta$  is initiated by randomly drawing from a uniform distribution. By default, `init_par` is ‘mclust’ and `init_omic.data.model` is ‘EEV’, so `lucid()` calls `Mclust()` to fit the mixture model of ‘EEV’ on omics data  $Z$  and use the mixture model estimates to initiate  $\mu$  and  $\Sigma$ , and then implement regression to initiate  $\gamma$ . Alternatively, the user can also specify `init_omic.data.model` to choose a certain mixture model. The available mixture models are listed in `mclust::mclustModelNames` ([Scrucca et al., 2016](#)). Note that `init_omic.data.model` can be ‘NULL’ to let `Mclust()` automatically choose the optimal mixture model. If not using `Mclust()` for initiation, the user can also specify `init_par` to be ‘random’ to initiate  $\mu$  and  $\Sigma$  by randomly drawing them from a uniform distribution and then implement regression to initiate  $\gamma$ . We recommend using the default setting to initiate the parameters by calling `Mclust()` for quick convergence and stable performance.

```
> # Fit LUCID model with spherical shape, equal volume
> fit3.1 <- lucid(G = exposome, Z = methylomics, Y = ck18,
+                  lucid_model = "early", family = "normal", K = 2,
+                  init_par = "mclust", init_omic.data.model = "EII")
> # Fit LUCID model with random guess
> fit3.2 <- lucid(G = exposome, Z = methylomics, Y = ck18,
+                  lucid_model = "early", family = "normal", K = 2,
+                  init_par = "random")
> # Fit LUCID model with ellipsoidal shape, varying volume, shape, and orientation
> fit4 <- lucid(G = exposome, Z = methylomics, Y = ck18,
+                  lucid_model = "early", family = "normal", K = 2,
+                  init_omic.data.model = "VVV")
```

### Supervised LUCID versus unsupervised LUCID

In Equation 1, the latent clusters  $X$  are jointly defined by genomic/environmental exposure  $G$ , other multi-omics data  $Z$ , and outcome  $Y$ . Because the outcome is explicitly incorporated in the likelihood function, the estimation process based on Equation 1 is similar to a supervised learning process. **LUCIDus** also allows for an unsupervised version of LUCID. In this unsupervised LUCID, the latent clusters are estimated only by  $G$  and  $Z$ . Specifically, the joint likelihood of unsupervised LUCID is written as

$$\log L(\Theta) = \sum_{i=1}^N \sum_{j=1}^K I(X_i = j) (\log S(X_i = j | G_i; \Theta) + \log \phi(Z_i | X_i = j; \Theta)) \quad (11)$$

and the corresponding responsibility based on Equation 11 is derived as

$$r_{ij}^{(t)} = \frac{S(X_i = j | G_i; \beta^{(t)}) \phi(Z_i | X_i = j; \mu_j^{(t)}, \Sigma_j^{(t)})}{\sum_{j=1}^K S(X_i = j | G_i; \beta^{(t)}) \phi(Z_i | X_i = j; \mu_j^{(t)}, \Sigma_j^{(t)})} \quad (12)$$

Estimations for unsupervised LUCID are also obtained by the EM algorithm discussed previously. The parameter `useY` in `lucid()` is a flag indicating a supervised or unsupervised LUCID model. By default, `useY = TRUE`, and `lucid()` fits supervised LUCID. To fit an unsupervised LUCID model, a user should set `useY = FALSE`.

```
> # Fit an unsupervised LUCID model
> fit5 <- lucid(G = exposome, Z = methylomics, Y = ck18, lucid_model = "early",
+                  family = "normal", K = 2, useY = FALSE)
```

The choice of a supervised LUCID or unsupervised LUCID model depends both on the study design and the research question. A supervised LUCID analysis may be appropriate, for example, if: (1) there is the belief that the cluster structure is defined jointly by  $G$ ,  $Z$ , and  $Y$  (for example, data collected from a cross-sectional design); or (2) the goal is to build a predictive model with data from one cohort to then apply directly to another cohort. If the aim is to obtain an unbiased estimate for the association between the latent cluster  $X$  and the outcome  $Y$ , an unsupervised LUCID model is more appropriate.

### Adjusting for covariates

**LUCIDus** also allows for the adjustment of covariates. According to the DAG in Figure 1 (a), covariates may be included for the association between the exposure and the latent cluster (referred to as  $G$  to  $X$  covariate) or for the association between the latent cluster and the outcome (referred to as  $X$  to  $Y$  covariate). Covariates in the  $G$  to  $X$  relationship act more like predictors of the latent cluster while covariates in the  $X$  to  $Y$  relationship may be interpreted more in the context of an adjustment for a potential confounding effect. A variable can serve as both a  $G$  to  $X$  covariate and a  $X$  to  $Y$  covariate. In the `lucid()` function,  $G$  to  $X$  covariates are specified by parameter `CoG` and  $X$  to  $Y$  covariates are specified by parameter `CoY`.

```
> # Include covariates as G to X covariates
> fit6 <- lucid(G = exposome, Z = methylomics, Y = ck18, lucid_model = "early", K = 2,
+                  family = "normal", CoG = covs)
> # Include covariates as X to Y covariates
> fit7 <- lucid(G = exposome, Z = methylomics, Y = ck18, lucid_model = "early", K = 2,
+                  family = "normal", CoY = covs)
```

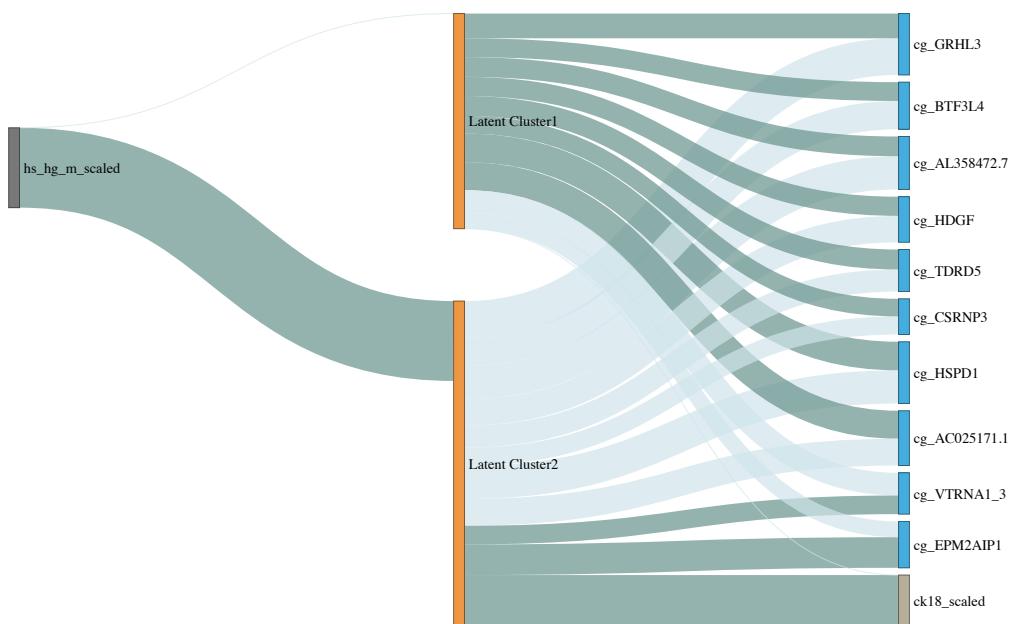
### 3.2 Interpreting LUCID

Since LUCID is an integrated analysis framework consisting of clustering, regression, and multiple data components, we provide two utility functions, `summary()` and `plot()` to summarize the results of LUCID and to facilitate interpretation.

#### Summarizing LUCID in tables by `summary()`

A direct call of `summary()` with the input of a model returned by `lucid()` prints three tables in the console, corresponding to associations among  $G$ ,  $Z$ , and  $Y$ . We take the LUCID model fitted with the continuous BMI as an example.

```
> # summarize a simple lucid model with a continuous outcome
> summary(fit1)
-----Summary of the LUCID Early Integration model-----
K = 2 , log likelihood = -6721.801 , BIC = 14808.7
(1) Y (continuous outcome): effect size of Y for each latent cluster
      Gamma
cluster1 0.0000000
cluster2 0.5040367
(2) Z: mean of omics data for each latent cluster
      mu_cluster1 mu_cluster2
cg_GRHL3    0.2452019 -0.3627827
cg_BTF3L4    0.1863510 -0.2798288
cg_AL358472.7 0.1954543 -0.3290283
cg_HDGF     0.1897587 -0.2614602
cg_TDRD5    0.1977435 -0.2186209
cg_CSRNP3   0.1742812 -0.1794667
cg_HSPD1    0.2829609 -0.3282916
cg_EPM2AIP1 -0.1565522  0.3031461
```



**Figure 3:** An example of using the Sankey diagram to visualize a LUCID model. The dark grey nodes represent the exposure, the light grey node represents the outcome, the blue nodes are omics data, and the orange nodes are latent clusters. The width of links and nodes corresponds to effect size. Light-colored links represent negative associations while dark-colored links indicate positive associations.

```
cg_AC025171.1  0.2769363 -0.2656729
```

```
cg_VTRNA1_3   -0.2250234  0.1847493
```

(3) E: odds ratio of being assigned to each latent cluster for each exposure

beta            OR

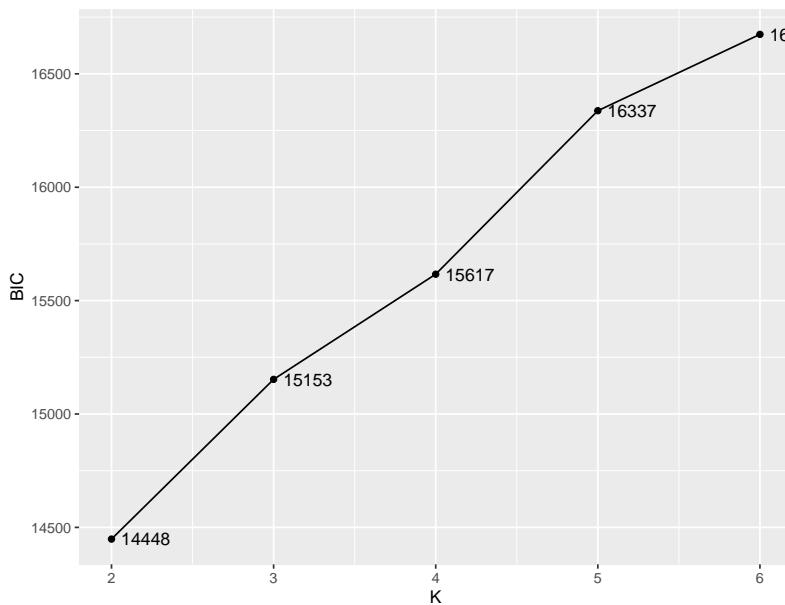
```
hs_hg_m_scaled.cluster2 0.7909168 2.205417
```

The first table summarizes the association between the latent cluster and the outcome (cytokeratin 18). The estimated effect size of cytokeratin 18 for cluster 1 is 0 since it is the reference cluster while that for cluster 2 is 0.504. In the case of a binary outcome, the coefficient in the first table is interpreted as the log odds for the variable for that specific cluster vs. the reference cluster. The second table characterizes each cluster by its omics signature. The omics signature is a cluster-specific mean for each variable in the omics data. For instance, latent cluster 1 is characterized by low levels of cg\_EPM2AIP1 ( $\mu_1(\text{cg\_EPM2AIP1}) = -0.157$ ) and cg\_VTRNA1\_3 ( $\mu_1(\text{cg\_VTRNA1\_3}) = -0.225$ ), and high levels of all other omic features; for example, cg\_GRHL3 ( $\mu_1(\text{cg\_GRHL3}) = 0.2452019$ ). On the contrary, latent cluster 2 features high levels of cg\_EPM2AIP1 ( $\mu_2(\text{cg\_EPM2AIP1}) = 0.303$ ) and cg\_VTRNA1\_3 ( $\mu_2(\text{cg\_VTRNA1\_3}) = 0.185$ ), and low levels of all other features including cg\_GRHL3 ( $\mu_2(\text{cg\_GRHL3}) = -0.363$ ). The third table relates the exposures to the latent clusters. For instance, 'hs\_hg\_m\_scaled' represents the scaled level of maternal mercury exposure. The coefficient OR for 'hs\_hg\_m\_scaled' is 0.791, meaning that for each doubling of the scaled level of maternal mercury exposure, the odds ratio of being assigned to latent cluster 2 is 2.205. Since latent cluster 2 is associated with a higher child level of cytokeratin 18, these results indicate that exposure to mercury is associated with a higher child level of cytokeratin 18 and thus higher risks of MAFLD.

### Visualizing LUCID by plot()

Visualization is another imperative way to interpret statistical models. Here we use a Sankey diagram (Schmidt, 2008) for the visualization of the relationships among different components in LUCID. `plot()` takes a fitted LUCID model as input and creates a Sankey diagram in `html` format. It also accepts a user-defined color palette.

```
> # Visualize LUCID model via a Sankey diagram
> plot(fit1)
> # Change the node color
> plot(fit1, G_color = "yellow")
```



**Figure 4:** Choosing an optimal number of latent clusters K based on BIC.

```
> # Change the link color
> plot(fit1, pos_link_color = "red", neg_link_color = "green")
```

Figure 3 shows an example of a Sankey diagram for the example data. Each node represents a component in the LUCID model. Different components are labeled by user-defined colors. Each link represents the statistical association between two nodes. The width of the links reflects the effect size of the association, and the color of the links indicates the direction of the association. By default, dark-colored links represent positive associations while light-colored links indicate negative associations. The output of the Sankey diagram is an interactive html file in which the user can drag and rearrange the layout of the elements to avoid overlapping. In practice, we recommend limiting the number of variables in the Sankey diagram to approximately 10 to facilitate interpretation.

### 3.3 Model selection and variable selection

To determine the number of latent clusters, K, LUCID implements model selection by fitting a series of models, each with a different value for K, and then the Bayesian Information Criteria (BIC) is used to choose the optimal model (Fan and Tang, 2013). The optimal model is the one with the lowest BIC. The BIC for LUCID is defined as

$$\text{BIC} = -2 \log L(\hat{\Theta}) + D \log N \quad (13)$$

where  $\hat{\Theta}$  is the maximum likelihood estimation and D is the number of parameters in LUCID. Specifically,  $D = P(K - 1) + KM + KM(M + 1)/2 + n_Y$ , where  $n_Y$  is the number of parameters dependent upon the type of outcome Y. If Y is a continuous outcome, then  $n_Y = 2K$ . If Y is a binary outcome, then  $n_Y = K$ . If a vector of K is used as input to the function `lucid()`, model selection is automatically performed and returns the model with optimal K.

```
> fit8 <- lucid(G = exposome, Z = methylomics, Y = ck18, lucid_model = "early",
+                   family = "normal", K = 2:6)
> # Check the optimal K
> fit8$K
[1] 2
```

Separately, users can use the auxiliary function `tune_lucid()` to explore model selection in more detail. This function returns the optimal model as well as a table recording the tuning process. Below is an example of tuning K and visualizing the tuning process by using `tune_lucid()`.

```
> # Look into the tuning process in more detail
> tune_K <- tune_lucid(G = exposome, Z = methylomics, Y = ck18, lucid_model = "early",
+                       family = "normal", K = 2:6)
> fit9 <- tune_K$best_model
```

```
> ggplot(data = tune_K$tune_list, aes(x = K, y = BIC, label = round(BIC, 0))) +
+   geom_point() +
+   geom_line() +
+   geom_text(hjust = -0.2)
```

Figure 4 shows the tuning process with  $K = 2$  resulting in the lowest BIC. Including additional or redundant variables in clustering models may increase model complexity, impair prediction accuracy, and boost computational time (Fop and Murphy, 2018). **LUCIDus** performs variable selection for both the genetic/environmental exposures  $G$  and for the other multi-omics data  $Z$  via  $L_1$ -norm penalty terms. For variable selection for  $G$ , we apply the LASSO regression to obtain sparse solutions for Equation 5, which is

$$\boldsymbol{\beta}_{\text{LASSO}}^{(t+1)} = \arg \max_{\boldsymbol{\beta}} \left\{ \sum_{i=1}^N \sum_{j=1}^K r_{ij}^{(t)} \log S(X_i = j | G_i; \boldsymbol{\beta}_j) - \lambda_{\beta} \sum_{j=1}^K \sum_{l=1}^P |\beta_{jl}| \right\} \quad (14)$$

To obtain sparse estimation for  $\boldsymbol{\mu}$  and  $\boldsymbol{\Sigma}$ , we implement the penalized model-based method (Zhou et al., 2009).  $\boldsymbol{\mu}$  is first updated by maximizing the following equation,

$$\boldsymbol{\mu}_j^{(t+1)} = \arg \max_{\boldsymbol{\mu}} \left\{ \sum_{i=1}^N \sum_{j=1}^K r_{ij}^{(t)} \log \phi(Z_i | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j) - \lambda_{\mu} \sum_{j=1}^K \sum_{l=1}^M |\mu_{jl}| \right\} \quad (15)$$

Denote  $\mathbf{W} = \boldsymbol{\Sigma}^{-1}$ . Then  $\boldsymbol{\Sigma}$  is updated by maximizing its inverse penalized by  $L_1$ -norm

$$\mathbf{W}_j^{(t+1)} = \arg \max_{\mathbf{W}} \left\{ \sum_{i=1}^N \sum_{j=1}^K r_{ij}^{(t)} (\det \mathbf{W}_j - \text{trace}(\mathbf{S}_j^{(t)} \mathbf{W}_j)) - \lambda_W \sum_{ls} |w_{jls}| \right\} \quad (16)$$

where  $\mathbf{S}_j$  is the empirical covariance matrix at each iteration, defined as

$$\mathbf{S}_j^{(t)} = \frac{\sum_{i=1}^N r_{ij} \left( Z_i - \boldsymbol{\mu}_j^{(t)} \right) \left( Z_i - \boldsymbol{\mu}_j^{(t)} \right)^T}{\sum_{i=1}^N r_{ij}^{(t)}} \quad (17)$$

To choose the optimal combination of the three  $L_1$ -norm penalties, we implement a grid search by adopting a modified BIC (Pan and Shen, 2007), defined as

$$\text{BIC}_p = -2 \log L(\hat{\boldsymbol{\Theta}}) + (D - D_G - D_Z) \log N \quad (18)$$

Where  $D_G$  is the number of exposure variables whose effect estimates are 0 across all latent clusters and  $D_Z$  is the number of omic variables whose effect estimates are 0 across all latent clusters. We can tune  $\lambda_{\beta}$ ,  $\lambda_{\mu}$ , and  $\lambda_W$  either separately or jointly. For instance, if only exposome data  $G$  is high-dimensional, we only need to tune  $\lambda_{\beta}$ . If variable selection is desired for both  $G$  and  $Z$ , then the three  $L_1$ -norm penalties should be tuned simultaneously.  $\lambda_{\beta}$ ,  $\lambda_{\mu}$ , and  $\lambda_W$  match the parameters `Rho_G`, `Rho_Z_Mu`, and `Rho_Z_Cov` in the `lucid()` function. Each parameter accepts a numeric vector or a scalar as input. For guidance, empirical experiments utilizing normalized multi-omics data suggest integer values in the range of 0 – 100 for `Rho_Z_Mu` and values in the range of 0 – 1 for `Rho_G` and `Rho_Z_Cov`. The higher the values of `Rho_Z_Mu`, `Rho_G`, and `Rho_Z_Cov`, the fewer omic features that will be selected. Below are examples for conducting variable selection for  $G$  and  $Z$ , separately and jointly.

```
> # Variable selection for G
> # Add 10 more noise variables to the exposome
> noise <- matrix(rnorm(420 * 10), nrow = 420)
> exposome_noise <- cbind(exposome, noise)
> fit10 <- lucid(G = exposome_noise, Z = methylomics, Y = ck18,
+                   lucid_model = "early", family = "normal", K = 2,
+                   Rho_G = seq(0, 0.4, by = 0.01), seed = 1008)
1/11 exposures are selected
> # Summary of optimal lucid model
> # summary(fit10)
> # Variable selection for Z
> # add 10 more noise variables to the methylomics
> methylomics_noise <- cbind(methylomics, noise)
> fit11 <- lucid(G = exposome, Z = methylomics_noise, Y = ck18,
+                   lucid_model = "early", family = "normal", K = 2,
+                   Rho_Z_Mu = 5:10, Rho_Z_Cov = seq(0.1, 0.4, by = 0.1), seed = 1008)
```

```

10/20 omics variables are selected
> # Summary of optimal lucid model
> # summary(fit11)
> # Variable selection for G and Z jointly
> fit12 <- lucid(G = exposome_noise, Z = methylomics_noise, Y = ck18,
+                   lucid_model = "early", family = "normal", K = 2,
+                   Rho_G = 0.01, Rho_Z_Mu = 10, Rho_Z_Cov = 0.5, seed = 123)
1/11 exposures are selected
11/20 omics variables are selected

```

Oftentimes, the number of features  $M$  for omics data  $Z$  is in the hundreds or even thousands. With a high dimensionality of  $Z$ , LUCID can still have stable performance and select the features with effect if the number of observations  $N$  is higher than  $M$ . Below is an example of conducting variable selection in  $Z$  with  $M = 210$  features (including 200 noise features) while  $N = 420$ . We can see that with higher  $Rho\_Z\_Mu$  and  $Rho\_Z\_Cov$ , LUCID successfully selects the features with effect.

```

> # Variable selection for Z (high number of features M)
> set.seed(1008)
> # Add 200 more noise features to the methylomics
> noise <- matrix(rnorm(420 * 200), nrow = 420)
> methylomics_noise <- cbind(methylomics, noise)
> fit13 <- lucid(G = exposome, Z = methylomics_noise, Y = ck18,
+                   family = "normal", lucid_model = "early", K = 2,
+                   Rho_Z_Mu = 20, Rho_Z_Cov = 0.6)
17/210 omics variables are selected
> # Summary of optimal lucid model
> # summary(fit13)

```

Via simulations, the run time for LUCIDus increases linearly as the number of features increases. While the integration of regularization allows for a high number of features to be analyzed, when  $M$  is higher than  $N$ , LUCID might result in unstable estimation; therefore, pre-screening approaches for the omic features such as the "meet-in-the-middle" become imperative to reduce dimensionality before fitting the LUCID model (Cadiou et al., 2021).

### 3.4 Deriving confidence intervals

We use the nonparametric bootstrap approach to derive confidence intervals (CIs) for the MLE estimates from a LUCID model (Davison and Hinkley, 1997). Nonparametric bootstrap draws observations from a sample with replacement. We index the bootstrap samples by  $b = 1, 2, \dots, B$ . For each bootstrap sample  $D^b$ , we fit the LUCID model and calculate MLE  $\Theta^b$ . Two types of bootstrap CIs are constructed based on the distribution of  $\Theta^b$ : (1) normal CIs and (2) percentile CIs. Given the confidence level  $1 - \alpha$ , the normal CIs are constructed by

$$(1 - \alpha)\%CI_{\text{normal}} = \Theta - \text{bias} \pm Z_{1-\alpha/2}\sigma(\Theta) \quad (19)$$

where  $\text{bias} = \bar{\Theta} - \Theta$ ,  $\bar{\Theta}$  is the mean of the bootstrap estimators and  $\sigma(\Theta)$  is the bootstrap standard error, which is

$$\sigma(\Theta) = \sqrt{\frac{1}{B-1} \sum_{b=1}^B (\Theta^b - \bar{\Theta})^2} \quad (20)$$

All calculations in Equations 19 and 20 are element-wise. In addition, the  $1 - \alpha$  percentile CIs are calculated by ordering the bootstrap estimators from smallest to largest and selecting the estimates at  $\alpha/2$  percentile and  $(1 - \alpha/2)$  percentile as the CIs. `boot_lucid()` constructs the two bootstrap CIs described above for inference. Users can specify the confidence level (`conf`) and the number of bootstrap replicates (`R`). While `boot_lucid()` is running, a progress bar is displayed in the R console. We recommend  $R \geq 200$  to estimate the bootstrap standard error for normal CIs and  $R \geq 800$  to estimate quantiles to construct percentile CIs. Below are examples of deriving bootstrap CIs for LUCID. By default, LUCID calculates 95% CIs.

```

> # Bootstrap to obtain 95% CI (by default) for LUCID
> set.seed(123)
> boot1 <- boot_lucid(G = exposome, Z = methylomics, Y = ck18, lucid_model = "early",
+                      model = fit1, R = 200)
# 90% CIs
> boot2 <- boot_lucid(G = exposome, Z = methylomics, Y = ck18, lucid_model = "early",
+                      model = fit1, R = 200, conf = 0.9)

```

We can obtain a more comprehensive summary table with bootstrap CIs by integrating `summary()` with output from `boot_lucid()`. The resulting summary table has 5 columns: ‘ $t_0$ ’ corresponds to parameters estimated from the observed data; ‘`norm_lower`’ and ‘`norm_upper`’ represent the lower and upper limit of normal CIs, respectively; ‘`perc_lower`’ and ‘`perc_upper`’ represent the percentile CIs.

```
> # Summary table with 95% bootstrap CIs
> summary(fit1, boot.se = boot1)
```

`boot_lucid()` is built upon the R library `boot` (Canty and Ripley, 2021). The original output from `boot` is also returned by `boot_lucid()` and can be used to evaluate the normality of the distribution of the bootstrap estimations.

### 3.5 Prediction

The prediction of LUCID includes two parts: prediction of the latent clusters  $X$  and prediction of the outcome  $Y$ . The `predict_lucid()` can perform both tasks. Prediction for the latent cluster is determined by IP using the optimal rule, which is

$$\hat{X}_i = \arg \max_j P(X_i = j | D, \Theta) \quad (21)$$

Prediction for the outcome follows a conventional regression framework. When making predictions, `predict_lucid()` for LUCID requires input of new  $G$  and  $Z$ , and optional input of  $Y$ . If  $Y$  is provided, we use Equation 3 to calculate IPs; otherwise, Equation 12 is applied. As a result, we can either fit a supervised LUCID model to make predictions in an unsupervised fashion or vice versa. This design allows for flexibility. For instance, we can build a LUCID model in one cohort in a supervised fashion and then make predictions on another new independent cohort, regardless of whether the outcome is measured or not. `predict_lucid()` for LUCID returns a list containing IPs for each observation, predicted cluster assignment, and predicted outcome. Below is an example code for prediction based on a supervised LUCID model, `fit1`.

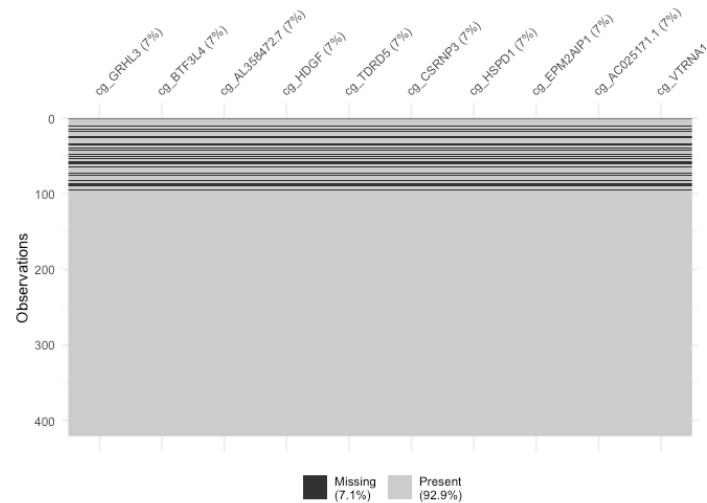
```
> # Predict cluster with information of Y
> pred1 <- predict_lucid(model = fit1, G = exposome, Z = methylomics, Y = ck18)
> # Predict cluster without information of Y
> pred2 <- predict_lucid(model = fit1, G = exposome, Z = methylomics)
> # Predicted cluster label
> table(pred1$pred.x)
  1   2
223 197
> # Predicted outcome
> pred1$pred.y[1:5]
[1] 0.2580164 -0.3806054  0.2504817  0.4123169 -0.1898255
```

### 3.6 Incorporating missingness in multi-omics data

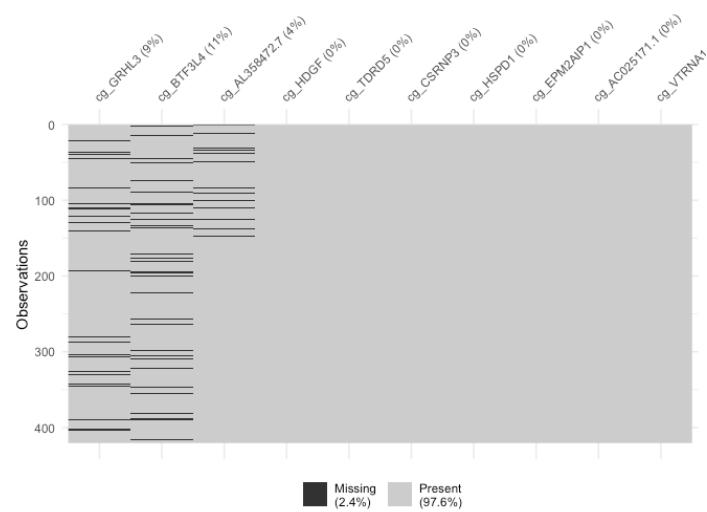
A major update in **LUCIDus** version 3 is the incorporation of missing data. Missingness is a major challenge in the integrative analysis of omics. In large cohort studies, it is common that some omics data are not available for all participants for various reasons such as budgetary constraints, low sample availability, or lack of consent for future use of biospecimens (Voillet et al., 2016). We refer to this type of missingness as a list-wise missingness pattern. In addition, even when omics data are measured for a given participant, it is common that some omic features are randomly missing due to the measurement process. We refer to this pattern as sporadic missingness. **LUCIDus** can deal with the two missing patterns mentioned above or the combination of the two. We assume the multi-omics data are missing completely at random (MCAR) for both missing patterns (i.e., for the list-wise missing pattern, the probability of multi-omics data being not available is the same for all subjects; for sporadic missing patterns, the probability of a certain omic feature being missing is the same for all omic features). For list-wise missingness, we use a modified LUCID model based on a likelihood partition. We divide the LUCID likelihood into two parts: observations with complete multi-omics data and observations without any measured multi-omics data. The likelihood of the former observations remains the same as Equation 2, denoted by  $l_o(\Theta | D)$ . The joint likelihood of the latter observations becomes

$$l_m(\Theta | D) = \sum_{l_o=1}^{N_o} \sum_{j=1}^K I(X_{i_o} = j) \left( \log S(X_{i_o} = j | G_{i_o}; \beta_j) + \log \phi(Y_{i_o} | \gamma_j, \sigma_j^2) \right) \quad (22)$$

## List-wise missing pattern



## Sporadic missing pattern



**Figure 5:** Two simulated missing patterns in methylomics data

with corresponding responsibility defined as

$$r_{ij}^{(t)} = \frac{S(X_i = j | G_i; \beta^{(t)}) \phi(Y_i | X_i = j; \gamma_j^{(t)}, \sigma_j^{2(t)})}{\sum_{j=1}^K S(X_i = j | G_i; \beta^{(t)}) \phi(Y_i | X_i = j; \gamma_j^{(t)}, \sigma_j^{2(t)})} \quad (23)$$

The joint likelihood of LUCID for all data becomes  $l(\Theta) = l_o(\Theta | D) + l_m(\Theta | D)$ . We then use the same EM algorithm described in Section [Overview of the original LUCID model \(early integration\)](#) to calculate the MLE of LUCID with list-wise missingness. For sporadic missingness, we use an integrated imputation method ([Zhang et al., 2021](#)). Each missing value in the omics data  $Z$  is treated as an unknown “parameter” to be optimized. Specifically in the M-step, after maximizing the parameters of the LUCID model with fixed  $Z$ , we implement an additional *imputation step* which maximizes  $Z$  with fixed parameters within LUCID. Details of the statistical derivation can be found in [Zhang et al. \(2021\)](#). To illustrate this new feature, we first simulate missing values in the HELIX methylomics data by randomly setting values to ‘NA’. Both list-wise and sporadic missing patterns are considered, as shown in Figure 5. We use the `vis_miss()` function from the R package `visdat` to visualize both missing patterns ([Tierney, 2023](#)).

```
> library(visdat)
> set.seed(1)
> methylomics_miss_sporadic <- methylomics_miss_listwise <- as.matrix(methylomics)
> index <- arrayInd(sample(1000, 0.1 * 1000), dim(methylomics))
> methylomics_miss_sporadic[index] <- NA # sporadic missing pattern
> methylomics_miss_listwise[sample(1:100, 30), ] <- NA # listwise missing pattern
> vis_miss(as.data.frame(methylomics_miss_sporadic))
> vis_miss(as.data.frame(methylomics_miss_listwise))
```

Below are examples of fitting the LUCID model with missingness in multi-omics data. `lucid()` will automatically examine missing patterns present in the multi-omics data and select a corresponding method to account for the missing data. For sporadic missing patterns, we use `mclust` to initialize missing values in the multi-omics data.

```
> fit14 <- lucid(G = exposome, Z = methylomics_miss_listwise, Y = ck18,
+                   lucid_model = "early", K = 2, family = "normal")
> fit15 <- lucid(G = exposome, Z = methylomics_miss_sporadic, Y = ck18,
+                   lucid_model = "early", K = 2, family = "normal")
> # summary(fit15)
```

### 3.7 LUCID in Parallel

In the previous sections, we focused on using the early integration strategy by concatenating each omic layer one by one and inputting a single data matrix into the LUCID model for estimation. However, researchers may be interested in modeling the correlation structure of each omic layer independently to investigate how multi-omics data may act in parallel with an outcome. In this section, we extend the LUCID model to incorporate multi-omics data by introducing multiple latent variables into the framework. Suppose we have a sample of size  $n$ , indexed by  $i = 1, 2, \dots, n$ . It has a collection of  $m$  omic layers, denoted by  $Z_1, \dots, Z_a, \dots, Z_m$  with corresponding dimensions  $p_1, \dots, p_a, \dots, p_m$ . Each omic layer  $Z_a$  is summarized by a latent categorical variable  $X_a$ , which contains  $k_a$  categories. Each category is interpreted as a latent cluster (or subgroup) for that particular omic layer. All latent variables are linked to the same exposure matrix  $G$  and the outcome  $Y$  for LUCID in parallel, as shown in Figure 1 (b). Suppose the latent variable  $X_{ai}$  is observed for  $a = 1, \dots, m$ . According to the DAG in Figure 1 (b), the distributions of  $f(Z_{ai} | G)$  are conditionally independent of each other for  $a = 1, \dots, m$ , and the distributions of the latent variable  $f(X_{ai} | G_i)$  are also conditionally independent of each other. Since  $X_{ai}$  is a discrete variable with  $k_i$  categories, we assume  $X_{ai}$  follows a multinomial distribution conditioning on  $G$ , denoted by the softmax function  $S(\cdot)$ . We assume multi-omics data  $Z_{ai}$  follows a multivariate Gaussian distribution conditioning on  $X_{ai}$ , denoted by  $\phi(Z_{ai} | X_{ai} = j_a; \mu_{ja}, \Sigma_{ja})$ , where  $\mu_{ja}$  and  $\Sigma_{ja}$  are cluster-specific means and variance-covariance matrices for omic layer  $a$ . The outcome  $Y$  can be either a continuous outcome or a binary outcome. For illustration, here, we discuss the situation that  $Y$  is a continuous variable that follows a Gaussian distribution. The relationship between latent variables  $X_a$  and outcome  $Y$  is formulated as

$$EY_i = \gamma_0 + \sum_{j_1=2}^{k_1} \gamma_{j_1} I(X_{1i} = j_1) + \cdots + \sum_{j_a=2}^{k_a} \gamma_{j_a} I(X_{ai} = j_a) + \cdots + \sum_{j_m=2}^{k_m} \gamma_{j_m} I(X_{mi} = j_m) \quad (24)$$

where the intercept  $\gamma_0$  is the expected value of  $Y_i$  given  $X_{ai} = 1$  for  $a = 1, 2, \dots, m$  (the reference cluster for all combinations of latent cluster assignments);  $\gamma_{j_a}$  is the change of  $Y$  if the latent variable  $X_a$  becomes  $j_a$  instead of 1. Let  $D$  be the generic notation for all observed data. The log-likelihood of LUCID in Parallel is constructed below:

$$\begin{aligned}
\log L(\Theta | D) &= \sum_{i=1}^n \log f(Z_{1i}, \dots, Z_{mi}, Y_i | G_i; \Theta) \\
&= \sum_{i=1}^n \log \left[ \prod_{j_1=1}^{k_1} \dots \prod_{j_m=1}^{k_m} f(Z_{1i}, \dots, Z_{mi}, X_{1i}, \dots, X_{mi}, Y_i | G_i; \Theta)^{I(X_{1i}=j_1, \dots, X_{mi}=j_m)} \right] \\
&= \sum_{i=1}^n \sum_{j_1=1}^{k_1} \dots \sum_{j_m=1}^{k_m} I(X_{1i}=j_1, \dots, X_{mi}=j_m) \log f(Z_{1i}, \dots, Z_{mi}, X_{1i}, \dots, X_{mi}, Y_i | G_i; \Theta) \\
&= \sum_{i=1}^n \sum_{j_1=1}^{k_1} \dots \sum_{j_m=1}^{k_m} I(X_{1i}=j_1, \dots, X_{mi}=j_m) \log \phi(Y_i | X_{1i}, \dots, X_{mi}, \gamma, \sigma^2) \\
&\quad + \sum_{i=1}^n \sum_{a=1}^m \sum_{j_1=1}^{k_1} \dots \sum_{j_m=1}^{k_m} I(X_{1i}=j_1, \dots, X_{mi}=j_m) \log \phi(Z_{ai} | X_{ai}=j_a, \mu_{a,j_a}, \Sigma_{a,j_a}) \\
&\quad + \sum_{i=1}^n \sum_{a=1}^m \sum_{j_1=1}^{k_1} \dots \sum_{j_m=1}^{k_m} I(X_{1i}=j_1, \dots, X_{mi}=j_m) \log S(X_{ai}=j_a | G_i, \beta_a)
\end{aligned} \tag{25}$$

The log-likelihood of LUCID in parallel is similar to that with LUCID early integration. It is natural to follow the same principles of the EM algorithm for LUCID early integration and estimate the parameters of LUCID in parallel with targeted modifications. For illustration, we continue with the previous data example using maternal exposure to mercury as the exposome  $G$  and childhood cytokeratin 18 level as the outcome  $Y$ . For a LUCID in parallel model, multi-omics data  $Z$  is a list of three  $420 \times p_a$  matrices of omic layers where  $a = 1, 2, 3$  (layer 1: methylome; layer 2: transcriptome; layer 3: miRNA). We fit LUCID in parallel models with continuous ck18 using `lucid()`. We specify the argument `lucid_model` to be 'parallel', and the argument `K` should be specified as a list of three integers (or sequences of integers if tuning) indicating the number of latent clusters for each omic layer. The variable selection feature is not yet available for LUCID in parallel, and we can only tune for `K`. We can use argument `CoG` and `CoY` to adjust for covariates for the  $G$  to  $X$  association and the  $X$  to  $Y$  association, respectively. Similar to LUCID early integration, the parameter `family` specifies the outcome type, which is 'normal' here, corresponding to a continuous outcome. For a binary outcome, `family` should be specified as 'binary'. `lucid()` returns an object of class 'lucid\_parallel' providing the MLE of the LUCID in parallel model as well as IPs. We can call `summary()` with the input of a 'lucid\_parallel' object to print three tables of the model results. Similar to the summary table of a LUCID early integration model, the first table summarizes the association between the latent clusters for each omic layer and the outcome (cytokeratin 18). The second table characterizes each cluster by its omics signature within each layer. The third table relates the exposures to the latent clusters for each layer. The nuance is that the LUCID in parallel model has multiple independent omic layers, and each layer corresponds to a separate latent cluster variable, thus all the related results are presented in the summary tables. See the two examples of fitting a LUCID in parallel model below.

```

> # Create a list of multi-omics data
> omics_lst <- simulated_data[-which(names(simulated_data) == "phenotype")]
> Z = omics_lst[c(1:3)]
> # LUCID in parallel, adjusting for the covariates for the E-X and X-Y associations
> fit16 <- lucid(G = exposome, Z = Z, Y = ck18, K = list(2, 2, 2), family = "normal",
+                   CoY = covs, CoG = covs,
+                   lucid_model = "parallel", useY = TRUE)
> # print the summary of the LUCID in parallel model
> summary(fit16)
-----Summary of the LUCID in Parallel model-----
K = 2 2 2 , log likelihood = -16413.64 , BIC = 36892.36
(1) Y (continuous outcome): effects of each non-reference latent cluster
    for each layer of Y
    (and effect of covariates if included)
        Gamma
cluster2Layer1      2.139337462
cluster2Layer2      -0.976575066
cluster2Layer3      1.617148373

```

```

e3_sex_None          0.023397524
hs_child_age_yrs_None -0.002297037
(2) Z: mean of omics data for each latent cluster of each layer
Layer 1
          mu_cluster1 mu_cluster2
cg_GRHL3      0.09782065 -0.2827826
cg_BTF3L4      0.12271079 -0.3059186
cg_AL358472.7  0.10019649 -0.3164803
cg_HDGF        0.15388979 -0.3322801
cg_TDRD5       0.10778098 -0.1832600
cg_CSRNP3       0.08692918 -0.1300857
cg_HSPD1        0.21255103 -0.3855858
cg_EPM2AIP1     -0.11651397  0.3691065
cg_AC025171.1   0.24747122 -0.3750450
cg_VTRNA1_3     -0.13739036  0.1515515
Layer 2
          mu_cluster1 mu_cluster2
tc_TC01006069_nc 0.098203482 -0.169267971
tc_SLC9A4        -0.025858454  0.046609093
tc_RAB6C_AS1      0.009567584 -0.087150028
tc_LOC100129029  0.032557846  0.078278431
tc_BRE           -0.039411493  0.073329913
tc_TC03001220_nc -0.068457436  0.013811537
tc_TC04002114_nc -0.030761261  0.033415671
tc_TC04002369_nc  0.297076410 -0.216123334
tc_BEND4         -0.128343956  0.074585000
tc_SLC9A3        0.139869536  0.002458175
Layer 3
          mu_cluster1 mu_cluster2
miR.101.3p      0.10748628  0.023648293
miR.125a.5p     -0.20581745  0.135670802
miR.125b.1.3p    0.07644648  0.019200103
miR.127.3p      -0.05542160  0.184175902
miR.140.5p       0.14979024  0.050720602
miR.142.3p       0.11157913 -0.017566242
miR.144.5p       0.07526727 -0.016027655
miR.19a.3p       0.09193308 -0.002059547
miR.19b.3p       0.08515311  0.032351416
miR.21.5p        0.05800681  0.029022417
(3) E: odds ratio of being assigned to each latent cluster for each exposure
      for each layer
Layer 1
          beta      OR
hs_hg_m_scaled.cluster2  0.7352341 2.0859703
e3_sex_None.cluster2     -0.4107395 0.6631596
hs_child_age_yrs_None.cluster2 -0.2657909 0.7665994
Layer 2
          beta      OR
hs_hg_m_scaled.cluster2  0.1313174 1.1403296
e3_sex_None.cluster2     0.1928708 1.2127261
hs_child_age_yrs_None.cluster2 -0.1066509 0.8988394
Layer 3
          beta      OR
hs_hg_m_scaled.cluster2  -0.4896515 0.6128399
e3_sex_None.cluster2     0.4855559 1.6250782
hs_child_age_yrs_None.cluster2  0.1042404 1.1098672
> # LUCID in parallel, tune for the number of clusters for methylomics
> fit17 <- lucid(G = exposome, Z = Z, Y = ck18, K = list(2, 2:3, 2), family = "normal",
+                  lucid_model = "parallel", useY = TRUE)
> # summary(fit17)

```

### 3.8 LUCID in Serial

In a more late integration framework, LUCID can also be extended to incorporate multiple latent variables in a serial fashion if researchers believe that given an exposure, multi-omics data act serially through a multistep process towards the outcome, as illustrated in Figure 1 (c). This framework can accommodate the following situations: (1) longitudinal measurements on the same multi-omics data type and (2) biological relationships of multi-omics data. The estimation of latent clusters for each omic layer can be formulated as an unsupervised LUCID early integration or LUCID in parallel sub-model. For a LUCID in serial model,  $Z$  includes  $m$  ordered omic layers in a sequential fashion. For illustration, assuming all the sub-models are LUCID early integration models. For a certain omic component  $a$ , the cluster variable  $X_{a-1}$  serves as the “ $G$ ” in the LUCID early integration model framework, and  $X_a$  is jointly estimated by  $X_{a-1}$  and  $Z_a$ . However, the special cases are for the first and the last omic layer, where the corresponding  $X_1$  is estimated jointly by  $G$  and  $Z_1$ , and the corresponding  $X_m$  is estimated jointly by  $X_{m-1}$ ,  $Z_m$ , and  $Y$  if supervised, respectively. Alternatively, for a certain omic component  $a$ , we can add another transition probability matrix component  $p$  in the joint likelihood relating  $X_{a-1}$  to  $X_a$  to describe the state change with  $P(X_a|X_{a-1}, C) = p_a$ , where  $p_a$  is a general notation for the inclusion probability of all the latent clusters for omic component  $a$  and  $C$  dynamically represents other variables that jointly estimate  $X_a$ . Note that a list of  $m$  ordered omic components for  $Z$  indicates  $m$  sub-models in a LUCID in serial model, and each sub-model can be either a LUCID Early Integration model or a LUCID in Parallel model. Note that the DAG in Figure 1 (c) only shows the scenario where all the sub-models are LUCID early integration models. Bold  $X_a$  and omic component  $Z_a$  indicate that they may include more than 1 latent variable or omic layer for each omic component if the corresponding sub-model is LUCID in parallel. Let  $m$  be the number of sub-models in a LUCID in serial model. For each sub-model  $a$ , given  $\Theta_a$  and  $D_a$ , its own  $\log L(\Theta_a | D_a)$  follows the log-likelihood of a LUCID Early Integration or a LUCID in Parallel model defined in the previous sections. Except for the sub-model  $m$ , all sub-models are unsupervised. For sub-models 2 to  $m$ , the  $\beta$  parameter defined previously becomes  $\delta$ . The log-likelihood of LUCID in Serial is constructed below,

$$\begin{aligned}
 \log L(\Theta | D) &= \sum_{a=1}^m \log L(\Theta_a | D_a) \\
 &= \sum_{i=1}^n \log f_1(Z_{1i}, G_i | \Theta_1) \\
 &\quad + \sum_{a=2}^{m-1} \sum_{i=1}^n \log f_a(Z_{ai}, p_{ai} | \Theta_a) + \sum_{i=1}^n \log f_m(Z_{mi}, p_{ai}, Y_i | \Theta_m) \\
 &= \sum_{i=1}^n \log f_1(Z_{1i}, G_i | \beta, \mu_1, \Sigma_1) \\
 &\quad + \sum_{a=2}^{m-1} \sum_{i=1}^n \log f_a(Z_{ai}, p_{ai} | \delta_a, \mu_a, \Sigma_a) \\
 &\quad + \sum_{i=1}^n \log f_m(Z_{mi}, p_{ai}, Y_i | \delta_m, \mu_m, \Sigma_m, \gamma)
 \end{aligned} \tag{26}$$

Here, we assume sub-models are independent of each other. The EM algorithm for estimating the parameters of each sub-model is exactly the same as estimating a single LUCID early integration model or a LUCID in parallel model. Using the previous data example for illustration, for a LUCID in serial model, multi-omics data  $Z$  can be an ordered list of three  $420 \times p_a$  matrices of omic layers where  $a = 1, 2, 3$  (layer 1: methylome; layer 2: transcriptome; layer 3: miRNA), representing that all the sub-models are LUCID early integration models and  $K$  is a list of three integers (or sequences of integers if tuning) indicating the number of latent clusters for each successively linked omic layer. We specify `lucid_model` to be ‘serial’. The use of arguments `CoG`, `CoY`, `useY`, and `family` is exactly the same as for a LUCID in parallel model. `lucid()` constructs the model and returns an object of class ‘`lucid_serial`’, and calling `summary()` with the input of a ‘`lucid_serial`’ object prints the summarizing tables of the model estimates. We have introduced summarizing tables for the LUCID early integration model and the LUCID in parallel model. Since the LUCID in serial model is a combination of LUCID early integration and LUCID in parallel sub-models, the summarizing tables for each sub-model follow a similar structure to what we have discussed for the two types of sub-models, but with targeted modifications to conform to the setup of the LUCID in serial model. See the summarizing tables for the first LUCID in serial model example below. Alternatively, besides an integer, the element of the ordered list of  $Z$  can also be a list, which indicates that the corresponding sub-model is a LUCID in parallel model. Here,  $K$  should contain a list of lists of integers. See the second LUCID in serial model example below.

```

> # LUCID in serial, adjusting for the covariates of the E-X and X-Y associations
> # Tune the number of clusters for methylome and transcriptome
> # All 3 sub-models are LUCID early integration models and each with two latent clusters for each layer
> fit18 <- lucid(G = exposome, Z = Z, Y = ck18,
+                   lucid_model = "serial", CoY = covs, CoG = covs,
+                   K = list(2:3, 2:3, 2), useY = TRUE, family = "normal")
> # Print the summary of the LUCID in serial model
> summary(fit18)
-----Summary of the First Component of the LUCID in Serial Model-----
-----Summary of the LUCID Early Integration Sub Model-----
K = 2 , log likelihood = -5977.843 , BIC = 13320.78
(1) Z: mean of omics data for each latent cluster
      mu_cluster1  mu_cluster2
cg_GRHL3     0.295741860 -0.410975212
cg_BTF3L4     0.021123889 -0.090792176
cg_AL358472.7 0.301523381 -0.439813049
cg_HDGF      0.170527566 -0.234364618
cg_TDRD5      0.192293276 -0.207240424
cg_CSRNP3      0.235578302 -0.242839045
cg_HSPD1      -0.002680386 -0.003986538
cg_EPM2AIP1      0.033545336  0.086638565
cg_AC025171.1   0.303125436 -0.287725955
cg_VTRNA1_3     -0.090000154  0.029899266
(2) E: odds ratio of being assigned to each latent cluster for each exposure
      beta        OR
hs hg m scaled.cluster2    1.8300519 6.2342102
e3 sex None cluster2      -0.1045747 0.9007075
hs child age yrs None cluster2 -0.4683347 0.6260439
-----Summary of the Middle Component of the LUCID in Serial Model-----
-----Summary of the LUCID Early Integration Sub Model-----
K = 2 , log likelihood = -6035.014 , BIC = 13435.13
(1) Z: mean of omics data for each latent cluster
      mu_cluster1  mu_cluster2
tc_TC01006069_nc -0.32822703  0.321670957
tc_SLC9A4        0.10853723 -0.106428335
tc_RAB6C_AS1     -0.24384843  0.195716238
tc_LOC100129029 0.10572458 -0.005932059
tc_BRE          -0.10399915  0.125807476
tc_TC03001220_nc -0.30053137  0.248463710
tc_TC04002114_nc  0.04931376 -0.060764599
tc_TC04002369_nc -0.38991227  0.586539770
tc_BEND4         -0.26952193  0.195503719
tc_SLC9A3        0.26488011 -0.110649505
(2) E: odds ratio of being assigned to each latent cluster for each cluster
      from the last sub model
      delta        OR
G1.cluster2 1.554258 4.731575
-----Summary of the Last Component of the LUCID in Serial Model-----
-----Summary of the LUCID Early Integration Sub Model-----
K = 2 , log likelihood = -4533.659 , BIC = 10444.5
(1) Y (continuous outcome): effect size of Y for each latent cluster
(and effect of covariates if included)
      Gamma
cluster1        0.0000000
cluster2        1.64219533
e3 sex None     0.06626907
hs child age yrs None -0.16063912
(2) Z: mean of omics data for each latent cluster
      mu_cluster1  mu_cluster2
miR.101.3p     0.12404049  0.028070959
miR.125a.5p    -0.19583201  0.088128200
miR.125b.1.3p   -0.02249475  0.064270997
miR.127.3p     -0.21746619  0.215300195
miR.140.5p     0.27069168  0.017292491

```

```

miR.142.3p      0.15068775 -0.015944126
miR.144.5p      0.25253765 -0.071952430
miR.19a.3p       0.15974513 -0.015887823
miR.19b.3p       0.20886420 -0.008073847
miR.21.5p        0.18868583 -0.017110888
(3) E: odds ratio of being assigned to each latent cluster for each cluster
      from the last sub model
          delta      OR
G1.cluster2 0.9911389 2.694301
-----Overall Summary of the LUCID in Serial model-----
log likelihood = -16546.52 , BIC = 37200.41
> # LUCID in serial with the first sub-model being LUCID in parallel
  and the second sub-model being LUCID early integration
> # Rearrange the omics list to match the new structure of the LUCID in serial model
> Z_new = list(list(omics_lst$methylome, omics_lst$transcriptome), omics_lst$miRNA)
> # Fit the LUCID in serial model with the new omics list and the new K list
> fit19 <- lucid(G = exposome, Z = Z_new, Y = ck18,
+                  lucid_model = "serial", CoY = covs, CoG = covs,
+                  K = list(list(2, 2), 2), useY = TRUE, family = "normal")
> # summary(fit19)

```

## 4 Conclusion

In this paper, we have introduced the **LUCIDus** package, version 3, to perform estimation of the LUCID model in R (Jia et al., 2024). **LUCIDus** focuses on integrative clustering analysis using multi-omics data, both with and without phenotypic traits. It consists of a set of toolkits for modeling, interpretation, visualization, inference, and prediction. Compared to the previous version, **LUCIDus** version 3 is faster (120 times faster than the previous version for a dataset with 20,000 observations), more stable, and includes several features for more functionality, including options for early, intermediate, and late integration of multi-omics data. The **LUCIDus** package is a useful tool for performing integrated clustering analysis and can potentially provide greater insights into environmental epidemiology studies with measured multi-omics data.

## References

- J. M. Albert, C. Geng, and S. Nelson. Causal mediation analysis with a latent mediator. *Biometrical Journal*, 58(3):535–548, 2016. [p5]
- B. Aslam, M. Basit, M. A. Nisar, M. Khurshid, and M. H. Rasool. Proteomics: technologies and their applications. *Journal of Chromatographic Science*, 55(2):182–196, 2017. [p4]
- J. D. Banfield and A. E. Raftery. Model-based gaussian and non-gaussian clustering. *Biometrics*, pages 803–821, 1993. [p7]
- R. D. Beger. A review of applications of metabolomics in cancer. *Metabolites*, 3(3):552–574, 2013. [p4]
- S. Cadiou, X. Basagaña, J. R. Gonzalez, J. Lepeule, M. Vrijheid, V. Siroux, and R. Slama. Performance of approaches relying on multidimensional intermediary data to decipher causal relationships between the exposome and health: A simulation study under various causal structures. *Environment International*, 153:106509, 2021. ISSN 0160-4120. doi: <https://doi.org/10.1016/j.envint.2021.106509>. URL <https://www.sciencedirect.com/science/article/pii/S0160412021001343>. [p15]
- A. Canty and B. Ripley. *boot: Bootstrap Functions (Originally by Angelo Canty for S)*, 2021. URL <https://CRAN.R-project.org/package=boot>. R package version 1.3-28. [p16]
- G. Celeux and G. Govaert. Gaussian parsimonious clustering models. *Pattern Recognition*, 28(5):781–793, 1995. [p7]
- C. Chan, F. Feng, J. Ottinger, D. Foster, M. West, and T. B. Kepler. Statistical mixture modeling for cell subtype identification in flow cytometry. *Cytometry Part A: The Journal of the International Society for Analytical Cytology*, 73(8):693–701, 2008. [p4]
- O. M. Crook, C. M. Mulvey, P. D. Kirk, K. S. Lilley, and L. Gatto. A bayesian mixture modelling approach for spatial proteomics. *PLoS Computational Biology*, 14(11):e1006516, 2018. [p4]

- C. Curtis, S. P. Shah, S.-F. Chin, G. Turashvili, O. M. Rueda, M. J. Dunning, D. Speed, A. G. Lynch, S. Samarajiwa, Y. Yuan, et al. The genomic and transcriptomic architecture of 2,000 breast tumours reveals novel subgroups. *Nature*, 486(7403):346–352, 2012. [p4]
- A. C. Davison and D. V. Hinkley. *Bootstrap methods and their application*. Cambridge University Press, 1997. [p15]
- A. Derkach, R. M. Pfeiffer, T.-H. Chen, and J. N. Sampson. High dimensional mediation analysis with latent variables. *Biometrics*, 75(3):745–756, 2019. [p5]
- Y. Fan and C. Y. Tang. Tuning parameter selection in high dimensional penalized likelihood. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 75(3):531–552, 2013. [p13]
- M. Fop and T. B. Murphy. Variable selection methods for model-based clustering. *Statistics Surveys*, 12: 18–65, 2018. [p14]
- J. R. González and A. Cáceres. *Omic association studies with R and Bioconductor*. CRC Press, 2019. [p4]
- S. Goodwin, J. D. McPherson, and W. R. McCombie. Coming of age: ten years of next-generation sequencing technologies. *Nature Reviews Genetics*, 17(6):333–351, 2016. [p4]
- B. P. Hejblum, J. Skinner, and R. Thiébaut. Time-course gene set analysis for longitudinal gene expression data. *PLoS Computational Biology*, 11(6):e1004310, 2015. [p4]
- B. P. Hejblum, C. Alkhassim, R. Gottardo, F. Caron, and R. Thiébaut. Sequential dirichlet process mixtures of multivariate skew t-distributions for model-based clustering of flow cytometry data. 2019. [p4]
- Q. Jia, Y. Zhao, D. Conti, and J. Goodrich. *LUCIDus: LUCID with Multiple Omics Data*, 2024. URL <https://CRAN.R-project.org/package=LUCIDus>. R package version 3.0.2. [p5, 23]
- R. Jin, R. McConnell, C. Catherine, S. Xu, D. I. Walker, N. Stratakis, D. P. Jones, G. W. Miller, C. Peng, D. V. Conti, et al. Perfluoroalkyl substances and severity of nonalcoholic fatty liver in children: an untargeted metabolomics approach. *Environment International*, 134:105220, 2020. [p5]
- J. Khan, J. S. Wei, M. Ringner, L. H. Saal, M. Ladanyi, F. Westermann, F. Berthold, M. Schwab, C. R. Antonescu, C. Peterson, et al. Classification and diagnostic prediction of cancers using gene expression profiling and artificial neural networks. *Nature Medicine*, 7(6):673–679, 2001. [p4]
- P. Kirk, J. E. Griffin, R. S. Savage, Z. Ghahramani, and D. L. Wild. Bayesian correlated clustering to integrate multiple datasets. *Bioinformatics*, 28(24):3290–3297, 2012. [p5]
- E. F. Lock and D. B. Dunson. Bayesian consensus clustering. *Bioinformatics*, 29(20):2610–2616, 2013. [p5]
- E. F. Lock, K. A. Hoadley, J. S. Marron, and A. B. Nobel. Joint and individual variation explained (jive) for integrated analysis of multiple data types. *The Annals of Applied Statistics*, 7(1):523, 2013. [p5]
- L. Maitre, J.-B. Guimbaud, C. Warembourg, N. Güil-Oumrait, T. E. D. C. P. Consortium, P. M. Petrone, M. Chadeau-Hyam, M. Vrijheid, J. R. Gonzalez, and X. Basagaña. State-of-the-art methods for exposure-health studies: results from the exposome data challenge event. *arXiv preprint arXiv:2202.01680*, 2022. [p5]
- K. Matta, T. Lefebvre, E. Vigneau, V. Cariou, P. Marchand, Y. Guitton, A.-L. Royer, S. Ploteau, B. Le Bizec, J.-P. Antignac, et al. Associations between persistent organic pollutants and endometriosis: A multiblock approach integrating metabolic and cytokine profiling. *Environment International*, 158: 106926, 2022. [p5]
- Q. Mo, S. Wang, V. E. Seshan, A. B. Olshen, N. Schultz, C. Sander, R. S. Powers, M. Ladanyi, and R. Shen. Pattern discovery and cancer gene identification in integrated cancer genomic data. *Proceedings of the National Academy of Sciences*, 110(11):4245–4250, 2013. [p5]
- F. Ozsolak and P. M. Milos. Rna sequencing: advances, challenges and opportunities. *Nature Reviews Genetics*, 12(2):87–98, 2011. [p4]
- W. Pan and X. Shen. Penalized model-based clustering with application to variable selection. *Journal of Machine Learning Research*, 8(May):1145–1164, 2007. [p14]
- C. Peng, J. Wang, I. Asante, S. Louie, R. Jin, L. Chatzi, G. Casey, D. C. Thomas, and D. V. Conti. A latent unknown clustering integrating multi-omics data (lucid) with phenotypic traits. *Bioinformatics*, 36 (3):842–850, 2020. [p5, 6]

- M. Pierre-Jean, J.-F. Deleuze, E. Le Floch, and F. Mauger. Clustering and variable selection evaluation of 13 unsupervised methods for multi-omics data integration. *Briefings in Bioinformatics*, 21(6):2011–2030, 2020. [p<sup>5</sup>]
- S. Prabhakaran, E. Azizi, A. Carr, and D. Pe'er. Dirichlet process mixture model for correcting technical variation in single-cell gene expression data. In *International Conference on Machine Learning*, pages 1070–1079. PMLR, 2016. [p<sup>4</sup>]
- N. Rappoport and R. Shamir. Multi-omic and multi-view clustering algorithms: review and cancer benchmark. *Nucleic Acids Research*, 46(20):10546–10562, 2018. [p<sup>4</sup>]
- M. Schmidt. The sankey diagram in energy and material flow management: part ii: methodology and current applications. *Journal of Industrial Ecology*, 12(2):173–185, 2008. [p<sup>12</sup>]
- L. Scrucca, M. Fop, T. B. Murphy, and A. E. Raftery. mclust 5: clustering, classification and density estimation using gaussian finite mixture models. *The R Journal*, 8(1):289, 2016. [p<sup>7</sup>, 10]
- A. Singh, C. P. Shannon, B. Gautier, F. Rohart, M. Vacher, S. J. Tebbutt, and K.-A. Lê Cao. Diablo: an integrative approach for identifying key molecular drivers from multi-omics assays. *Bioinformatics*, 35(17):3055–3062, 2019. [p<sup>5</sup>]
- Y. Song, X. Zhou, M. Zhang, W. Zhao, Y. Liu, S. L. Kardia, A. V. D. Roux, B. L. Needham, J. A. Smith, and B. Mukherjee. Bayesian shrinkage estimation of high dimensional causal mediation effects in omics studies. *Biometrics*, 76(3):700–710, 2020. [p<sup>5</sup>]
- N. Stratakis, D. V. Conti, R. Jin, K. Margetaki, D. Valvi, A. P. Siskos, L. Maitre, E. Garcia, N. Varo, Y. Zhao, et al. Prenatal exposure to perfluoroalkyl substances associated with increased susceptibility to liver injury in children. *Hepatology*, 72(5):1758–1770, 2020. [p<sup>5</sup>]
- I. Subramanian, S. Verma, S. Kumar, A. Jere, and K. Anamika. Multi-omics data integration, interpretation, and its application. *Bioinformatics and Biology Insights*, 14:1177932219899051, 2020. [p<sup>4</sup>]
- A. Tenenhaus, C. Philippe, V. Guillemot, K.-A. Le Cao, J. Grill, and V. Frouin. Variable selection for generalized canonical correlation analysis. *Biostatistics*, 15(3):569–583, 2014. [p<sup>5</sup>]
- R. Tibshirani. Regression selection and shrinkage via the lasso. *Journal of the Royal Statistical Society: Series B*, 58(1):267–288, 1996. [p<sup>5</sup>]
- N. Tierney. *visdat: Visualising Whole Data Frames Including Missing Values*, 2023. URL <https://cran.r-project.org/package=visdat>. R package version 0.5.3. [p<sup>18</sup>]
- G. Tini, L. Marchetti, C. Priami, and M.-P. Scott-Boyer. Multi-omics integration—a comparison of unsupervised clustering methodologies. *Briefings in Bioinformatics*, 20(4):1269–1279, 2019. [p<sup>4</sup>]
- V. Voillet, P. Besse, L. Liaubet, M. San Cristobal, and I. González. Handling missing rows in multi-omics data integration: multiple imputation in multiple factor analysis framework. *BMC Bioinformatics*, 17(1):1–16, 2016. [p<sup>16</sup>]
- M. Vrijheid, R. Slama, O. Robinson, L. Chatzi, M. Coen, P. Van den Hazel, C. Thomsen, J. Wright, T. J. Athersuch, N. Avellana, et al. The human early-life exposome (helix): project rationale and design. *Environmental Health Perspectives*, 122(6):535–544, 2014. [p<sup>4</sup>, 5]
- B. Wang, A. M. Mezlini, F. Demir, M. Fiume, Z. Tu, M. Brudno, B. Haibe-Kains, and A. Goldenberg. Similarity network fusion for aggregating data types on a genomic scale. *Nature Methods*, 11(3):333–337, 2014. [p<sup>5</sup>]
- C. P. Wild. Complementing the genome with an “exposome”: the outstanding challenge of environmental exposure measurement in molecular epidemiology. *Cancer Epidemiology, Biomarkers & Prevention*, 14(8):1847–1850, 2005. [p<sup>4</sup>]
- G. Yu. *dlstats: Download Stats of R Packages*, 2022. URL <https://CRAN.R-project.org/package=dlstats>. R package version 0.1.5. [p<sup>5</sup>]
- M. Yuan and Y. Lin. Model selection and estimation in regression with grouped variables. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 68(1):49–67, 2006. [p<sup>5</sup>]
- H. Zhang, Y. Zheng, Z. Zhang, T. Gao, B. Joyce, G. Yoon, W. Zhang, J. Schwartz, A. Just, E. Colicino, et al. Estimating and testing high-dimensional mediation effects in epigenetic studies. *Bioinformatics*, 32(20):3150–3154, 2016. [p<sup>5</sup>]

Y. Zhang, M. Li, S. Wang, S. Dai, L. Luo, E. Zhu, H. Xu, X. Zhu, C. Yao, and H. Zhou. Gaussian mixture model clustering with incomplete data. *ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM)*, 17(1s):1–14, 2021. [p18]

H. Zhou, W. Pan, and X. Shen. Penalized model-based clustering with unconstrained covariance matrices. *Electronic Journal of Statistics*, 3:1473, 2009. [p14]

*Yinqi Zhao*

*Division of Biostatistics, Department of Population and Public Health Sciences, University of Southern California*

*2001 N Soto St., Los Angeles, CA, 90032*

*United States*

*ORCID: 0000-0003-2413-732X*

[yinqiz@usc.edu](mailto:yinqiz@usc.edu)

*Qiran Jia*

*Division of Biostatistics, Department of Population and Public Health Sciences, University of Southern California*

*2001 N Soto St., Los Angeles, CA, 90032*

*United States*

*ORCID: 0000-0002-0790-5967*

[qiranjia@usc.edu](mailto:qiranjia@usc.edu)

*Jesse A. Goodrich*

*Division of Environmental Epidemiology, Department of Population and Public Health Sciences, University of Southern California*

*2001 N Soto St., Los Angeles, CA, 90032*

*United States*

*ORCID: 0000-0001-6615-0472*

[jagoodri@usc.edu](mailto:jagoodri@usc.edu)

*David V. Conti*

*Division of Biostatistics, Department of Population and Public Health Sciences, University of Southern California*

*2001 N Soto St., Los Angeles, CA, 90032*

*United States*

*ORCID: 0000-0002-2941-7833*

[dconti@usc.med.edu](mailto:dconti@usc.med.edu)

# Current State and Prospects of R-Packages for the Design of Experiments

by Emi Tanaka and Dewi Amaliah

**Abstract** Re-running an experiment is generally costly and, in some cases, impossible due to limited resources; therefore, the design of an experiment plays a critical role in increasing the quality of experimental data. In this paper, we describe the current state of R-packages for the design of experiments through an exploratory data analysis of package downloads, package metadata, and a comparison of characteristics with other topics. We observed that experimental designs in practice appear to be sufficiently manufactured by a small number of packages, and the development of experimental designs often occurs in silos. We also discuss the interface designs of widely utilized R packages in the field of experimental design and discuss their future prospects for advancing the field in practice.

## 1 Introduction

The critical role of data collection is well captured in the expression “garbage in, garbage out” – in other words, if the collected data are rubbish, then no analysis, however complex it may be, can make something out of it. Therefore, a carefully crafted data-collection scheme is critical for optimizing the information from the data. The field of experimental design is specifically devoted to planning the collection of experimental data, largely based on the founding principles of Fisher (1935) or an optimization framework like those described in Pukelsheim (2006). These experimental designs are often constructed with the aid of statistical software such as R (R Core Team 2021), Python (Rossum 1995), and SAS (SAS Institute 1985); thus the use of experimental design software can inform us about some aspects of experimental designs in practice.

Methods for data collection can be dichotomized by the type of data collected – namely, experimental or observational – or alternatively, categorized as experimental design (including quasi-experimental design) or survey design. This dichotomization, to a great extent, is seen in the [Comprehensive R Archive Network \(CRAN\)](#) task views (a volunteer maintained list of R-packages by topic) where R-packages for experimental design are in the [ExperimentalDesign](#) task view and R-packages for survey designs are in the [OfficialStatistics](#) task view. A full list of available topics is provided in Table S1 in the Supplementary Materials. A subset of experimental designs is segregated into the [ClinicalTrials](#) task view, where the focus is on clinical trials with primary interest in sample size calculations. This paper focuses on packages in the [ExperimentalDesign](#) task view, henceforth referred to as “DoE packages”.

From the [ExperimentalDesign](#) task view, there are 105 R packages for the experimental design and analysis of data from experiments. The sheer quantity and variation of experimental designs in the R-packages are arguably unmatched by any other programming languages; for example, in Python, only a handful of packages that generate design of experiments exist (namely `pyDOE`, `pyDOE2`, `dexpy`, `experimenter`, and `GPdoemd`) with a limited type of design. Thus, the study of DoE packages, based on quantitative and qualitative data, can provide an objective view of the state of current experimental designs in practice.

The utility of the software can also be described by its design to facilitate the clear expression and interpretation of the desired experimental design. Certain programming language designs can hinder or discourage the development of reliable programs (Wasserman 1975). The immense popularity of [tidyverse](#) (a collection of R-packages for various stages of data analysis that places enormous emphasis on the interface design by Wickham et al. 2019) is a testament to the impact that an interface design can have in practice. The practice of experimental design can be advanced by adopting similar interface design principles across the DoE packages.

The remainder of this paper is organized as follows. [Data](#) briefly describes the data source used for the analysis; [Exploratory data analysis](#) presents some insights into the state of the current DoE packages by the exploratory data analysis of package download data, text descriptions, and comparisons with other CRAN task views; [Interface design](#) discusses the interface designs of widely used DoE packages, and we conclude with a discussion in [Discussion](#) of future prospects in the software development of experimental designs.

## 2 Data

To study the DoE packages, we analyze data using three sources of data as described below.

### 2.1 RStudio CRAN download logs

The Comprehensive R Archive Network (CRAN) is a network of servers located across the world that stores mirrored versions of R and R packages. The most popular network is the RStudio mirror (the default server for those that use the RStudio IDE). The RStudio mirror is also the only server that provides comprehensive daily download logs of R and R packages since October 2012. The summary data can be easily accessed using the [cranlogs](#) package (Csárdi 2019). This paper uses the data from the beginning of 2013 to the end of 2021 (a total of nine years) for the packages in the CRAN task views.

### 2.2 Package descriptions

All CRAN packages have a title, description, package connections (suggests, depends, and imports of other packages), and other meta-information in the DESCRIPTION file. We use text data from the title and description (accessed on 2022-12-12).

### 2.3 CRAN task views

CRAN task views are volunteer-maintained lists of R-packages on CRAN relevant to the corresponding topic. There were 39 CRAN task views in total. Table S1 in the Supplementary Materials lists the available topics from the [ctv](#) package (Zeileis 2005). The list of packages in each CRAN task view (as of 2022-12-12) is used to contrast the characteristics of the DoE packages.

## 3 Exploratory data analysis

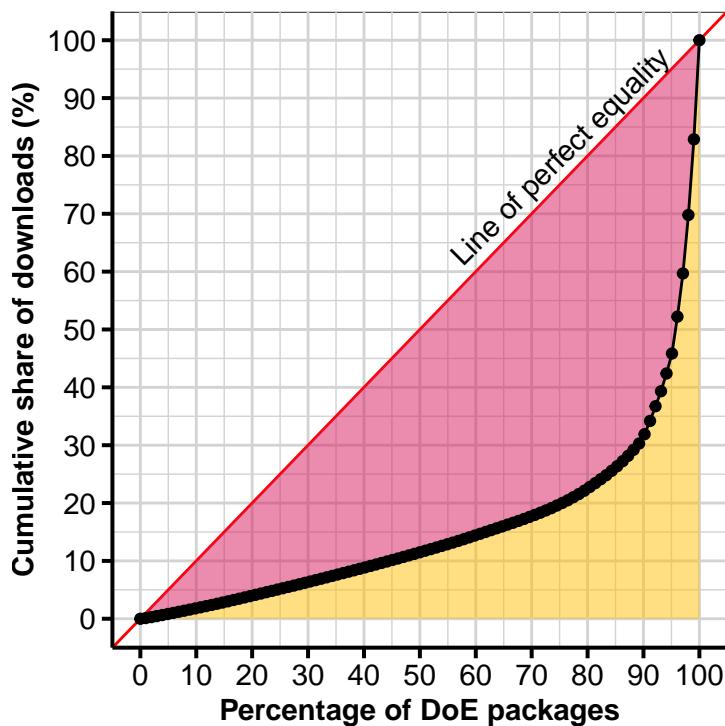
In this section, we derive some conjectures based on an analysis of the data described in [Data](#). All results presented are from exploratory data analysis of observational data; consequently, all interpretations are somewhat speculative and may not be indicative of the true state of the field of experimental design. In particular, any analysis over time is confounded by the fact that the nature of users and package management has changed over the years. It should be noted that some DoE packages may have been archived or removed from the task view over the years; therefore, any cross-sectional analysis presented may not reflect the set of all DoE packages at that particular time period (although we assume such incidents are low).

A subset of DoE packages is not primarily about the design of experiments but about the analysis of experimental data. A complete delineation of these packages is difficult, as there is almost always at least one function that can aid decisions or constructions of experimental designs (and any categorization is prone to our subjective bias); therefore, we opted not to remove any DoE packages in the analysis.

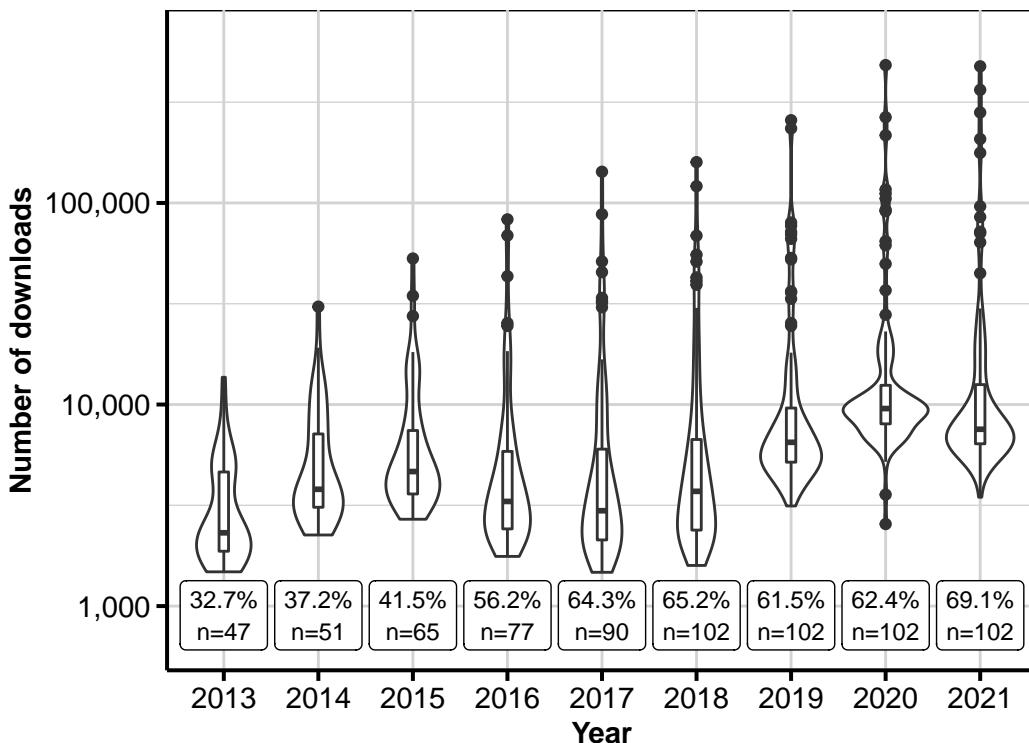
### 3.1 Small, but diverse, set of packages are sufficient for most experimental designs in practice

There are at least 50 DoE packages since 2013, but most of the downloads are concentrated in only a handful of packages. For example, Figure 1 shows a Lorenz curve (Lorenz 1905) for the total package downloads in 2021 for 102 DoE packages (first released prior to 2021). We can see from Figure 1 that the bottom 90% of DoE packages (in terms of total download count in 2021) only share approximately 32% of total downloads across all DoE packages; in other words, 68% of the total downloads are due to 10 packages (10% of the DoE packages).

If we consider package downloads as a measure of “wealth”, then we can consider using the Gini index (Gini 1921) as a measure of download inequality across packages. The ratio of the red region to the total colored regions in Figure 1 corresponds to the Gini index for 2021. A Gini index of 0% indicates equality in downloads across packages, whereas a value of 100% indicates maximal inequality (all downloads are due to one package). In Figure 2, we see that the distributions of the package downloads each year have a heavy right tail, with the Gini index ranging from 32.7% to 69.1%



**Figure 1:** Lorenz curve of the total download count for DoE packages in 2021. The red line corresponds to the line of perfect equality. The yellow region shows the area under the Lorenz curve and the red region shows the area of the gap in equality.



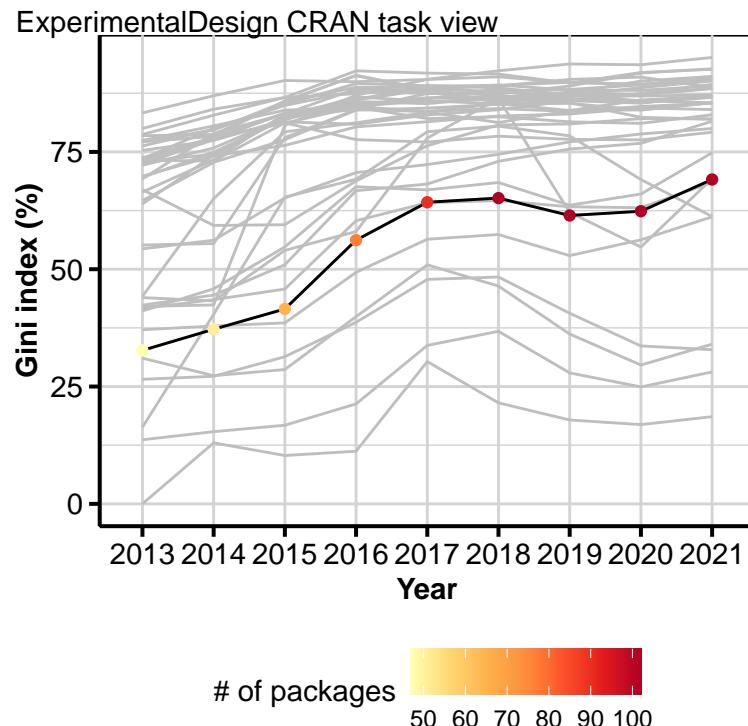
**Figure 2:** Distribution of number of downloads for DoE packages by year. Packages were removed in any year if they were released in that year or later so that each download count was for the full year. The label at the bottom of the plot shows the Gini index for downloads and the number of packages with a full download count in the corresponding year. In the last six years, the Gini index has consistently exceeded 60%, indicating that most downloads are due to a relatively small number of packages.

across the years 2013 to 2021, indicating that there is a high level of inequality in package downloads, particularly with more pronounced inequality in the last six years.

An increase in the number of packages that are not highly downloaded may mean that *there are more packages to construct niche experimental designs*. Some examples of these packages include `qltDesign`, `PwrGSD` and `Crossover` made for QTL experiments, group sequential designs and crossover trials, respectively. These packages would naturally have fewer potential users. Counterfactual to this, the increase could be due to other external factors, such as an increase in the number of skilled developers (and thus more package contributions), a change in CRAN policy or management to add packages (either to CRAN and/or task view), and/or the fact that new packages are still yet to amass users. While there is an argument that low download counts are due to the low utility and/or quality of packages, packages in CRAN task views are selected by expert maintainers. We can reasonably assume that any package listed in the CRAN task view has an acceptable utility and quality.

If the downloads are reflective of the experimental designs used in practice, *a small set of packages appears to be sufficient for most users to construct the full set of designs of experiments they need in practice*. Packages of course evolve, and the top downloaded packages have had regular updates that may have broadened their scope from their previous releases.

While in absolute terms the Gini index is high for the `ExperimentalDesign` task view (32.7% to 69.1%), the inequality is not as severe as in other CRAN task views, as shown in Figure 3. We can see in Figure 3 that the Gini index generally increases over time for DoE packages (as is generally the case for other CRAN task views, as shown in Figure S1 in the Supplementary Materials), but most other CRAN task views have a Gini index of over 75%. This suggests that other CRAN task views may have dominant standards, and in comparison to other topics, there are more *diverse approaches to designing experiments*, and thus, no single DoE package is dominant. However, this observation does not consider other approaches to generate experimental designs, such as the proprietary software CycDesignN (Whittaker, Williams, and John 2022), which may be widely used.

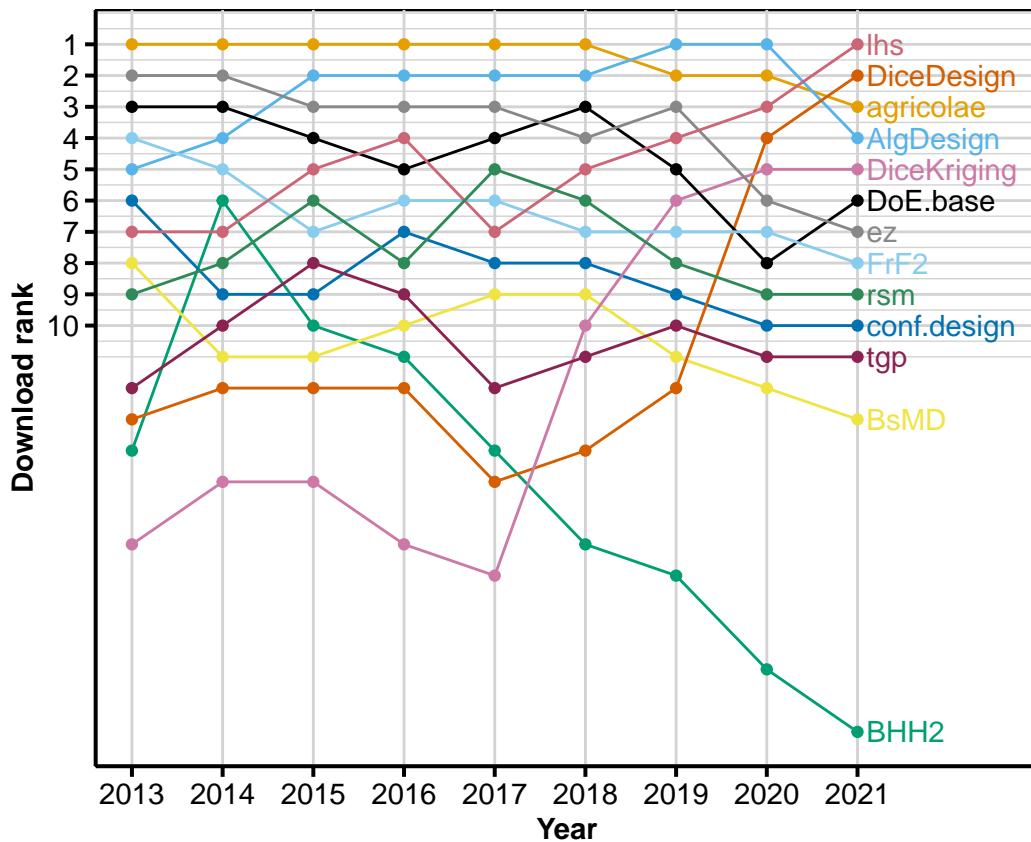


**Figure 3:** The points show the Gini index of the download counts by year for the `ExperimentalDesign` task view with the color showing the number of packages. The gray lines show the line plots of the Gini index across years for all other CRAN task views. See Figure S1 in the Supplementary Material for the line graph of the Gini index across years for each CRAN task view.

### 3.2 The field of experimental design is slow-changing

We can see in Figure 4 that most of the top 10 ranking packages have been in the top 10 for the last nine years, with `lhs` steadily climbing up the ranks in the last few years. It should be noted that the

download of one package can prompt the download of another package; the most notable package connection is `AlgDesign` and `agricolae`, where the former is an import for the latter. The full network of package connections within the DoE packages is shown in Figure 8.



**Figure 4:** The plot shows the rank of the top 10 packages downloaded by year. Packages that did not appear in the top 10 for at least two periods were omitted from the plot. Most packages are consistently in the top 10 for the period shown.

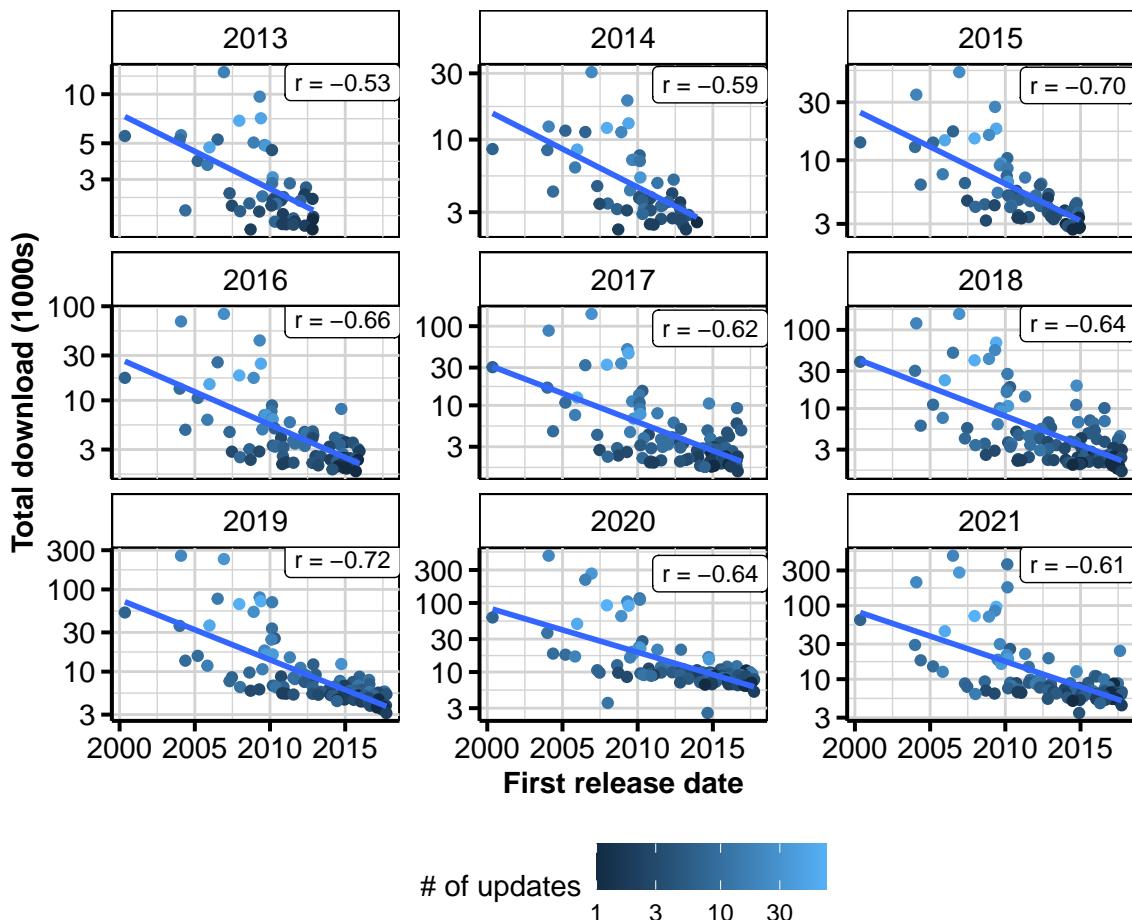
Figure 5 shows a moderate negative correlation between the first release date and the (log of) total download counts of DoE packages for any given year from 2013 to 2021. This suggests that packages released earlier are more likely to be used today (possibly for legacy reasons or the general inertia to adopt new packages). We can also see in Figure 5 that most downloaded packages were released from 2004 to 2010.

The consistency in the top 10 ranking packages (Figure 4) and the fact that most downloaded DoE packages were first released more than 10 years ago (Figure 5) indicate that either the existing packages fulfilled the needs of the masses in practice or no new packages were compelling for many to switch their practice. However, we also see that the top downloaded packages generally have more updates (see Figure 5); therefore, it is possible that the packages have improved or broadened the scope of their usage.

### 3.3 Optimal designs are of interest

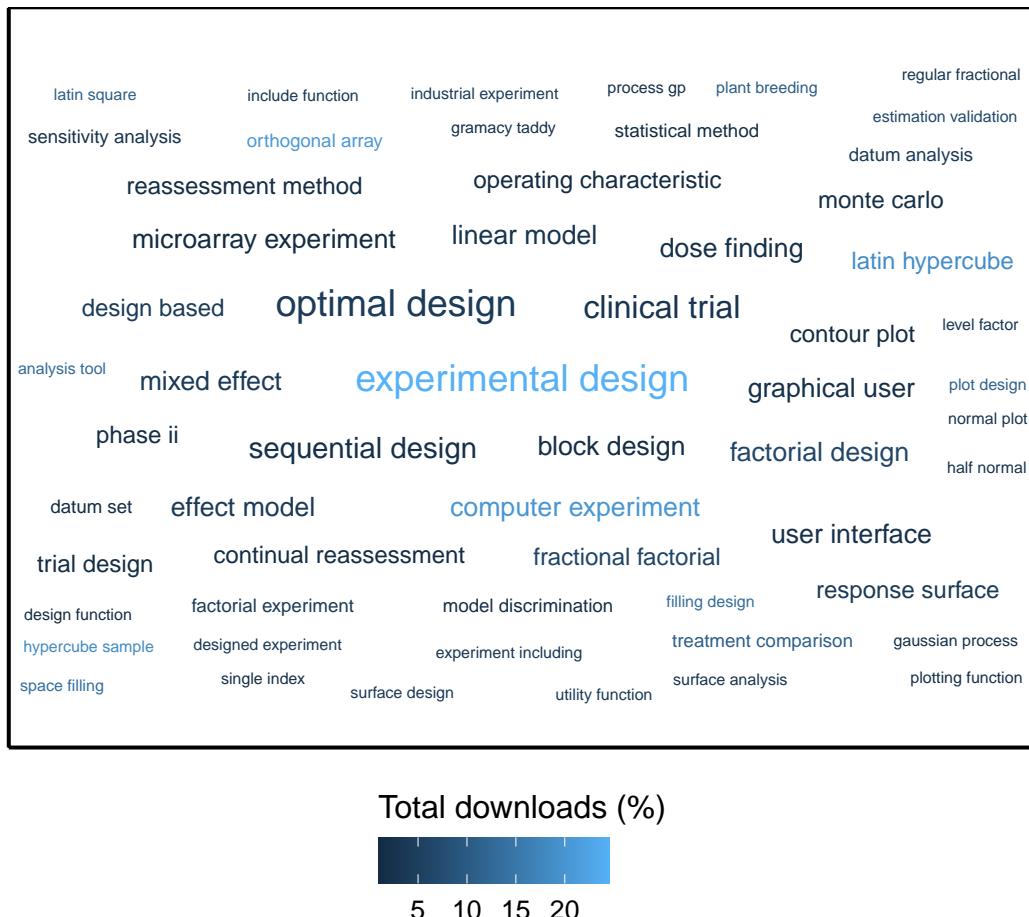
Figure 6 shows some of the common purposes of DoE packages, based on bigrams in the package title and description. We only show the bigrams as unigrams were not insightful, and there were not many trigrams common across packages. To count the bigrams, we processed the text data as follows:

1. We standardized the words to lower case and removed pluralization.
2. Multiple mentions of the same bigram within a package were counted as one (for example, `AlgDesign` mentions “experimental design” four times in the title and the description, but this is counted as one).
3. Bigrams consisting of stop words were removed. The stop words are sourced from the lexicons in `tidytext::stop_words` in addition to other words we deemed irrelevant, e.g., “provide”, “e.g.”, “calculate” and so on – the full list is shown in the code provided in the link under Acknowledgment.



**Figure 5:** The above figure shows the total download (in log scale) of a package in the corresponding year against the first release date of the package. The blue line corresponds to the least-squares fit of the simple linear regression model. The label in the upper right-hand corner shows the sample correlation coefficient between the first release date and log (with base 10) of the total download count. The high leverage point on the far left belongs to 'conf.design', authored by one of the earlier contributors to R.

Unsurprisingly, the bigram “experimental design” was the most common. More interestingly, “optimal design” and “sequential design” appeared across different packages (indicated by the size of the word in Figure 6), and the bigrams “latin hypercube” and “computer experiment” are used across a few packages that are downloaded frequently (indicated by the color of the word in Figure 6). Sequential design, Latin hypercube sampling, and computer experiments (which generally include space-filling designs such as Latin hypercube sampling) generally operate by optimizing a user-selected criterion and can be classified as optimal designs.



**Figure 6:** The above figure shows the word cloud of bigrams from the title and descriptions of the DoE packages. The size shows how often the bigram appears across the DoE packages and the color is relative to the total download count in 2021 for the packages that contain the bigram.

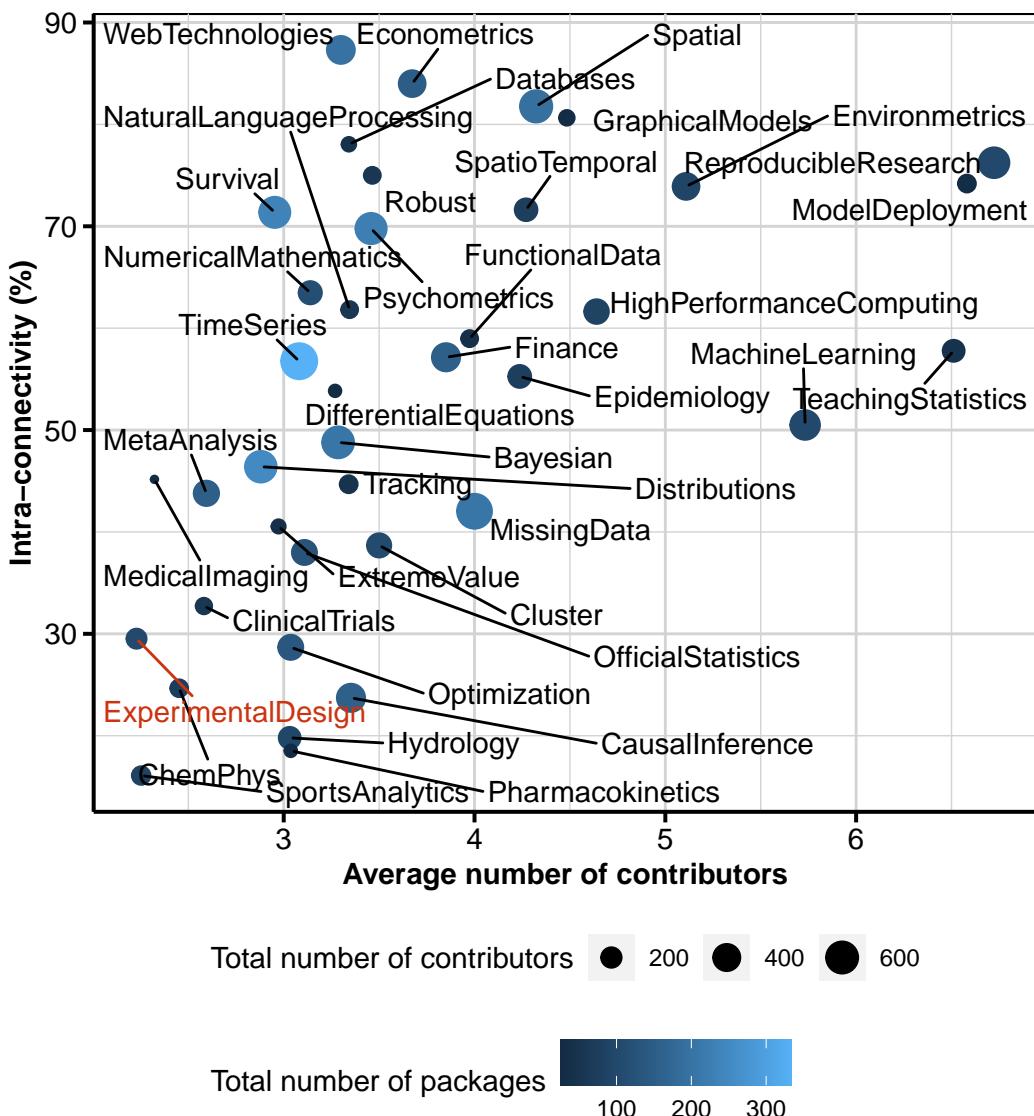
Although there exists a separate *ClinicalTrials* task view, the DoE packages clearly include some packages that are of interest to clinical trials, as shown by the size of the bigram “clinical trial” (and related bigrams like “dose finding” and “phase ii”) in Figure 6.

### 3.4 Development of experimental designs occurs in silos

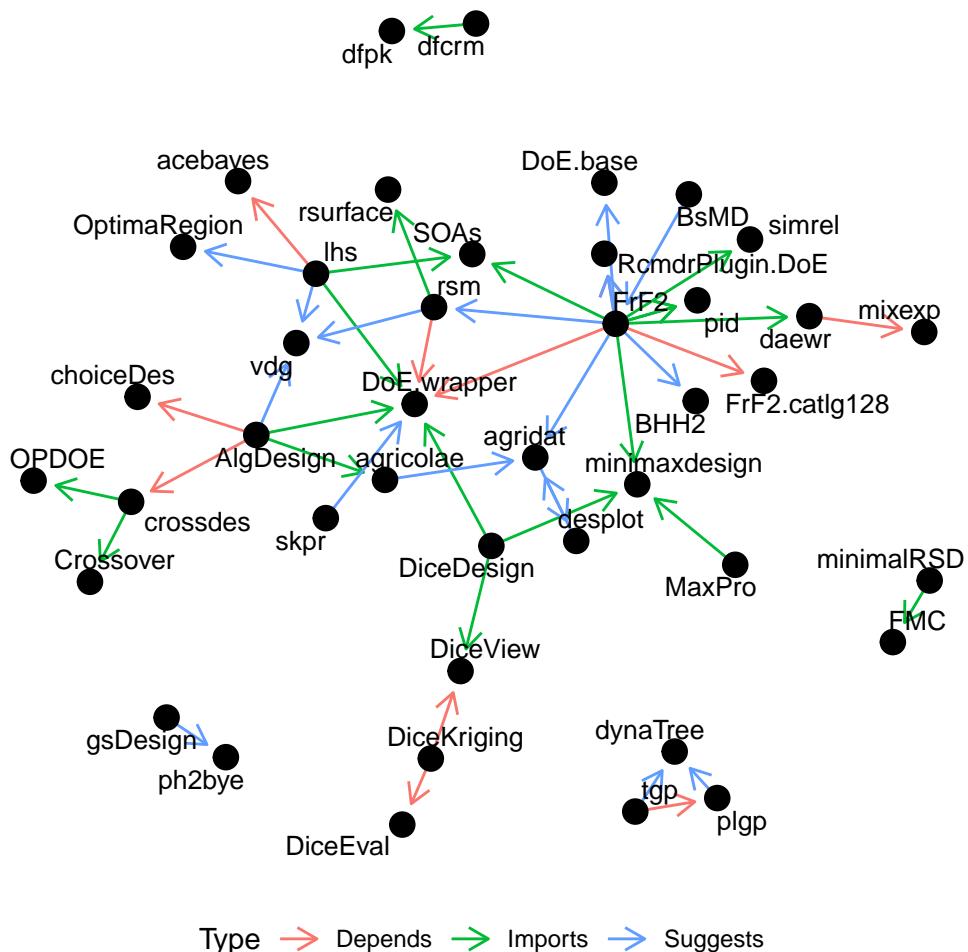
Figure 7 shows that the *ExperimentalDesign* task view has the lowest average number of contributors among all the 39 CRAN task views. In addition, we can also see in Figure 7 that the *ExperimentalDesign* task view has one of the least intra-connectivity (the percentage of packages that make use of other packages within the same task view). The full connection between DoE packages is shown in Figure 8. These observations suggest that *experimental design is one of the least collaborative fields* and package development generally occurs in silos.

## 4 Interface design

In software design, there are two interface designs to consider: user interface (UI) and application programming interface (API). The UI is concerned with the interaction of the software by the user, while the API is concerned with how different programs interact and is predominantly of interest to



**Figure 7:** The above figure is a scatterplot of intra-connectivity (the percentage of packages that depend, suggest, or import at least one other package within the same task view) and the average number of contributors for each CRAN task view. Low intra-connectivity suggests that development within the topic mostly occurs in silos, while high intra-connectivity suggests that there are more interactions within the topic. The color shows the number of packages, the size of the point corresponds to the total number of contributors, and the text labels show the CRAN task view names. The label of the ExperimentalDesign task view is colored in red. The task views in the bottom-left corner are topics that are more indicative of contributors working in silos. The actual numerical values are listed in Table S1 in the Supplementary Material.



**Figure 8:** Package connections (depends, suggests, and imports) within DoE packages. The direction of the arrow shows the connection of packages, where the package on the tail of the arrow is a dependency, suggestion, or import for the package on the head of the arrow. DoE packages that do not depend, suggest, or import another DoE package are not shown.

the developer. The UI design is an abstraction that specifies the desired experimental design, and its choices enable how a user expresses the specification of an experimental design. The API design aids other developers in leveraging existing systems.

In this section, we discuss the interface designs of functions that output an experimental design based on three broad areas: factorial, recipe, and augmenting designs. The discussion is exclusive to the top downloaded packages (shown in Figure 4), with the exception of `ez` and `DiceKriging`, as the former is predominantly visualization of experimental data and the latter is about the analysis of computer experiments in addition to belonging to the same suite of packages as `DiceDesign` (Dupuy, Helbert, and Franco 2015).

#### 4.1 The case of factorial designs

Factorial experiments offer a challenge in allocating the treatment factors to experimental units, where the full set of factorial treatments cannot be administered (and replicated), and/or the experimental units have a grouping structure. The effort to address this challenge is reflected in the number of packages that focus on the construction of factorial designs, as described next.

The `DoE.base` package (Grömping 2018) can construct full factorial and (regular and irregular) orthogonal array designs via `fac.design()` and `oa.design()`, respectively. The `FrF2` package (Grömping 2014) constructs regular fractional 2-level factorial designs using `FrF2()` and `FrF2Large()`, with the latter for large designs. The `conf.design` package (Venables 2013) constructs symmetric confounded factorial designs via `conf.design()` and the `BHH2` package (Barrios 2016) generates a full or fractional 2-level factorial design matrix via `ffDesMatrix()`.

While the argument names and underlying algorithms of these functions to generate factorial designs differ, it generally requires the users to input the number of:

- treatment factors,
- experimental runs (or replications), and
- levels for each factor if the design was allowed to vary in the number of levels.

The output of the design is either a special class of list (e.g., the `design` class for `DoE.base` and `FrF2`) or a `data.frame` for `conf.design` or a matrix for `BHH2`, such that an element or column corresponds to a treatment factor with each value corresponding to one experimental run. The treatment factors are generally assigned pseudo names (e.g., letters of the alphabet) or an argument exists for users to input treatment names as a character vector.

Some forms of factorial design are known by other names. For example, *response surface designs* are factorial designs where the treatment factors are discrete levels of continuous variables. Two types of response surface designs can be constructed using the `rsm` package (Lenth 2009): Box-Behnken design (Box and Behnken 1960) and central-composite designs (Box and Wilson 1951) via functions `bbd()` and `ccd()`, respectively, where the minimum required input is the number of factors. Another form of factorial design is the *saturated designs*, where higher-order interaction effects of treatment factors are typically confounded with the main effects. Plackett-Burman designs (Plackett and Burman 1946) are a type of saturated design that can be generated by the function `pb()` in the `FrF2` package, where the user provides the number of experimental runs and the number of treatment factors.

#### 4.2 The case of recipe designs

The `agricolae` package (de Mendiburu 2021) is the prime example of constructing designs based on a set of so-called “recipe functions”, where each function corresponds to a single class of experimental design. For example, `design.crd()`, `design.rcbd()`, and `design.split()` construct completely randomized, randomized complete block, and split-plot designs, respectively. Users typically supply treatment labels (or the number of treatments in the case of `design.split()`) and the number of replications as arguments for these functions. The output is a list with one element corresponding to a `data.frame` that contains the design in a table such that the row corresponds to the experimental run and the columns correspond to the experimental variables (we refer to this format simply as “table format” henceforth).

The use of recipe functions is not limited to classical experimental designs; the `AlgDesign` (Wheeler 2022) package offers three primary functions for generating optimal designs: `optBlock()`, `optFederov()`, and `optMonteCarlo()`. In general, these functions require data and formulas in terms of the supplied data variables, along with the choice of the criterion (e.g., the D-criterion), with the output as a list with one element corresponding to the design in a table format. The difference between these functions lies in the underlying search strategy for optimal designs, and the name of the function is a surrogate for the search algorithm.

Computer experiments, which generally involve space-filling designs, are implemented in packages such as `lhs` (Carnell 2022) and `DiceDesign` (Dupuy, Helbert, and Franco 2015). For the `lhs` package, functions such as `randomLHS()`, `optimumLHS()`, and `maximinLHS()` require users to specify the sample size ( $n$ ) and the number of variables ( $p$ ). It generates a Latin hypercube sample (McKay, Beckman, and Conover 1979) based on different optimization schemes (in this case, random, S-optimal, and maxmin criteria, respectively; see package documentation for more details). Similarly, for `DiceDesign`, there is a comprehensive list of space-filling designs such as `dmaxDesign()`, `lhsDesign()`, and `wspDesign()` with input of  $n$  and  $p$  as before (among additional parameters for some) that implement algorithms that either maximize the entropy (Shewry and Wynn 1987), produce a random Latin hypercube sample, or use the WSP algorithm (Santiago, Claeys-Bruno, and Sergent 2012), respectively. These designs output an  $n \times p$  matrix with values between 0 and 1. Again, these functions employ a recipe style, in which each function has a name that corresponds to a certain search strategy to generate the experimental design.

### 4.3 The case of augmenting designs

Some functions in the DoE packages require the input of existing experimental designs to produce new designs. For example, the `DoE.base` package contains some experimental functions, `cross.design()` and `param.design()`, to combine designs, with the former taking a Cartesian product of the input designs, while the latter uses a Taguchi style (Taguchi 1986) to aggregate the designs with the inner and outer arrays. The `lhs` package contains a function, `augmentLHS()`, to add additional samples to the existing Latin hypercube sample.

Another class of augmenting design is *sequential design* (also called *adaptive sampling*), which is best represented by `tgp` (Gramacy and Taddy 2010). This requires prior information that is used to inform the next experimental design using `tgp.design()` and `dopt.gp()`. The user is required to supply candidate samples to subsample from and a model or a prior experimental design. Follow-up experiments, which can also be classified as sequential designs, are implemented by `BsMD` (Barrios 2020) using a model-discriminant approach with `MD()`.

## 5 Discussion

Through the exploratory data analysis of the three data sources (package download logs, package metadata, and CRAN task views) outlined in `Data`, we observed in `Exploratory data analysis` that the total download of DoE packages is concentrated only on a handful of R-packages, although these represent a diverse set in comparison with other CRAN task views. Furthermore, the data suggest that experimental design is the least collaborative field.

There are a number of limitations and shortcomings to our exploratory data analysis. First, CRAN task views are volunteer maintained, so some experimental design packages may not be included in the DoE packages. Second, we used only the RStudio CRAN mirror download, which may have biased our observations. Third, our analysis was limited to R-packages alone, and many practitioners may use other methods to construct experimental designs. Finally, all our statements should be treated as speculative rather than conclusive; the data are all observational, so no conclusive, generalizable statement is possible. Regardless, the data-driven nature of our analysis provides objective insight into the field of experimental design.

The interface design (discussed in `Interface design`) reveals that the most widely used DoE packages generally have functions that 1) focus on certain aspects of experimental design (e.g., factorial structure or augmenting design), 2) are in a recipe format (i.e., the name of the function is a surrogate for a single class of the design or optimal search algorithm), and 3) context is often a second thought – many inputs are a single integer corresponding to the number of factors, levels, or experimental runs (or sample size). The function will often assign pseudo-factor names, or there is an optional argument to input a character vector that corresponds to the factor names. These interface designs require users to have processed the experiment in statistical terms (often stripping the experimental context away) and simultaneously, users must choose the generation mechanism by selecting an appropriate function. Arguably, the current dominant interface designs are not aligned with the way practitioners cognitively design their experiments. Often, the critical part of designing an experiment is to understand the experimental structure, and the experimental context can govern or guide the choice of algorithm to allocate treatments. In addition, the nature of recipe functions can obscure the understanding and relation of designs (e.g., how do you go from an unstructured factorial design to a split plot design?). Each new method for generating an experimental design appears to correspond to a completely new function, and intermediate results are often not easily accessible. These factors may contribute to why developers often work in silos in the field of experimental design.

A consistent cognitive interface design that leverages existing developments and can be easily extended to new methods will conceivably be of great help to practitioners. Some efforts to this end are seen in `DoE.wrapper` (Grömping 2020), which contains wrapper recipe functions for other DoE packages such as `lhs`, `AlgDesign`, and `FrF2` (see Figure 8), and is also the subject of the developmental package `edibble` (Tanaka 2021). Undoubtedly, no single developer or package can cater to all experimental designs; therefore, any unifying interface should consider how other developers can contribute or add their methods. Future research could benefit from further exploratory data analysis, expanding the study beyond R-packages, and discussing other aspects of interface designs.

## 6 Acknowledgment

This paper uses `targets` (Landau 2021) and `renv` (Ushey 2022) for reproducibility, `knitr` (Xie 2015) and `rmarkdown` (Xie, Allaire, and Grolemund 2018) for creating reproducible documents, `ggplot2` (Wickham 2016), `ggraph` (Pedersen 2021), `ggwordcloud` (Le Pennec and Slowikowski 2019) and `colorspace` (Zeileis et al. 2020) for visualization, `kableExtra` (Zhu 2021) for customizing the table in the Supplementary Material, and `tidyverse` (Wickham et al. 2019), `tidytext` (Silge and Robinson 2016), `pluralize` (Rudis and Embrey 2020), and `ineq` (Zeileis 2014) for data processing and manipulation, and `cranlogs` (Csárdi 2019) and `cvg` (Zeileis 2005) for extracting data. All codes used to reproduce this paper are found in <https://github.com/emitanaka/paper-DoE-review>.

## References

- Barrios, Ernesto. 2016. *BHH2: Useful Functions for Box, Hunter and Hunter II*. <https://CRAN.R-project.org/package=BHH2>.
- . 2020. *BsMD: Bayes Screening and Model Discrimination*. <https://CRAN.R-project.org/package=BsMD>.
- Box, G E P, and D W Behnken. 1960. “Some New Three Level Designs for the Study of Quantitative Variables.” *Technometrics: A Journal of Statistics for the Physical, Chemical, and Engineering Sciences* 2 (4): 455–75. <https://doi.org/10.2307/1266454>.
- Box, G E P, and K B Wilson. 1951. “On the Experimental Attainment of Optimum Conditions.” *Journal of the Royal Statistical Society. Series B, Statistical Methodology* 13 (1): 1–45.
- Carnell, Rob. 2022. *Lhs: Latin Hypercube Samples*. <https://CRAN.R-project.org/package=lhs>.
- Csárdi, Gábor. 2019. *Cranlogs: Download Logs from the 'RStudio' 'CRAN' Mirror*. <https://CRAN.R-project.org/package=cranlogs>.
- de Mendiburu, Felipe. 2021. *Agricolae: Statistical Procedures for Agricultural Research*. <https://CRAN.R-project.org/package=agricolae>.
- Dupuy, Delphine, Céline Helbert, and Jessica Franco. 2015. “DiceDesign and DiceEval: Two R Packages for Design and Analysis of Computer Experiments.” *Journal of Statistical Software* 65 (11): 1–38. <https://www.jstatsoft.org/v65/i11/>.
- Fisher, Ronald. 1935. *The Design of Experiments*. Oliver; Boyd.
- Gini, Corrado. 1921. “Measurement of Inequality of Incomes.” *The Economic Journal* 31 (121): 124–26. <https://doi.org/10.2307/2223319>.
- Gramacy, Robert B., and Matthew Taddy. 2010. “Categorical Inputs, Sensitivity Analysis, Optimization and Importance Tempering with tgp Version 2, an R Package for Treed Gaussian Process Models.” *Journal of Statistical Software* 33 (6): 1–48. <https://doi.org/10.18637/jss.v033.i06>.
- Grömping, Ulrike. 2014. “R Package FrF2 for Creating and Analyzing Fractional Factorial 2-Level Designs.” *Journal of Statistical Software* 56 (1): 1–56. <https://www.jstatsoft.org/v56/i01/>.
- . 2018. “R Package Doe.base for Factorial Experiments.” *Journal of Statistical Software* 85 (5): 1–41. <https://doi.org/10.18637/jss.v085.i05>.
- . 2020. *DoE.wrapper: Wrapper Package for Design of Experiments Functionality*. <https://CRAN.R-project.org/package=DoE.wrapper>.
- Landau, William Michael. 2021. “The Targets r Package: A Dynamic Make-Like Function-Oriented Pipeline Toolkit for Reproducibility and High-Performance Computing.” *Journal of Open Source Software* 6 (57): 2959. <https://doi.org/10.21105/joss.02959>.
- Le Pennec, Erwan, and Kamil Slowikowski. 2019. *Ggwordcloud: A Word Cloud Geom for 'Ggplot2'*. <https://CRAN.R-project.org/package=ggwordcloud>.
- Lenth, Russell V. 2009. “Response-Surface Methods in R, Using rsm.” *Journal of Statistical Software* 32 (7): 1–17. <https://doi.org/10.18637/jss.v032.i07>.
- Lorenz, M O. 1905. “Methods of Measuring the Concentration of Wealth.” *Publications of the American Statistical Association* 9 (70): 209–19. <https://doi.org/10.2307/2276207>.
- McKay, M D, R J Beckman, and W J Conover. 1979. “Comparison of Three Methods for Selecting Values of Input Variables in the Analysis of Output from a Computer Code.” *Technometrics*:

- A Journal of Statistics for the Physical, Chemical, and Engineering Sciences* 21 (2): 239–45. <https://doi.org/10.1080/00401706.1979.10489755>.
- Pedersen, Thomas Lin. 2021. *Ggraph: An Implementation of Grammar of Graphics for Graphs and Networks*. <https://CRAN.R-project.org/package=ggraph>.
- Plackett, R L, and J P Burman. 1946. "The Design of Optimum Multifactorial Experiments." *Biometrika* 33 (4): 302–25.
- Pukelsheim, Friedrich. 2006. *Optimal Design of Experiments*.
- R Core Team. 2021. *R: A Language and Environment for Statistical Computing*. Vienna, Austria: R Foundation for Statistical Computing. <https://www.R-project.org/>.
- Rossum, G. van. 1995. "Python Tutorial." CS-R9526. Amsterdam: Centrum voor Wiskunde en Informatica (CWI).
- Rudis, Bob, and Blake Embrey. 2020. *Pluralize: Pluralize and 'Singularize' Any (English) Word*. <https://CRAN.R-project.org/package=pluralize>.
- Santiago, J, M Claeys-Bruno, and M Sergent. 2012. "Construction of Space-Filling Designs Using WSP Algorithm for High Dimensional Spaces." *Chemometrics and Intelligent Laboratory Systems* 113 (April): 26–31. <https://doi.org/10.1016/j.chemolab.2011.06.003>.
- SAS Institute. 1985. *SAS User's Guide: Statistics*. Vol. 2. Sas Inst.
- Shewry, M C, and H P Wynn. 1987. "Maximum Entropy Sampling." *Journal of Applied Statistics* 14 (2): 165–70. <https://doi.org/10.1080/02664768700000020>.
- Silge, Julia, and David Robinson. 2016. "Tidytext: Text Mining and Analysis Using Tidy Data Principles in r." *JOSS* 1 (3). <https://doi.org/10.21105/joss.00037>.
- Taguchi, Genichi. 1986. *Introduction to Quality Engineering: Designing Quality into Products and Processes*. Quality Resources.
- Tanaka, Emi. 2021. "edibble R-package." GitHub Repository. <https://github.com/emitanaka/edibble>; GitHub.
- Ushey, Kevin. 2022. *Renv: Project Environments*. <https://CRAN.R-project.org/package=renv>.
- Venables, Bill. 2013. *Conf.design: Construction of Factorial Designs*. <https://CRAN.R-project.org/package=conf.design>.
- Wasserman, Anthony I. 1975. "Issues in Programming Language Design - an Overview." In.
- Wheeler, Bob. 2022. *AlgDesign: Algorithmic Experimental Design*. <https://CRAN.R-project.org/package=AlgDesign>.
- Whittaker, D., E. R. Williams, and J. A John. 2022. *CycDesign: A Package for the Computer Generation of Experimental Designs*, (version 7.0). <https://vsni.co.uk/software/cycdesign>.
- Wickham, Hadley. 2016. *Ggplot2: Elegant Graphics for Data Analysis*. Springer-Verlag New York. <https://ggplot2.tidyverse.org>.
- Wickham, Hadley, Mara Averick, Jennifer Bryan, Winston Chang, Lucy D'agostino McGowan, Romain Francois, Garrett Grolemund, et al. 2019. "Welcome to the Tidyverse." *Journal of Open Source Software* 4 (43): 1686.
- Xie, Yihui. 2015. *Dynamic Documents with R and Knitr*. 2nd ed. Boca Raton, Florida: Chapman; Hall/CRC. <https://yihui.org/knitr/>.
- Xie, Yihui, J. J. Allaire, and Garrett Grolemund. 2018. *R Markdown: The Definitive Guide*. Boca Raton, Florida: Chapman; Hall/CRC. <https://bookdown.org/yihui/rmarkdown>.
- Zeileis, Achim. 2005. "CRAN Task Views." *R News* 5 (1): 39–40. <https://CRAN.R-project.org/doc/Rnews/>.
- . 2014. *Ineq: Measuring Inequality, Concentration, and Poverty*. <https://CRAN.R-project.org/package=ineq>.
- Zeileis, Achim, Jason C. Fisher, Kurt Hornik, Ross Ihaka, Claire D. McWhite, Paul Murrell, Reto Stauffer, and Claus O. Wilke. 2020. "colorspace: A Toolbox for Manipulating and Assessing Colors and Palettes." *Journal of Statistical Software* 96 (1): 1–49. <https://doi.org/10.18637/jss.v096.i01>.
- Zhu, Hao. 2021. *kableExtra: Construct Complex Table with 'Kable' and Pipe Syntax*. <https://CRAN.R-project.org/package=kableExtra>.

*Emi Tanaka*  
Monash University  
Department of Econometrics and Business Statistics  
Monash University, Clayton, VIC 3800  
<https://emitanaka.org/>  
ORCID: 0000-0002-1455-259X  
[emi.tanaka@monash.edu](mailto:emi.tanaka@monash.edu)

*Dewi Amaliah*  
Monash University  
Department of Econometrics and Business Statistics

*Monash University, Clayton, VIC 3800*

# BCClong: An R Package for Bayesian Consensus Clustering for Multiple Longitudinal Features

Zhiwen Tan, Chang Shen, Zihang Lu

**Abstract** It is very common nowadays for a study to collect multiple longitudinal features, and appropriately integrating these features simultaneously for defining individual subgroups (i.e., clusters) becomes increasingly crucial to understanding population heterogeneity and predicting future outcomes. The aim of this paper is to describe a new package, **BCClong**, which implements a Bayesian consensus clustering (BCC) model for multiple longitudinal features. Compared to existing packages, several key features make the **BCClong** package appealing: (a) it allows simultaneous clustering of mixed-type (e.g., continuous, discrete, and categorical) longitudinal features, (b) it allows each longitudinal feature to be collected from different sources, with measurements taken at distinct sets of time points (known as irregularly sampled longitudinal data), and (c) it relaxes the assumption that all features have the same clustering structure by estimating the feature-specific (local) clusterings and consensus (global) clustering. Using two real data examples, we provide a tutorial with step-by-step instructions on how to use the package.

## 1 Introduction

Cluster analysis has been widely used to identify subgroups of the study population, and it is also known as an unsupervised learning method in the machine learning literature. When longitudinal data is available, identifying distinct developmental patterns is a common goal in many research studies. Traditional statistical models often focus on clustering a single longitudinal feature, which may ignore the complexity of the clustering structures and the dynamic relationship when multiple longitudinal features are available. It is very common nowadays for a study to collect multiple features, and appropriately integrating multiple longitudinal features simultaneously for defining individual subgroups (i.e., clusters) becomes increasingly crucial to understanding population heterogeneity and predicting future outcomes. In particular, compared to clustering a single longitudinal feature, simultaneously clustering multiple longitudinal features (also known as joint clustering) facilitates the discovery of co-occurring trajectory patterns for multiple features of interest. Several R packages have been developed to cluster longitudinal data with multiple features. For model-based clustering, the **gbmt** package (Magrini, 2022a) and the **flexmix** package (Leisch, 2004) are developed to implement the group-based multi-trajectory analysis (GBMT) (Nagin et al., 2018; Magrini, 2022b). Both packages implement the EM algorithm, allowing unbalanced data and missing values, and can only be applied to continuous features. Notably, a SAS procedure PROC TRAJ is also developed for this model and has been widely used (Nagin et al., 2018). The GBMT assumes that the features and repeated measures are independent conditional on the cluster membership. Alternatively, the **lcmm** and **mixAK** packages are developed under the framework of mixed-effect models. The dependence between features and repeated measurements is captured by joint modeling of the random effects. Specifically, the **lcmm** package is designed for performing analyses such as latent class mixed models for a single longitudinal feature, joint latent class mixed models for multiple longitudinal features, as well as joint latent class mixed models for longitudinal features and time-to-event outcomes. While this package allows modeling multiple continuous or categorical longitudinal data of the same type, it does not support clustering mixed data types (e.g., continuous and categorical) simultaneously. On the other hand, the **mixAK** package can handle same-type (e.g., multiple continuous) or mixed-type (e.g., continuous, discrete, and categorical) features. For nonparametric clustering, the **kml3d** package (Genolini et al., 2013) is developed to implement the K-means clustering for multiple longitudinal features, which is an extension of the **kml** package (Genolini and Falissard, 2009) for a single longitudinal feature. However, this package can only be applied to continuous features and does not allow for missing data. See Lu et al. (2023) for a review and comparison of these packages. The aim of this paper is to describe a new package, the **BCClong** package, which implements the BCC model for multiple longitudinal features (Lu and Lou, 2022b; Tan et al., 2022). Compared to existing packages, several key features make the **BCClong** package appealing: (a) it allows simultaneous clustering of mixed-type (e.g., continuous, discrete, and categorical) longitudinal features, (b) it allows each longitudinal feature to be collected from different sources, with measurements taken at distinct sets of time points (known as irregularly sampled longitudinal data), and (c) it relaxes the assumption that all features have the same clustering structure by estimating the feature-specific (local) clusterings and overall (global) clustering. This model is very flexible, which facilitates the interpretation of clustering results and

enhances its practical utilities.

## 2 Models and software

The methodology of the BCC model for clustering multiple longitudinal features implemented in the **BCClong** package is described in detail in [Lu and Lou \(2022b\)](#) (for multiple continuous features) and in [Tan et al. \(2022\)](#) (for multiple mixed-type features). In this section, we provide a brief overview of this model and the **BCClong** package. Without loss of generality, we use the term "features" to represent variables of interest used for the cluster analysis; however, these variables can also be a patient's traits, biomarkers, risk factors, exposures, and outcomes of interest in different contexts and studies.

### 2.1 Bayesian Consensus Clustering for Multiple Longitudinal Features

The BCC model is originally proposed as a flexible approach to model the dependence and heterogeneity of the data collected from multiple sources ([Lock and Dunson, 2013](#)). Instead of arriving at a single overall clustering structure, the BCC model allows each feature to follow feature-specific (local) clustering and these local clusterings are aggregated to find a consensus (global) clustering. Let  $L_{i,r} = 1, \dots, K$  denote the feature-specific (local) cluster label for individual  $i$  and feature  $r$  for  $r = 1, \dots, R$ , and  $C_i = 1, \dots, K$  denote the consensus (global) cluster label for individual  $i$ , where  $K$  denotes the number of clusters and  $R$  denotes the number of longitudinal features. Let  $\mathbf{y}_{i,r} = (y_{i1,r}, \dots, y_{in_{i,r},r})^\top$  denote the measurement for individual  $i$  and feature  $r$ , where  $y_{ij,r}$  denotes the observation  $j$  of individual  $i$  for feature  $r$ , and  $n_{i,r}$  is the number of measurements, for  $i = 1, \dots, N, j = 1, \dots, n_{i,r}, r = 1, \dots, R$ . The BCC model assumes that  $\mathbf{L}_r = (L_{1,r}, \dots, L_{N,r})$  is dependent on  $\mathbf{C} = (C_1, \dots, C_N)$  through  $P(L_{i,r} = k | C_i) = \vartheta(k, C_i, \alpha_r)$ , where  $\alpha_r$  adjusts the dependence function  $\vartheta(\cdot)$ . The conditional model can be specified as

$$P(L_{i,r} = k | \mathbf{y}_{i,r}, C_i, \gamma_{k,r}, \beta_{ik}) \propto \vartheta(k, C_i, \alpha_r) f_{k,r}(\mathbf{y}_{i,r} | \gamma_{k,r}, \beta_{ik,r}) \quad (1)$$

where  $\gamma_{k,r}$  and  $\beta_{ik,r}$  denote the fixed effects and random effects for individual  $i$  belonging to cluster  $k$ , where  $k = 1, \dots, K$ . To specify  $\vartheta(k, C_i, \alpha_r)$ , the package uses the following function ([Lock and Dunson, 2013](#)).

$$\vartheta(k, C_i, \alpha_r) = \begin{cases} \alpha_r & \text{if } C_i = L_{i,r} \\ (1 - \alpha_r)/(K - 1) & \text{otherwise} \end{cases} \quad (2)$$

Therefore,  $\alpha_r$  (for  $r = 1, \dots, R$ ) can also be viewed as an adherence parameter representing the degree of agreement between  $L_{i,r}$  and  $C_i$ . Moreover,  $f_{k,r}(\mathbf{y}_{i,r} | \gamma_{k,r}, \beta_{ik,r})$  is a distribution from the exponential family with the dispersion parameter  $\phi_{k,r}$  and the fully specified mean function is given by

$$h_{k,r}^{-1}(E(\mathbf{y}_{i,r} | \beta_{ik,r}, \gamma_{k,r})) = \eta_{ik,r} = \mathbf{x}_{i,r}^\top \gamma_{k,r} + \mathbf{Z}_{i,r}^\top \beta_{ik,r} \quad (3)$$

where  $h_{k,r}^{-1}$  is a canonical link function for the mean of the feature  $r$  in cluster  $k$ . Also,  $\eta_{ik,r} = (\eta_{i1,r}, \dots, \eta_{in_{i,r},r})^\top$  is the linear predictor for feature  $r$ ,  $\mathbf{x}_{i,r}$  is a  $p_r \times n_{i,r}$  vector of predictors,  $\gamma_{k,r}$  is the corresponding vector of coefficients,  $\mathbf{Z}_{i,r}$  is a  $q_r \times n_{i,r}$  vector of predictors, and  $\beta_{ik,r} | \cdot \sim \text{MVN}(\mathbf{0}, \Sigma_{k,r})$ . In addition, the distribution of  $\mathbf{y}_{i,r}$  can be written as

$$P(\mathbf{y}_{i,r} | \beta_{ik,r}, \gamma_{k,r}, \phi_{k,r}) = \prod_{j=1}^{n_{i,r}} \exp \left\{ \frac{y_{ij,r} \eta_{ijk,r} - g_{k,r}(\eta_{ijk,r})}{\phi_{k,r}} + w_{k,r}(y_{ij,r}, \phi_{k,r}) \right\} \quad (4)$$

where  $g(\eta)$  and  $w(y, \phi)$  are appropriate distribution-specific functions. For example, for features with Gaussian distributions,  $g(\eta) = \eta^2/2$  and  $w(y, \phi) = \frac{\log(2\pi\phi)}{2} - \frac{y^2}{2\phi}$ . Furthermore, for three features with Gaussian, Poisson, and Binomial distributions, respectively, the model can be specified as:

- (1):  $E(\mathbf{y}_{i,1} | \gamma_{k,1}, \beta_{ik,1}) = \eta_{ik,1} = \mathbf{x}_{i,1}^\top \gamma_{k,1} + \mathbf{Z}_{i,1}^\top \beta_{ik,1}$
- (2):  $\log(E(\mathbf{y}_{i,2} | \gamma_{k,2}, \beta_{ik,2})) = \eta_{ik,2} = \mathbf{x}_{i,2}^\top \gamma_{k,2} + \mathbf{Z}_{i,2}^\top \beta_{ik,2}$
- (3):  $\text{logit}(E(\mathbf{y}_{i,3} | \gamma_{k,3}, \beta_{ik,3})) = \eta_{ik,3} = \mathbf{x}_{i,3}^\top \gamma_{k,3} + \mathbf{Z}_{i,3}^\top \beta_{ik,3}$

where  $i = 1, \dots, N, j = 1, \dots, n_{i,r}$  and  $r = 1, 2, 3$ . The model also involves dispersion parameters  $\phi_{k,r} = \sigma_{k,r}^2$  for Gaussian distribution. The corresponding dispersion parameters for logistic regression and Poisson regression are both 1, for  $r = 1, \dots, R$  and  $k = 1, \dots, K$ . Given this proposed model, the

complete-data likelihood can be written as

$$\mathbb{L}^{Bayes}(\Theta | \mathbf{y}_i, \boldsymbol{\beta}_i; \mathbf{x}_i) = \sum_{k_1, \dots, k_R} \left( \sum_k \prod_{r=1}^R \vartheta(L_{i,r} = k_r, C_i = k, \alpha_r) \pi_k \right) \prod_{r=1}^R f(\mathbf{y}_{i,r} | \gamma_{ik_r, r}, \boldsymbol{\beta}_{ik_r, r}; \mathbf{x}_{i,r}) \quad (5)$$

where  $\mathbf{y}_i = (\mathbf{y}_{i,1}, \dots, \mathbf{y}_{i,R})$ ,  $\mathbf{x}_i = (\mathbf{x}_{i,1}, \dots, \mathbf{x}_{i,R})$ ,  $\boldsymbol{\beta}_i = (\boldsymbol{\beta}_{i,1}, \dots, \boldsymbol{\beta}_{i,R})$ ,  $\pi_k = P(C_i = k)$  and  $\Theta = (\boldsymbol{\alpha}, \boldsymbol{\pi}, \boldsymbol{\gamma}, \boldsymbol{\Sigma}, \boldsymbol{\phi})$ .

## 2.2 Bayesian Inference

For computational convenience, the **BCClong** package uses independent conjugate priors for the model parameters. Specifically, for the adherence parameters, the prior distribution is  $\alpha_r \sim \text{TBeta}(\delta_{1,r}, \delta_{2,r}, 1/K)$  for  $r = 1, \dots, R$ , where TBeta denotes a truncated Beta distribution ranged from  $[1/K, 1]$ . For the cluster probabilities, the prior distribution is  $\boldsymbol{\pi} \sim \text{Dirichlet}(\boldsymbol{\varphi}_0)$ . For the fixed effect coefficients, the prior distribution is  $\boldsymbol{\gamma}_{k,r} \sim \text{MVN}(\mathbf{0}, \mathbf{V}_{0k,r})$ , where  $\mathbf{V}_{0k,r}$  is an  $m_{k,r} \times m_{k,r}$  variance-covariance matrix, and  $m_{k,r}$  is the dimension of  $\boldsymbol{\gamma}_{k,r}$ . The dispersion parameter  $\phi_{k,r}$  is not constant only when the feature is a Gaussian distribution. In such a case,  $\phi_{k,r} = \sigma_{k,r}^2$ . The prior distribution is  $\sigma_{k,r}^2 \sim \text{IG}(a_{0k,r}, b_{0k,r})$ , where  $a_{0k,r}$  and  $b_{0k,r}$  are the parameters of an inverse gamma distribution. By definition,  $\phi_{k,r} = 1$  for features that follow a Poisson or Binomial distribution, for  $k = 1, \dots, K$  and  $r = 1, \dots, R$ . For the variance-covariance matrix of the random effect, the prior distribution is  $\boldsymbol{\Sigma}_{k,r}^{-1} \sim \text{Wishart}(\lambda_{0k,r}, (\lambda_{0k,r} \boldsymbol{\Lambda}_{0k,r})^{-1})$ , where the prior for the Wishart distribution is parametrized such that the mean is  $\boldsymbol{\Lambda}_{0k,r}^{-1}$ . In the special case when  $\boldsymbol{\Sigma}_{k,r}$  is a diagonal matrix, that is,  $\boldsymbol{\Sigma}_{k,r} = \text{diag}(\xi_{k1,r}^2, \dots, \xi_{kq_r,r}^2)$ , the prior is an inverse gamma distribution, i.e.,  $\xi_{km,r}^2 \sim \text{IG}(c_{0k,r}, d_{0k,r})$ , for  $m = 1, \dots, q_r$ . For posterior computation, a Gibbs sampling approach to update the model parameters is described as follows. After initializing the parameters, the algorithm repeats the following steps until convergence. At the  $s$  step of the iteration,

- Update local cluster membership  $L_{i,r}^{(s)}$  given  $\{\mathbf{y}_{i,r}, \boldsymbol{\Theta}_r^{(s-1)}, \alpha_r^{(s-1)}, C_i^{(s-1)}, \boldsymbol{\beta}_r^{(s-1)}\}$ , for  $i = 1, \dots, N$  and  $r = 1, \dots, R$ .
- Update  $\alpha_r^{(s)}$  given  $\{\mathbf{C}^{(s-1)}, \mathbf{L}^{(s)}\}$ , for  $r = 1, \dots, R$ .
- Update  $C_i^{(s)}$  given  $\{\mathbf{L}^{(s)}, \boldsymbol{\alpha}^{(s)}, \boldsymbol{\pi}^{(s-1)}\}$ , for  $i = 1, \dots, N$ .
- Update  $\boldsymbol{\pi}^{(s)}$  given  $\mathbf{C}^{(s)}$ .
- Update cluster-specific parameters  $\boldsymbol{\Theta}^{(s)} = (\boldsymbol{\gamma}^{(s)}, \boldsymbol{\Sigma}^{(s)}, \sigma^{2(s)})$  and  $\boldsymbol{\beta}^{(s)}$ . In particular,  $\boldsymbol{\gamma}^{(s)}$  and  $\boldsymbol{\beta}^{(s)}$  are updated via the Metropolis-Hastings algorithm.

Each MCMC iteration will produce a realization of cluster membership for  $C, L_1, \dots, L_R$ . The package uses the mode over all the MCMC samples (after burn-in and thinning) as the point estimates for both feature-specific and global clustering. Label switching is a common phenomenon in mixture models. This problem arises due to both the likelihood and posterior distributions being invariant to permutations of the parameters. The **BCClong** package applies a post-processing algorithm, via the **label.switching** package (Papastamoulis, 2016), to reorder the labels based on Kullback-Leibler divergence (Stephens, 2000). The computation of the BCC model can be carried out using the main function in the **BCClong** package:

```
BCC.multi(mydat, dist, id, time, formula, num.cluster,
          hyper.par, initials, initial.cluster.membership,
          input.initial.local.cluster.membership,
          input.initial.global.cluster.membership,
          burn.in, thin, per, max.iter, ...)
```

The meaning of the key arguments of the **BCC.multi()** function is described as follows.

- **mydat**: list of  $R$  longitudinal features (i.e., with a length of  $R$ ), where  $R$  is the number of features. The data should be prepared in a long format (each row is one time point per individual).
- **dist**: a character vector (with a length of  $R$ ) that determines the distribution for each feature. Possible values are "gaussian" for a continuous feature, "poisson" for a discrete feature (e.g., count data) using a log link, and "binomial" for a dichotomous feature (0/1) using a logit link. A single value (i.e., a length of 1) is recycled if necessary.
- **id**: a list (with a length of  $R$ ) of vectors of the study ID of individuals for each feature. A single value (i.e., a length of 1) is recycled if necessary.
- **time**: a list (with a length of  $R$ ) of vectors of time (or age) at which the feature measurements are recorded.

- **formula**: a list (with a length of  $R$ ) of formulas for each feature. Each formula is a two-sided linear formula object describing both the fixed-effects and random effects part of the model, with the response (i.e., longitudinal feature) on the left of a `~` operator and the terms, separated by `+` operations, to the right. Random-effects terms are distinguished by vertical bars (`|`) separating expressions for design matrices from grouping factors. See the `formula` argument from the **lme4** package.
- **num.cluster**: number of clusters  $K$ . A minimum value of  $K$  is 2.
- **hyper.par**: hyper-parameters of the prior distributions for the model parameters. The default hyper-parameters values will result in weakly informative prior distributions.
- **initials**: a list of initial values for model parameters.
- **initial.cluster.membership**: a character, which can be "random", "mixAK" or "input". For "random" (default), the local cluster membership is randomly generated from 1 to  $K$  for each feature. For "mixAK", the local cluster membership is obtained by fitting the Bayesian mixture model for longitudinal data to each feature through the **mixAK** package for the given  $K$ . For both methods, the global cluster membership is set to be equal to the local cluster membership of the first feature. For "input", the user can input initial values for local and global cluster memberships, using the argument `input.initial.local.cluster.membership` and `input.initial.global.cluster.membership`.
- **input.initial.local.cluster.membership**: if `initial.cluster.membership = "input"`, then this argument must not be empty. It is a list with a length of  $R$  and each element of the list corresponds to the local clustering of each feature. See the Illustration section for an example.
- **input.initial.global.cluster.membership**: if `initial.cluster.membership = "input"`, by default the `input.initial.global.cluster.membership` is the local cluster membership of the first feature, and therefore this argument need not be supplied, but the user can also use this argument to supply a different initial cluster membership for the global clustering. See the Illustration section for an example.
- **burn.in**: the number of samples discarded. This value must be smaller than `max.iter`.
- **thin**: the number of thinning. For example, if `thin = 10`, then the MCMC chain will keep one sample every 10 iterations.
- **per**: specify how often the MCMC chain will print the iteration number.
- **max.iter**: the number of MCMC iterations.

In the **BCClong** package, the default hyper-parameter values for the prior distributions are set to reflect no prior information regarding the values of these parameters, that is, uninformative priors are used. The argument `hyper.par` allows setting these hyper-parameter values. The default values are:

```
hyper.par = list(a.star = 1, b.star = 1, delta = 1,
                 aa0 = 0.001, bb0 = 0.001, cc0 = 0.001, dd0 = 0.001,
                 vv0 = 1000)
```

This indicates that for  $\delta_{1,r} = \delta_{2,r} = 1$  for  $r = 1, \dots, R$  (corresponding to `a.star = 1, b.star = 1`). The elements of  $\varphi$  are set to 1 (corresponding to `delta = 1`), and  $a_{0k,r} = b_{0k,r} = 0.001$  for  $k = 1, \dots, K$  and  $r = 1, \dots, R$  (corresponding to `aa0 = 0.001, bb0 = 0.001`). Also,  $V_{0k,r}$  is set to be a diagonal matrix with elements of 1000 (corresponding to `vv0 = 1000`). Finally,  $\Sigma_{k,r}$  is set to be a diagonal matrix with elements following an inverse gamma distribution (corresponding to `cc0 = 0.001, dd0 = 0.001`). Customized or informative priors can be used by modifying the values in `hyper.par`. By default, the `initials` is set to `NULL` and the **BCClong** will automatically generate initial values for both the local and global cluster memberships and the model parameters. Convergence of MCMC is diagnosed using visual inspection of trace plots as well as using the Geweke statistics (Geweke, 1991). The **BCClong** package provides a function called `traceplot()` to produce trace plots for model parameters of interest. The `traceplot()` function has the following arguments:

```
traceplot(fit, cluster.indx, feature.indx, parameter)
```

The meaning of these arguments is described as follows.

- **fit**: an objective output from `BCC.multi()` function.
- **cluster.indx**: a numeric value. For cluster-specific parameters, specifying `cluster.indx` will generate the trace plot for the corresponding cluster.
- **feature.indx**: a numeric value. For cluster-specific parameters, specifying `feature.indx` will generate the trace plot for the corresponding cluster.
- **parameter**: a character value. Specify which parameter for which the trace plot will be generated. The value can be "PPI" for  $\pi$ ,  $\alpha$  for  $\alpha$ , "GA" for  $\gamma$ , "SIGMA.SQ.U" for  $\Sigma$  and "SIGMA.SQ.E" for  $\sigma$ .

In addition, the output of summary statistics from the **BCC.multi()** function also includes the Geweke statistics calculated using the **coda** package (Plummer et al., 2006). Briefly, the Geweke statistics determine whether a Markov chain converges or not based on a test for equality of the means of the first and last parts of a Markov chain (by default the first 10% and the last 50%). If the samples are drawn from the stationary distribution of the chain, the two means are equal and Geweke's statistic has an asymptotically standard normal distribution. Practically, if the estimated Geweke statistics fall between -1.96 and 1.96, we consider the model parameter has converged.

### 2.3 Missing Data, Measurement Time and Frequency

The proposed BCC model for longitudinal data and the **BCClong** package allow the modeling of missing data under the assumption that the data are missing at random (Lu and Lou, 2022b; Tan et al., 2022). No extra steps are required from the users to handle the missing data. The package will keep individuals with at least one observation for any input features and use them for the analysis, whereas the package will remove individuals with no data on all features. The package also allows the number of observations to be different between individuals and features, and these observations can also be measured at time points that are different between individuals and features.

### 2.4 Determining the Number of Clusters

Determining the number of clusters in model-based clustering is challenging and there are no widely accepted approaches. A review of different approaches within the Bayesian framework can be found in Celeux et al. (2006); Nasserinejad et al. (2017); Merkle et al. (2019). The **BCClong** package offers three criteria to determine the number of clusters, namely the mean adjusted adherence (Lock and Dunson, 2013), the Deviance Information Criterion (DIC) (Spiegelhalter et al., 2002; Celeux et al., 2006), and the widely Applicable Information Criterion (WAIC) (Watanabe, 2010; Gelman et al., 2014). We briefly describe these three approaches in this subsection. For consensus clustering, Lock and Dunson (2013) proposed an empirical separability criterion, which selects the value of  $K$  that gives maximum adherence to an overall clustering. For each  $K \geq 2$ , the estimated adherence parameters  $\alpha_r \in [\frac{1}{K}, 1]$  are mapped to the unit interval by the linear transformation  $\alpha_r^* = \frac{K\alpha_r - 1}{K-1}$ , thus  $\alpha_r^* \in [0, 1]$ . One then selects the value of  $K$  that results in the highest mean adjusted adherence, which is defined as  $\bar{\alpha}^* = \frac{1}{K} \sum_{m=1}^R \alpha_m^*$ . This approach selects a model with which the feature-specific clusterings contribute the most information (on average) to the global clustering. A model with a higher  $\bar{\alpha}^*$  is preferred. DIC and its variants have also been widely used for clustering multivariate longitudinal data. See references (Lu and Lou, 2019; Neelon et al., 2011; Leiby et al., 2009; Elliott et al., 2005; Frühwirth-Schnatter and Pyne, 2010) for example. In the **BCClong** package, the DIC is computed based on the complete likelihood (also known as DIC<sub>4</sub> in Celeux et al. (2006)), which is defined as  $\text{DIC} = -4E_{\Theta, C, L}[\log f(\mathbf{y}, C, L|\Theta)|\mathbf{y}] + 2E_{C, L}[\log f(\mathbf{y}, C, L|E[\Theta|\mathbf{y}, C, L])|\mathbf{y}]$ , where  $\mathbf{y} = (y_1, \dots, y_N)$  and  $\Theta$  is defined in equation (5). The effective number of parameters and the posterior mean of deviance are defined as  $p_D = -E_{\Theta, C, L}[2 \log f(\mathbf{y}, C, L|\Theta)|\mathbf{y}] + E_{C, L}[2 \log f(\mathbf{y}, C, L|E[\Theta|\mathbf{y}, C, L])|\mathbf{y}]$  and  $\overline{D(\Theta)} = -E_{\Theta, C, L}[2 \log f(\mathbf{y}, C, L|\Theta)|\mathbf{y}]$ , respectively. A model with a lower DIC is preferred. WAIC is an extension of the Akaike Information Criterion (AIC) that is more fully Bayesian than the DIC. Similar to DIC, WAIC estimates the effective number of parameters to adjust for overfitting. The **BCClong** package uses the **LaplaceDemon** package (Hall, 2021) to calculate the WAIC, which is defined as  $\text{WAIC} = -2(lppd - p_{\text{WAIC}})$ , where lppd is the log pointwise predictive density, which is equivalent to  $lppd = \log \prod_{i=1}^N p_{\text{post}}(y_i | C, L)$ . Also,  $p_{\text{WAIC}} = \sum_{i=1}^N \text{var}_{\text{post}}(\log p(y_i | C, L))$  is an approximation to the number of unconstrained and uninformed parameters, where a parameter counts as 1 when estimated without constraint or any prior information, 0 if fully constrained or all information comes from the prior distribution, or an intermediate number if both the data and prior are informative. The  $p_{\text{WAIC}}$  defined here corresponds to  $p_{\text{WAIC}2}$  in Gelman et al. (2014), which is recommended because its results are closer in practice to the results of leave-one-out cross-validation. A model with a lower WAIC is preferred. In the **BCClong** package, the **model.selection.criteria()** function is used to compute the DIC and WAIC:

```
model.selection.criteria(fit, fast_version)
```

- **fit**: an objective output from **BCC.multi()** function.
- **fast\_version**: if **fast\_verion** = TRUE (default), then compute the DIC and WAIC using the first 100 MCMC samples (after burn-in and thinning). If **fast\_version** = FALSE, then compute the DIC and WAIC using all MCMC samples (after burn-in and thinning).

## 2.5 Goodness of Fit

To assess the model goodness of fit, one can compute the posterior predictive check (Gelman et al., 1996). This approach is performed by comparing the observed data with data replicated from the posterior predictive distribution. If the model fits the data well, the replicated data, denoted as  $\mathbf{y}^{rep}$ , should have a similar distribution as the observed data  $\mathbf{y}$ . Let  $T = T(\mathbf{y}, \Phi)$  denote a discrepancy measure, where  $\Phi = (\Theta, \beta, L, C)$  denotes all model parameters (including random effects and cluster labels) and  $T$  is sample quantiles or residual-based measures. Following Gelman et al. (1996), a natural discrepancy measure is a  $\chi^2$  measure. For the BCC model, it is defined as

$$T^{obs}(\mathbf{y}, \Phi) = \sum_{k=1}^K \sum_{i=1}^N \sum_{j=1}^{n_{i,r}} \sum_{r=1}^R \frac{z_{ik,r} ||\mathbf{y}_{i,r} - h_{k,r}^{-1}(E(\mathbf{y}_{i,r} | \boldsymbol{\beta}_{ik,r}, \gamma_{k,r}))||^2}{\sigma_{k,r}^2} \quad (6)$$

where  $z_{ik,r} = 1$  if  $L_{i,r} = k$ , and 0 otherwise. This measure is computed at each MCMC step by treating cluster indicator  $z_{ik,r}$  and random effects  $\beta_{ik,r}$  as observed data. The Bayesian predictive p-value is the probability that the discrepancy measure based on a predictive sample,  $T^{rep}(\mathbf{y}^{rep}, \Phi)$  is more extreme than the observed measure  $T^{obs}(\mathbf{y}, \Phi)$ . This quantity is estimated by computing the proportion of draws in which  $T^{rep} > T^{obs}$ . A p-value close to 0.5 indicates the model provides a good fit to the data, whereas a p-value close to 0 or 1 indicates a poor fit. In the **BCClong** package, **BayesT0** is used for computing the posterior predictive check.

`BayesT(fit)`

- `fit`: an objective output from **BCC.multi()** function.

## 2.6 Visualization of Trajectory Patterns by Clusters

To plot the longitudinal trajectory of features by local and global clusterings derived from the BCC model, the **trajplot()** function can be used. This function uses the **ggplot2** package internally.

`trajplot(fit, feature.ind, which.cluster)`

- `fit`: an objective output from the **BCC.multi()** function.
- `feature.ind`: a numeric value indicating which feature to plot. The number indicates the order of the feature specified in the `mydat` argument of the **BCC.multi()** function.
- `which.cluster`: a character value: "global" or "local", indicating whether to plot the trajectory by global cluster or local cluster indices.

## 3 Illustrations

In this section, we demonstrate the utility of the **BCClong** package (version 1.0.3) using two real data examples. The first example was quality of life data following epilepsy drug treatment, which included three continuous features. This dataset can be found in the **joineRML** package (Hickey et al., 2018). The second example was the Mayo Clinic primary biliary cholangitis (PBC) data, which included three features of mixed types (i.e., continuous, discrete, and categorical). This dataset can be found in **mixAK** package (Komárek and Komárková, 2014).

### 3.1 Example 1: Quality of Life Data Following Epilepsy Drug Treatment

The data called **epileptic.qol** came from the Standard and New Antiepileptic Drugs (SAND) study (Marson et al., 2007). This is a randomized control trial of standard and new antiepileptic drugs, comparing effects on long-term clinical outcomes. Quality of life (QoL) data were collected by mail at baseline, 3 months, and at 1 and 2 years using validated measures. The first measurement for each individual was the baseline measurement; however, there was variability in the time taken to return the questionnaires. Similarly, the second, third, and fourth follow-up times, which were scheduled for 3 months, 1-year, and 2 years, respectively, also had variability in completion times. The data can be loaded by

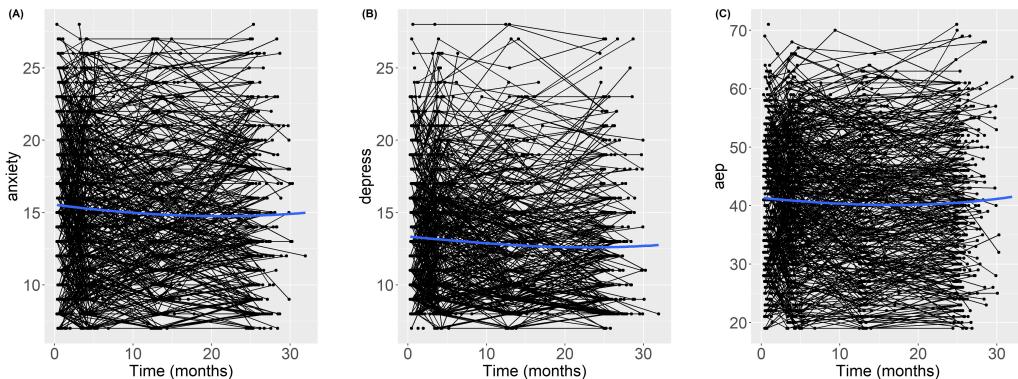
```
R> library("joineRML")
R> head(epileptic.qol[,c(1,5,6,7,8)])
   id time anxiety depress aep
1  1  147      11      14    43
```

2	1	259	12	12	51
3	1	519	20	21	63
4	1	906	17	20	53
5	2	134	19	13	45
6	2	258	21	16	50

This is a longitudinal dataset with one row per visit. At each visit, measurements of three continuous features ( $R = 3$ ) are recorded, namely the anxiety score (anxiety), depression score (depress), and Liverpool Adverse Events Profile (aep) defined according to the NEWQOL (Newly Diagnosed Epilepsy Quality of Life) assessment. One could use a spaghetti plot to visualize the trajectories of these features (Figure 1). The goal of the analysis was to identify distinct longitudinal QoL patterns based on these three features simultaneously. Based on the observed trajectories, we used a linear form with a random intercept for each feature with a Gaussian distribution. The model is written as:

$$\begin{aligned} \text{anxiety: } E(\mathbf{y}_{i,1} | \gamma_{k,1}, \beta_{ik,1}) &= \eta_{ik,1} = \mathbf{x}_{i,1}^\top \gamma_{k,1} + \mathbf{Z}_{i,1}^\top \beta_{ik,1} \\ \text{depress: } E(\mathbf{y}_{i,2} | \gamma_{k,2}, \beta_{ik,2}) &= \eta_{ik,2} = \mathbf{x}_{i,2}^\top \gamma_{k,2} + \mathbf{Z}_{i,2}^\top \beta_{ik,2} \\ \text{aep: } E(\mathbf{y}_{i,3} | \gamma_{k,3}, \beta_{ik,3}) &= \eta_{ik,3} = \mathbf{x}_{i,3}^\top \gamma_{k,3} + \mathbf{Z}_{i,3}^\top \beta_{ik,3} \end{aligned}$$

where  $\mathbf{x}_{i,r} = (1, t_{i,r})^\top$  for  $r = 1, 2, 3$  and  $t_{i,r} = (t_{i,1,r}, \dots, t_{i,n_{i,r},r})^\top$  is a vector of time values at which the feature  $r$  is recorded for individual  $i$ . The corresponding coefficients are  $\gamma_{k,r} = (\gamma_{k,r1}, \gamma_{k,r2})^\top$ . Also,  $\mathbf{Z}_{i,r} = \mathbf{1}_{n_{i,r}}$  for  $r = 1, 2, 3$ . For model-based clustering based on a finite mixture model, it is necessary



**Figure 1:** Longitudinal trajectories of three features for the epileptic.qol dataset. A locally weighted scatterplot smoothing curve is overlaid on each panel to provide an estimate of the overall trend. (A) longitudinal trajectories of Anxiety score (anxiety), (B) longitudinal trajectories of depression score (depress), (C) longitudinal trajectories of Liverpool Adverse Events Profile (aep).

to normalize the data within each feature in order to model the growth pattern, not the level. This is because the level could dominate the sample variability so that the clusters will be mainly determined by the level and will not necessarily be homogeneous in trend (Heggeseth and Jewell, 2018). Therefore, in the current analysis, the three features were normalized using the `scale` function before entering the model, using the following commands,

```
R> epileptic.qol$anxiety_scale <- scale(epileptic.qol$anxiety)
R> epileptic.qol$depress_scale <- scale(epileptic.qol$depress)
R> epileptic.qol$aep_scale <- scale(epileptic.qol$aep)
```

### Determining the Number of Clusters

The following codes fit the BCC model with the number of clusters ranging from 2 to 5, i.e.,  $K = 2, \dots, 5$ . For each value of  $K$ , we saved the mean adjusted adherence (`alpha.adjust`). For each model, we generated 2000 samples and discarded the first 1000 samples. Note that the following part of the codes may be computationally intensive; if users just want to get familiar with the codes and outputs, a smaller number of iterations and burn-in can be used to reduce the computational time, for example, by setting `burn.in = 10`, `per = 1`, and `max.iter = 20`.

```
R> dat <- epileptic.qol
R> alpha.adjust <- NULL
R> for (k in 2:5){
+   fit.BCC <- BCC.multi (
```

```

+   mydat = list(dat$anxiety_scale, dat$depress_scale, dat$aep_scale),
+   dist = c("gaussian"),
+   id = list(dat$id),
+   time = list(dat$time),
+   formula = list(y ~ time + (1|id)),
+   num.cluster = k,
+   burn.in = 1000,
+   thin = 1,
+   per = 100,
+   max.iter = 2000)
+   alpha.adjust <- c(alpha.adjust, fit.BCC$alpha.adjust)
+
}

```

An alternative specification for the `formula` argument is

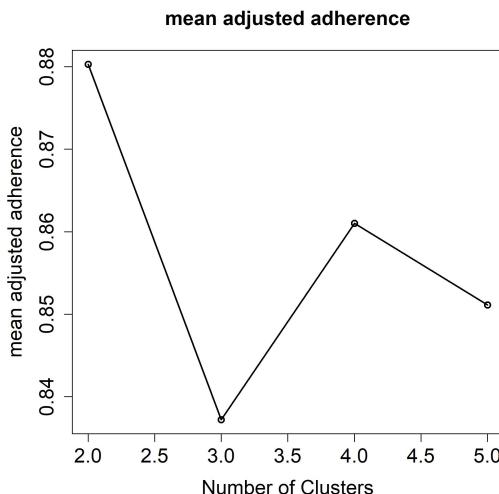
```
formula = list(y ~ time + (1|id), y ~ time + (1|id), y ~ time + (1|id))
```

This specification can be used if the users would like to use different functional forms for different features. Each element within the list corresponds to a model for a feature listed in the `mydat` argument. To compare the model under different numbers of clusters, one can plot the mean adjusted adherence index using the following commands:

```

R> num.cluster <- 2:5
R> plot(num.cluster, alpha.adjust, type = "o",
+       cex.lab = 1.5, cex.axis = 1.5, cex.main = 1.5, lwd = 2,
+       xlab = "Number of Clusters",
+       ylab = "mean adjusted adherence", main="mean adjusted adherence")

```



**Figure 2:** Mean adjusted adherence for  $K = 2$  to  $5$  for the epileptic.qol dataset.

From Figure 2, it was clear that the model with the number of clusters  $K = 2$  had the largest mean adjusted adherence (0.88), suggesting that a two-cluster model provided the best fit for the data. To determine the number of clusters based on DIC and WAIC, one can also use the `model.selection.criteria()` function to compute these indices. Next, we used the following commands to fit the final model with  $K = 2$ .

```

R> fit.BCC2 <- BCC.multi (
+   mydat = list(dat$anxiety_scale, dat$depress_scale, dat$aep_scale),
+   dist = c("gaussian"),
+   id = list(dat$id),
+   time = list(dat$time),
+   formula = list(y ~ time + (1|id)),
+   num.cluster = 2,
+   burn.in = 1000,
+   thin = 1,
+   per = 100,
+   max.iter = 2000)

```

One can apply `print()` and `summary()` functions to the object returned from `BCC.multi()` to obtain summary information of the result, for example,

```
> print(fit.BCC2)
> summary(fit.BCC2)
```

Both functions will return result summaries such as the number of individuals, number of features included in the analysis, and tabulation for global and local clusterings.

### Posterior Summary Statistics and Trace Plots

One can print the summary statistics for all the model parameters using the following commands (results not shown):

```
print(fit.BCC2$summary.stat)
```

As an illustration, we print the estimated cluster probabilities  $\pi = (\pi_1, \pi_2)$  as follows:

```
R> print(fit.BCC2$summary.stat$PPI)
```

	[,1]	[,2]
mean	0.43608373	0.56391627
sd	0.02367214	0.02367214
2.5%	0.39326189	0.51584563
97.5%	0.48415437	0.60673811
geweke.stat	-0.78671117	0.78671117

The summary statistics provided the posterior mean, standard deviation (sd), 2.5%tile, 97.5%tile, and the Geweke statistics. For example, the posterior mean for  $\pi_1$  and  $\pi_2$  were 0.44 and 0.56, respectively. This suggested that the estimated cluster proportions for Clusters 1 and 2 are 44% and 56%, respectively. The Geweke statistics for  $\pi_1$  and  $\pi_2$  fell between -2 and 2, suggesting that the two parameters converge to a stationary distribution. One can also use the following commands to print the adherence parameters for the three features, i.e.,  $\alpha = (\alpha_1, \alpha_2, \alpha_3)$ :

```
R> print(fit.BCC2$summary.stat$ALPHA)
```

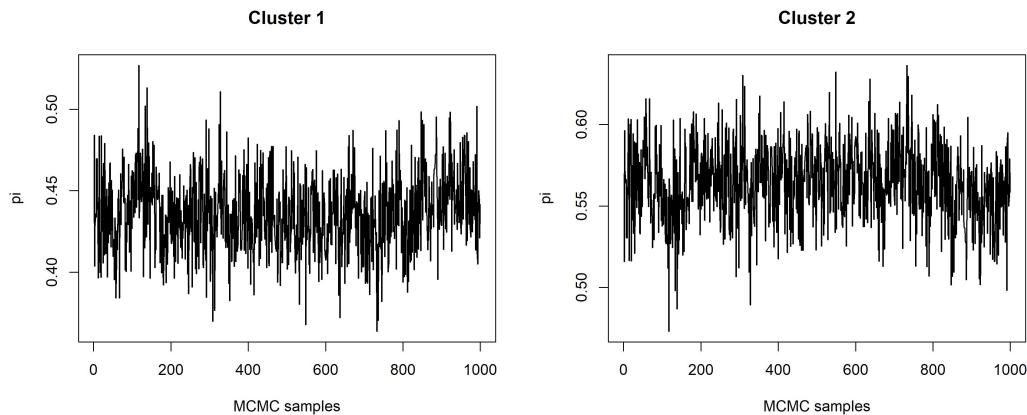
	[,1]	[,2]	[,3]
mean	0.97572877	0.91727462	0.92427049
sd	0.01230074	0.01764104	0.01771331
2.5%	0.94682192	0.88115409	0.88730582
97.5%	0.99620462	0.95091755	0.95581869
geweke.stat	0.10674171	0.05585855	-1.13303537

The output indicated that the posterior means for the adherence parameters were 0.98, 0.92, and 0.92, respectively. This suggested that all three features highly adhered to the global clustering. The Geweke statistics suggested that all the parameters converge to a stationary distribution. One can print similar information for the fixed effect coefficients  $\gamma_k$  using commands `print(fit.BCC2$summary.stat$GA)`, for the variance-covariance matrix of the random effects  $\Sigma_k$  using commands `print(fit.BCC2$summary.stat$SIGMA.SQ.U)`, for the residual variance of continuous features  $\sigma_k$  using commands `print(fit.BCC2$summary.stat$SIGMA.SQ.E)`, for  $k = 1, \dots, K$ . To visualize the MCMC chain for model parameters, one can use the `traceplot()` function to generate the trace plots. For example, to generate a trace plot for cluster probabilities  $\pi = (\pi_1, \pi_2)$  and the adherence parameters  $\alpha = (\alpha_1, \alpha_2, \alpha_3)$ , the following commands were used:

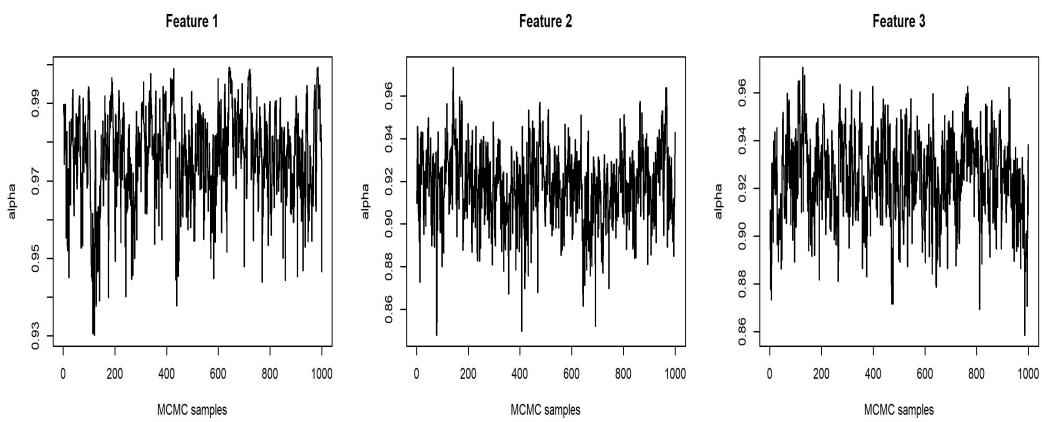
```
R> traceplot(fit = fit.BCC2, parameter = "PPI", ylab = "pi", xlab = "MCMC samples")
R> traceplot(fit = fit.BCC2, parameter = "ALPHA", ylab = "alpha", xlab = "MCMC samples")
```

The outputs are displayed in Figure 3 and Figure 4, respectively. These trace plots suggested that  $\pi_1, \pi_2, \alpha_1, \alpha_2$ , and  $\alpha_3$  are mixing well and converging. In addition, one can generate trace plots for cluster-specific parameters, for example, the fixed effect regression coefficients,  $\gamma_{k,r} = (\gamma_{k,r1}, \gamma_{k,r2})$ , with the first element as the intercept and the second element as the slope. The following commands generate trace plots (results not shown) for Cluster 1 of the first feature (anxiety), second (depress), and third (aep) features, that is,  $\gamma_{1,1} = (\gamma_{1,11}, \gamma_{1,12})$ ,  $\gamma_{1,2} = (\gamma_{1,21}, \gamma_{1,22})$ , and  $\gamma_{1,3} = (\gamma_{1,31}, \gamma_{1,32})$ .

```
R> traceplot(fit = fit.BCC2, cluster.indx = 1, feature.indx = 1,
+            parameter = "GA", ylab = "GA", xlab = "MCMC samples")
R> traceplot(fit = fit.BCC2, cluster.indx = 1, feature.indx = 2,
+            parameter = "GA", ylab = "GA", xlab = "MCMC samples")
R> traceplot(fit = fit.BCC2, cluster.indx = 1, feature.indx = 3,
+            parameter = "GA", ylab = "GA", xlab = "MCMC samples")
```



**Figure 3:** Trace plots for cluster probabilities  $\pi_1$  and  $\pi_2$ .



**Figure 4:** Trace plots for adherence parameters  $\alpha_1, \alpha_2, \alpha_3$ .

Also, the following commands generated trace plots (results not shown) for Cluster 2 of the first feature (anxiety), second (depress), and third (aep) features, that is,  $\gamma_{2,1} = (\gamma_{2,11}, \gamma_{2,12})$ ,  $\gamma_{2,2} = (\gamma_{2,21}, \gamma_{2,22})$ , and  $\gamma_{2,3} = (\gamma_{2,31}, \gamma_{2,32})$ .

```
R> traceplot(fit = fit.BCC2, cluster.indx = 2, feature.indx = 1,
+            parameter = "GA", ylab = "GA", xlab="MCMC samples")
R> traceplot(fit = fit.BCC2, cluster.indx = 2, feature.indx = 2,
+            parameter = "GA", ylab = "GA", xlab = "MCMC samples")
R> traceplot(fit = fit.BCC2, cluster.indx = 2, feature.indx = 3,
+            parameter = "GA", ylab = "GA", xlab = "MCMC samples")
```

Similarly, one can use `parameter = "SIGMA.SQ.U"` to generate the elements of the variance-covariance matrix ( $\Sigma_k$ ) for the random effects, and `parameter = "SIGMA.SQ.E"` for the residual variance of continuous features  $\sigma_k$ , for  $k = 1, \dots, K$ .

### Longitudinal Profiles by Local and Global Clusters

The package uses the mode over all the MCMC samples (after burn-in and thinning) as the point estimates for both feature-specific and global clusterings. The following commands computed the number of individuals belonging to each cluster for both feature-specific and global clusterings.

```
R> table(fit.BCC2$cluster.local[[1]])
  1   2
231 309
R> table(fit.BCC2$cluster.local[[2]])
  1   2
223 317
```

```
R> table(fit.BCC2$cluster.local[[3]])
  1   2
267 273
R> table(fit.BCC2$cluster.global)
  1   2
236 304
```

That is, for clustering based on the first feature (anxiety)  $L_1$ , 231 individuals were assigned to Cluster 1 and 309 were assigned to Cluster 2. For clustering based on the second feature (depress)  $L_2$ , 223 individuals were assigned to Cluster 1 and 317 were assigned to Cluster 2. For clustering based on the third feature (aep)  $L_3$ , 267 individuals were assigned to Cluster 1 and 273 were assigned to Cluster 2. For global clustering,  $C$ , 236 individuals were assigned to Cluster 1 and 304 were assigned to Cluster 2. To plot the longitudinal trajectory of features by local and global clusterings, the **trajplot()** function can be used. Example codes for plotting the trajectory of the features in the **epileptic.qol** data are provided as follows:

```
R> gp1 <- trajplot(fit = fit.BCC2, feature.ind = 1,
+   which.cluster = "local.cluster",
+   title= bquote(paste("Local Clustering (",
+   hat(alpha)[1] ==.(round(fit.BCC2$alpha[1],2)),")")),
+   xlab = "time (months)", ylab = "anxiety", color = c("#00BA38", "#619CFF"))
R> gp2 <- trajplot(fit = fit.BCC2, feature.ind = 2,
+   which.cluster = "local.cluster",
+   title = bquote(paste("Local Clustering (",
+   hat(alpha)[2] ==.(round(fit.BCC2$alpha[2],2)),")")),
+   xlab = "time (months)", ylab = "depress", color = c("#00BA38", "#619CFF"))
R> gp3 <- trajplot(fit = fit.BCC2, feature.ind = 3,
+   which.cluster = "local.cluster",
+   title = bquote(paste("Local Clustering (",
+   hat(alpha)[3] ==.(round(fit.BCC2$alpha[3],2)),")")),
+   xlab = "time (months)", ylab = "aep", color = c("#00BA38", "#619CFF"))
R> gp4 <- trajplot(fit = fit.BCC2, feature.ind = 1,
+   which.cluster = "global.cluster",
+   title = "Global Clustering", xlab = "time (months)", ylab = "anxiety",
+   color = c("#00BA38", "#619CFF"))
R> gp5 <- trajplot(fit = fit.BCC2, feature.ind = 2,
+   which.cluster = "global.cluster",
+   title = "Global Clustering", xlab = "time (months)", ylab = "depress",
+   color = c("#00BA38", "#619CFF"))
R> gp6 <- trajplot(fit = fit.BCC2, feature.ind = 3,
+   which.cluster = "global.cluster",
+   title = "Global Clustering", xlab = "time (months)", ylab = "aep",
+   color = c("#00BA38", "#619CFF"))
```

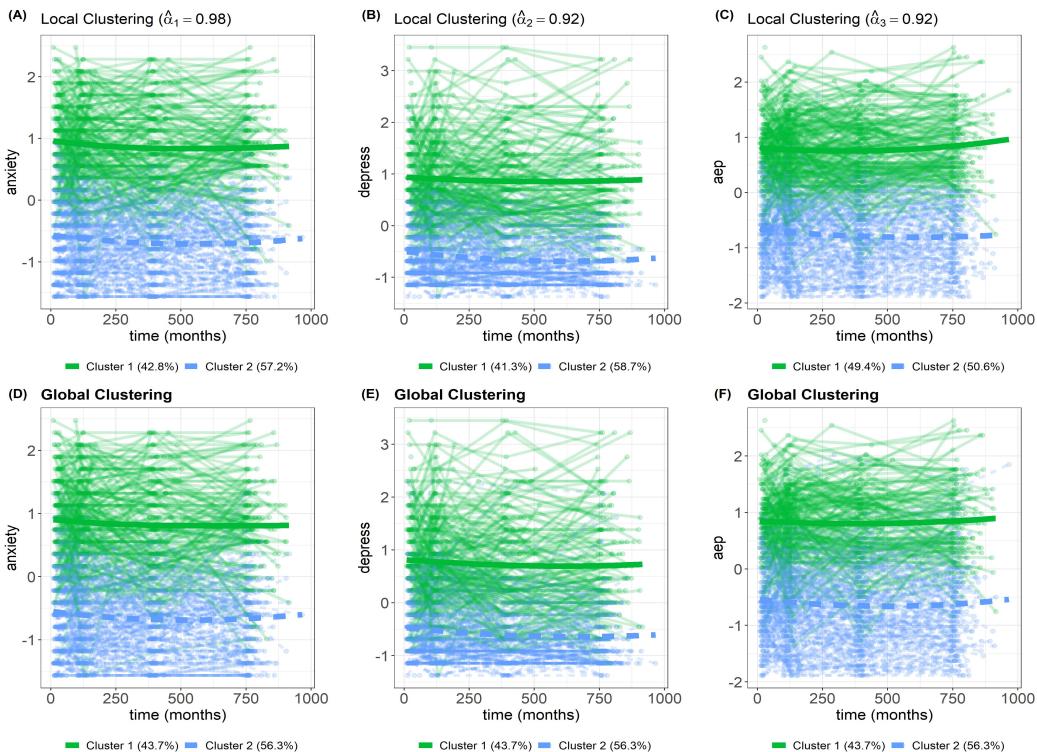
The **cowplot** package ([Wilke, 2024](#)) is used here to combine all six plots together into one figure.

```
R> library("cowplot")
R> dev.new(width = 180, height = 120)
R> plot_grid(gp1, gp2, gp3, gp4, gp5, gp6,
+   labels = c("(A)", "(B)", "(C)", "(D)", "(E)", "(F)" ),
+   ncol = 3, align = "v" )
```

The results are displayed in Figure 5. The top panel (A, B, and C) shows the three features plotted by local clustering indices, whereas the bottom panel (D, E, F) shows the three features plotted by global clustering indices. Individuals in Cluster 1 had lower anxiety scores, depression scores, and aep scores, which represented a better health condition compared to those in Cluster 2. The local clustering based on the anxiety score highly adhered to the global clustering ( $\hat{\alpha}_1 = 0.98$ ), followed by aep ( $\hat{\alpha}_2 = 0.92$ ) and depress scores ( $\hat{\alpha}_3 = 0.92$ ), respectively. This suggested that all three features highly adhered to the global clustering, and the anxiety score contributed the most information to determine the global clustering.

## Goodness of Fit

To assess the model goodness of fit, one can use **BayesT0** for computing the posterior predictive check. The commands of the two-cluster model for **epileptic.qol** data are as follows:



**Figure 5:** Longitudinal trajectories for features by local and global clusterings **epileptic.qol** data. Locally weighted scatterplot smoothing curves are overlaid on each panel to provide an estimate of the overall trend. (A) Longitudinal trajectories of anxiety plotted by local clustering  $L_1$ . (B) Longitudinal trajectories of depress plotted by local clustering  $L_2$ . (C) Longitudinal trajectories of aep plotted by local clustering  $L_3$ . (D) Longitudinal trajectories of anxiety plotted by global clustering  $C$ . (E) Longitudinal trajectories of depress plotted by global clustering  $C$ . (F) Longitudinal trajectories of aep plotted by global clustering  $C$ .

```
R> res <- Bayest(fit = fit.BCC2)
R> plot(log(res$T.obs), log(res$T.rep), xlim = c(8.45,8.7), ylim = c(8.45,8.7), cex = 1.5,
+       xlab = "Observed T statistics (in log scale)",
+       ylab = "Predicted T statistics (in log scale)")
R> abline(0, 1, lwd = 2,col = 2)
```

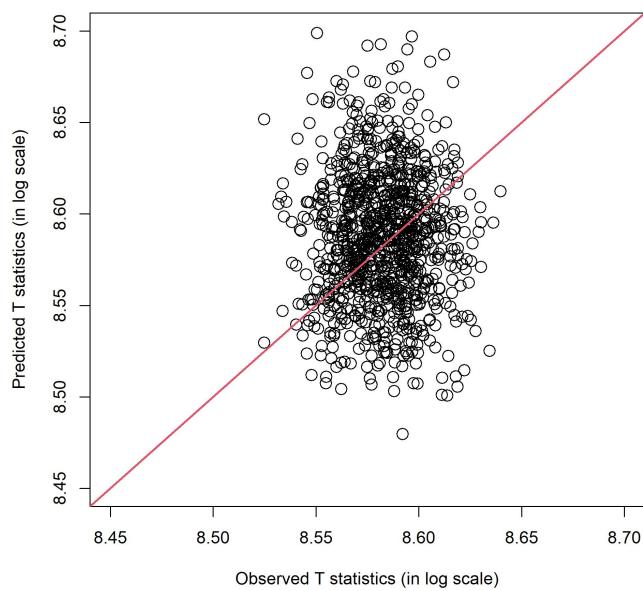
The result is displayed in Figure 6. This figure suggested that there was no systematic pattern (e.g., consistent overestimation or underestimation) between the observed  $T$  (in log scale) and predicted  $T$  (in log scale) statistics, indicating that the model provided a good fit to the data. To further evaluate the goodness of fit using an objective measure, the Bayesian p-value can be calculated, using the following commands,

```
R> p.value <- sum(res$T.rep > res$T.obs)/length(res$T.rep)
R> p.value
[1] 0.559
```

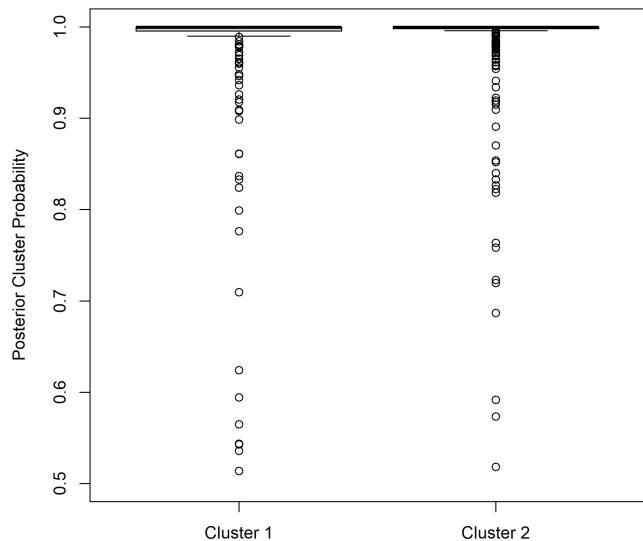
The Bayesian p-value of 0.559 corresponds to the proportion of samples above the diagonal, further suggesting that there is no clear evidence of model misspecification or that the model provides a poor fit to the data. In addition, the output also provided the posterior cluster probability, which is a measure of cluster uncertainty. The following commands generated a boxplot for the posterior cluster probabilities by clusters.

```
R> fit.BCC2$cluster.global <- factor(fit.BCC2$cluster.global,
+                                       labels = c("Cluster 1", "Cluster 2"))
R> boxplot(fit.BCC2$postprob ~ fit.BCC2$cluster.global, ylim = c(0.5, 1),
+           xlab = "", ylab = "Posterior Cluster Probability")
```

The results were displayed in Figure 7. The posterior cluster probabilities for a majority of individuals in both Clusters 1 and 2 were close to 1, except for a few individuals. This suggested that the cluster membership for individuals is robust.



**Figure 6:** Posterior predictive check: observed  $T$  (in log scale) and predicted  $T$  (in log scale) statistics.



**Figure 7:** Boxplot for the posterior cluster probabilities by clusters.

### Using a Different Set of Initial Values for Local and Global Cluster Memberships

It is a good practice to fit the model with different starting values to ensure the model converges and the clustering results are stable. By default, the initial values of the local cluster membership are randomly generated from 1 to  $K$  for each feature, and the global cluster membership is set to be equal to the local cluster membership of the first feature. One can also supply the local cluster membership by fitting the Bayesian mixture model for longitudinal data to each feature through the mixAK package for the given  $K$ . To do so, the argument `initial.cluster.membership = "mixAK"` can be added to the BCC.multi function. One can also provide user input initial values for the cluster membership, by using `initial.cluster.membership = "input"`, and then supply the initial values using the `input.initial.local.cluster.membership` and `input.initial.global.cluster.membership` arguments.

ments. For example, if the user would like to refit the model using clustering results generated from the previous model (`fit.BCC2`), the following code can be used:

```
> fit.BCC2a <- BCC.multi (
+   mydat = list(dat$anxiety_scale, dat$depress_scale, dat$aep_scale),
+   dist = c("gaussian"),
+   id = list(dat$id),
+   time = list(dat$time),
+   formula = list(y ~ time + (1|id)),
+   initial.cluster.membership = "input",
+   input.initial.local.cluster.membership = list(fit.BCC2$cluster.local[[1]],
+   fit.BCC2$cluster.local[[2]], fit.BCC2$cluster.local[[3]]),
+   input.initial.global.cluster.membership = fit.BCC2$cluster.global,
+   num.cluster = 2,
+   burn.in = 1000,
+   thin = 1,
+   per = 100,
+   max.iter = 2000)
```

### Comparing Clustering Results to Other Model Specifications

To evaluate the robustness of the resulting clusters, one can compare the clustering results to more complicated models. We fit a BCC model with a linear form using a random intercept and slope, and a BCC model with a quadratic form using a random intercept and slope. The commands are as follows:

```
R> fit.BCC2b <- BCC.multi (
+   mydat = list(dat$anxiety_scale, dat$depress_scale, dat$aep_scale),
+   dist = c("gaussian"),
+   id = list(dat$id),
+   time = list(dat$time),
+   formula = list(y ~ time + (1 + time|id)),
+   num.cluster = 2,
+   print.info = "FALSE",
+   burn.in = 1000,
+   thin = 1,
+   per = 100,
+   max.iter = 2000)
R> fit.BCC2c <- BCC.multi (
+   mydat = list(dat$anxiety_scale, dat$depress_scale, dat$aep_scale),
+   dist = c("gaussian"),
+   id = list(dat$id),
+   time = list(dat$time),
+   formula = list(y ~ time + time2 + (1 + time|id)),
+   num.cluster = 2,
+   print.info = "FALSE",
+   burn.in = 1000,
+   thin = 1,
+   per = 100,
+   max.iter = 2000)
> fit.BCC2b$cluster.global <- factor(fit.BCC2b$cluster.global,
+   labels = c("Cluster 1", "Cluster 2"))
> table(fit.BCC2b$cluster.global, fit.BCC2b$cluster.global)
```

	Cluster 1	Cluster 2
Cluster 1	223	13
Cluster 2	1	303

The agreement between the first model (linear model with a random intercept) and the second model (linear model with a random intercept and random slope) was  $(223 + 303)/540 = 97\%$ .

```
R> fit.BCC2c$cluster.global <- factor(fit.BCC2c$cluster.global,
+   labels = c("Cluster 1", "Cluster 2"))
R> table(fit.BCC2c$cluster.global, fit.BCC2c$cluster.global)

      Cluster 1 Cluster 2
Cluster 1      227       9
```

Cluster 2	2	302
-----------	---	-----

The agreement between the first model (linear model with a random intercept) and the second model (quadratic model with a random intercept and random slope) was  $(226 + 301)/540 = 98\%$ . This suggested that fitting a more complicated model does not substantially change the individual cluster membership, and further attests that the estimated cluster membership was robust to different model specifications.

### 3.2 Example 2: Mayo Clinic Primary Biliary Cholangitis Data

The second dataset was from a randomized placebo-controlled trial of the drug D-penicillamine for Primary Biliary Cholangitis conducted between 1974 and 1984 at the Mayo Clinic Primary Biliary Cholangitis (PBC). The dataset used here contains multiple laboratory results collected longitudinally on 312 randomized PBC patients. Example features (variables) include serum bilirubin (mg/dl), serum albumin (mg/dl), platelet count, serum glutamic-oxaloacetic transaminase and the presence of blood vessel malformations in the skin. The dataset was analyzed previously by Komárek and Komárková (2013) using the **mixAK** package. The analysis goal was to identify subgroups of patients with similar profiles based on multiple features which may serve as critical information regarding patients' prognostics. To illustrate the clustering process of mixed-type longitudinal features, we followed Komárek and Komárková (2013) analysis and chose three features of interest, namely the *lbili* (log of serum bilirubin), platelet count, and *spiders* (presence of blood vessel malformations in the skin). Also, following Komárek and Komárková (2013), we included  $N = 260$  individuals known to be alive at 910 days of follow-up, and only the longitudinal measurements up to this point were considered (known as the **PBC910** data). The median (min, max) number of observations for these markers was 4 (1, 5). Given a small number of observations available for each individual, we considered a model with only a random intercept, i.e., using the following specification:

```
formula = list(y ~ time + (1|id))
```

For a complete analysis, one can follow the steps outlined in Example 1 to determine the number of clusters and to evaluate the model performance using several different approaches (e.g., posterior predictive check). Here, as an illustration, and for the purpose of comparison, we only fitted a model with two clusters (i.e.,  $K = 2$ ) as in Komárek and Komárková (2013). The data is available through the **mixAK** package. The following commands prepared the dataset for analysis.

```
R> library("mixAK")
R> data(PBC910)
```

Using the argument `dist = c("gaussian", "poisson", "binomial")` to specify the distributions for the three features, the following commands performed a BCC model with  $K = 2$ .

```
R > fit.BCC2 <- BCC.multi(
+   mydat = list(PBC910$lbili, PBC910$platelet, PBC910$spiders),
+   dist = c("gaussian", "poisson", "binomial"),
+   id = list(PBC910$id),
+   time = list(PBC910$month),
+   formula = list(y ~ time + (1|id)),
+   num.cluster = 2,
+   burn.in = 10000,
+   thin = 10,
+   per = 1000,
+   max.iter = 20000)
```

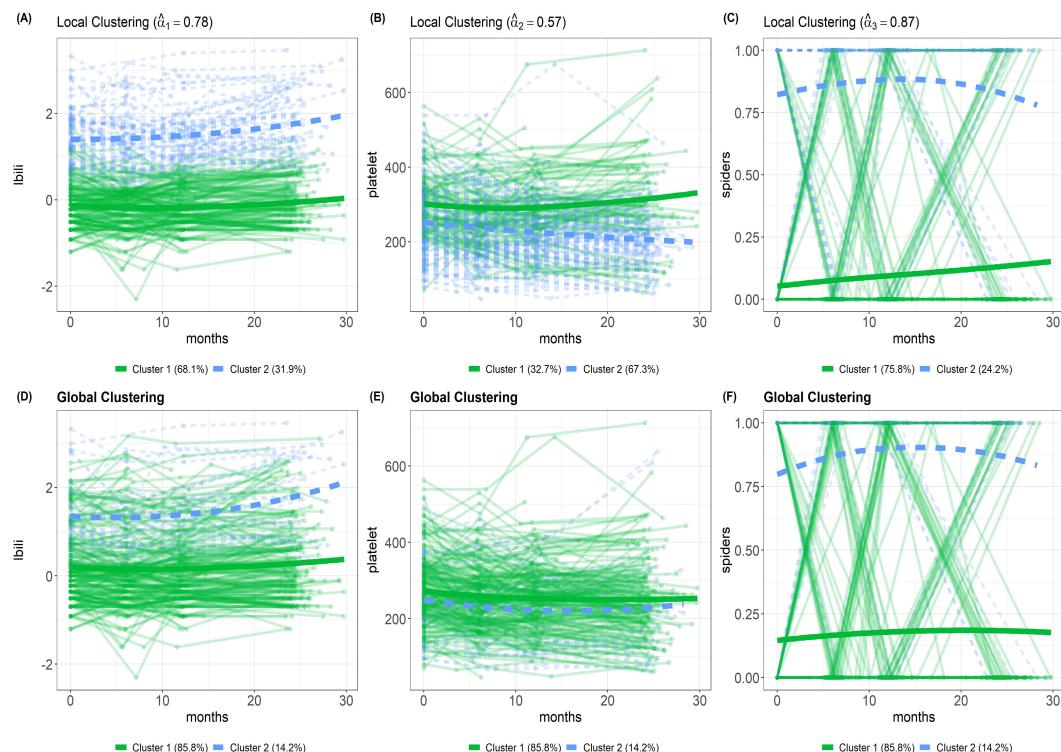
The following commands generated the longitudinal profile plots for the three features (*lbili*, *platelet*, *spiders*).

```
R> gp1 <- trajplot(fit = fit.BCC2, feature.ind = 1,
+   which.cluster = "local.cluster",
+   title = bquote(paste("Local Clustering (",
+   hat(alpha)[1] ==.(round(fit.BCC2$alpha[1],2)),")")),
+   xlab = "months", ylab = "lbili", color = c("#00BA38", "#619CFF"))
R> gp2 <- trajplot(fit = fit.BCC2, feature.ind = 2,
+   which.cluster = "local.cluster",
+   title = bquote(paste("Local Clustering (",
+   hat(alpha)[2] ==.(round(fit.BCC2$alpha[2],2)),")")),
+   xlab = "months", ylab = "platelet", color = c("#00BA38", "#619CFF"))
R> gp3 <- trajplot(fit = fit.BCC2, feature.ind = 3,
```

```

+     which.cluster = "local.cluster",
+     title = bquote(paste("Local Clustering (",
+     hat(alpha)[3] ==.(round(fit.BCC2$alpha[3],2)),")")),
+     xlab = "months", ylab = "spiders", color = c("#00BA38", "#619cff"))
R> gp4 <- trajplot(fit = fit.BCC2, feature.ind = 1,
+     which.cluster = "global.cluster",
+     title = "Global Clustering",
+     xlab = "months", ylab = "lbili", color = c("#00BA38", "#619cff"))
R> gp5 <- trajplot(fit = fit.BCC2, feature.ind = 2,
+     which.cluster = "global.cluster",
+     title = "Global Clustering",
+     xlab = "months", ylab = "platelet", color = c("#00BA38", "#619cff"))
R> gp6 <- trajplot(fit = fit.BCC2, feature.ind = 3,
+     which.cluster = "global.cluster",
+     title = "Global Clustering",
+     xlab = "months", ylab = "spiders", color = c("#00BA38", "#619cff"))

```



**Figure 8:** Longitudinal trajectories for features by local and global clusterings for PBC910 data. Locally weighted scatterplot smoothing curves are overlaid on each panel to provide an estimate of the overall trend. (A) Longitudinal trajectories of lbili plotted by local clustering  $L_1$ . (B) Longitudinal trajectories of platelet plotted by local clustering  $L_2$ . (C) Longitudinal trajectories of spiders plotted by local clustering  $L_3$ . (D) Longitudinal trajectories of lbili plotted by global clustering  $C$ . (E) Longitudinal trajectories of platelet plotted by global clustering  $C$ . (F) Longitudinal trajectories of spiders plotted by global clustering  $C$ .

To compare the results with an existing model, we fit a Bayesian mixture model using the **mixAK** package using the following commands (see (Komárek and Komárková, 2014) for details of using this package):

```

R> mod <- GLMM_MCMC(y = PBC910[, c("lbili", "platelet", "spiders")],
+     dist = c("gaussian", "poisson(log)", "binomial(logit"),
+     id = PBC910[, "id"],
+     x = list(lbili = "empty", platelet = "empty", spiders = PBC910[, "month"]),
+     z = list(lbili = PBC910[, "month"], platelet = PBC910[, "month"], spiders = "empty"),
+     random.intercept = rep(TRUE, 3), prior.b = list(Kmax = 2),
+     nMCMC = c(burn = 100, keep = 1000, thin = 10, info = 100),
+     parallel = FALSE)
R> mod <- NMixRelabel(mod, type = "stephens", keep.comp.prob = TRUE)

```

```
R> cluster.mixAK <- apply(mod[[1]]$poster.comp.prob, 1, which.max)
R> table(mixAK = cluster.mixAK, BCCLong = fit.BCC2$cluster.global)

BCCLong
mixAK   1   2
  1 158   3
  2  65  34
```

The agreement between the two models for the **PBC910** data was  $(158 + 34)/260 = 74\%$ .

## 4 Summary and discussion

In this paper, we described a BCC model for clustering longitudinal data with multiple features. Using two real-life and open-access data, we provided step-by-step guidance on using the **BCCLong** package and interpreting the results. The results also demonstrated that the **BCCLong** package is a useful tool for clustering longitudinal data and it enhanced our understanding of the heterogeneity underlying each feature and across features. Of note, the reason we require the user to create three separate lists of data vectors (instead of a data frame) when using `BCC.multi()`, in specifying the arguments `mydat`, `id`, and `time` is to allow more flexibility in modeling longitudinal features with distinct structures. This is particularly useful when these features are collected from multiple data sources (with different study designs), and therefore the time scale and data collection frequency could differ between features. Additional details regarding the functions and their arguments can be found in the package vignettes. The BCC model implemented using the **BCCLong** package yields both feature-specific (local) clusterings and consensus (global) clustering. In practice, global clustering is often of greater interest, as it encompasses information from all features. It can be used to associate with exposures and predict long-term outcomes. Conceptually, feature-specific clustering can be viewed as clustering based on a single feature alone, whereas global clustering can be viewed as a weighted average clustering across all features, with weights being the adherence parameters. In order to capture the complexity of the underlying population heterogeneity, it is of great importance to consider multiple features of interest simultaneously, as we often do for cross-sectional data. Several directions can be considered in future studies. For example, incorporating variable selection (e.g., spike-and-slab and shrinkage) priors (Lu and Lou, 2022a, 2021) for fixed and random effects will enhance the flexibility of the package and allow it to be applied to high-dimensional settings. In addition, determining the number of clusters can be achieved by using a Dirichlet process mixture model (Escobar, 1994; Lu and Chandra, 2024; Lu et al., 2024). The package does not support predicting the cluster membership of a new individual or future trajectories for a given individual. These functions will be considered in future versions of the package. Finally, the current version of the **BCCLong** package relies on a standard MCMC algorithm to estimate the model, which is generally slow and requires particular care for the user to check convergence. While the computational time is reasonable for our current applications, to scale up the model to a large dataset, other computational algorithms such as distributed stochastic gradient MCMC (Ahn et al., 2014) and variational inference (Blei et al., 2017) can be considered in future studies.

## Computational details

Most of the functions in the **BCCLong** package are written in C++. The output shown in this article was obtained using R version 4.2.1 (June, 2022), **BCCLong** 1.0.3, and the following contributed packages which are the dependencies or imports of the **BCCLong** package: `cluster` (2.1.4) (Maechler et al., 2022), `coda` (0.19-4) (Plummer et al., 2006), `ggplot2` (3.4.0) (Wickham, 2016), `label.switching` (1.8) (Papastamoulis, 2016), `LaplacesDemon` (16.1.6) (Hall, 2021), `lme4` (1.1-31) (Bates et al., 2015), `MASS` (7.3-58.1) (Ripley et al., 2013), `mclust` (6.0.0) (Scrucca et al., 2023), `MCMCpack` (1.6-3) (Martin et al., 2011), `mixAK` (5.5) (Komárek and Komárová, 2014), `mvtnorm` (1.1-3) (Genz et al., 2021), `nnet` (7.3-18) (Venables and Ripley, 2002), `Rcpp` (1.0.9) (Eddelbuettel and François, 2011), `Rmpfr` (0.8-9) (Maechler, 2022), `truncdist` (1.0-2) (Novomestky et al., 2016). R itself and all packages used are available from the Comprehensive R Archive Network (CRAN) at <https://CRAN.R-project.org/>.

## Acknowledgments

ZT is supported by a Dean's Doctoral Award from Queen's University. ZL is supported by a Discovery Grant funded by the Natural Sciences and Engineering Research Council of Canada.

## References

- S. Ahn, B. Shahbaba, and M. Welling. Distributed stochastic gradient MCMC. In *International conference on machine learning*, pages 1044–1052. PMLR, 2014. URL <https://proceedings.mlr.press/v32/ahn14.html>. [p57]
- D. Bates, M. Mächler, B. Bolker, and S. Walker. Fitting linear mixed-effects models using lme4. *Journal of Statistical Software*, 67(1):1–48, 2015. doi: 10.18637/jss.v067.i01. R package version 1.1-31. [p57]
- D. M. Blei, A. Kucukelbir, and J. D. McAuliffe. Variational inference: A review for statisticians. *Journal of the American Statistical Association*, 112(518):859–877, 2017. URL <https://doi.org/10.1080/01621459.2017.1285773>. [p57]
- G. Celeux, F. Forbes, C. P. Robert, D. M. Titterington, et al. Deviance information criteria for missing data models. *Bayesian Analysis*, 1(4):651–673, 2006. doi: 10.1214/06-BA122. [p45]
- D. Eddelbuettel and R. François. Rcpp: Seamless R and C++ integration. *Journal of Statistical Software*, 40(8):1–18, 2011. doi: 10.18637/jss.v040.i08. R package version 1.0.9. [p57]
- M. R. Elliott, J. J. Gallo, T. R. Ten Have, H. R. Bogner, and I. R. Katz. Using a bayesian latent growth curve model to identify trajectories of positive affect and negative events following myocardial infarction. *Biostatistics*, 6(1):119–143, 2005. doi: 10.1093/biostatistics/kxh022. [p45]
- M. D. Escobar. Estimating normal means with a Dirichlet process prior. *Journal of the American Statistical Association*, 89(425):268–277, 1994. URL <https://doi.org/10.2307/2291223>. [p57]
- S. Frühwirth-Schnatter and S. Pyne. Bayesian inference for finite mixtures of univariate and multivariate skew-normal and skew-t distributions. *Biostatistics*, 11(2):317–336, 2010. URL <https://doi.org/10.1093/biostatistics/kxp062>. [p45]
- A. Gelman, X.-L. Meng, and H. Stern. Posterior predictive assessment of model fitness via realized discrepancies. *Statistica Sinica*, pages 733–760, 1996. URL <https://www.jstor.org/stable/24306036>. [p46]
- A. Gelman, J. Hwang, and A. Vehtari. Understanding predictive information criteria for Bayesian models. *Statistics and Computing*, 24(6):997–1016, 2014. URL <https://doi.org/10.1007/s11222-013-9416-2>. [p45]
- C. Genolini and B. Falissard. Kml: k-means for longitudinal data. *Computational Statistics*, 25(2):317–328, 2009. doi: 10.1007/s00180-009-0178-4. [p41]
- C. Genolini, J.-B. Pingault, T. Driss, S. Côté, R. E. Tremblay, F. Vitaro, C. Arnaud, and B. Falissard. Kml3d: a non-parametric algorithm for clustering joint trajectories. *Computer Methods and Programs in Biomedicine*, 109(1):104–111, 2013. URL <https://doi.org/10.1016/j.cmpb.2012.08.016>. [p41]
- A. Genz, F. Bretz, T. Miwa, X. Mi, F. Leisch, F. Scheipl, and T. Hothorn. mvtnorm: Multivariate Normal and t Distributions, 2021. URL <https://CRAN.R-project.org/package=mvtnorm>. R package version 1.1-3. [p57]
- J. Geweke. *Evaluating the accuracy of sampling-based approaches to the calculation of posterior moments*, volume 196. Federal Reserve Bank of Minneapolis, Research Department Minneapolis, MN, USA, 1991. [p44]
- B. Hall. LaplacesDemon: Complete Environment for Bayesian Inference, 2021. URL <https://CRAN.R-project.org/package=LaplaceDemon>. R package version 16.1.6. [p45, 57]
- B. C. Heggeseth and N. P. Jewell. How Gaussian mixture models might miss detecting factors that impact growth patterns. *The Annals of Applied Statistics*, 12(1):222–245, 2018. URL <https://doi.org/10.1214/17-aos1066>. [p47]
- G. L. Hickey, P. Philipson, A. Jorgensen, and R. Kolamunnage-Dona. joineRML: a joint model and software package for time-to-event and multivariate longitudinal outcomes. *BMC Medical Research Methodology*, 18(1):1–14, 2018. URL <https://doi.org/10.1186/s12874-018-0502-1>. [p46]
- A. Komárek and L. Komárová. Clustering for multivariate continuous and discrete longitudinal data. *The Annals of Applied Statistics*, 7(1):177–200, 2013. doi: 10.1214/12-AOAS580. [p55]
- A. Komárek and L. Komárová. Capabilities of R package mixak for clustering based on multivariate continuous and discrete longitudinal data. *Journal of Statistical Software*, 59(12):1–38, 2014. URL <https://doi.org/10.18637/jss.v059.i12>. R package version 5.5. [p46, 56, 57]

- B. E. Leiby, M. D. Sammel, T. R. Ten Have, and K. G. Lynch. Identification of multivariate responders and non-responders by using Bayesian growth curve latent class models. *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, 58(4):505–524, 2009. URL <https://doi.org/10.1111/j.1467-9876.2009.00663.x>. [p45]
- F. Leisch. FlexMix: A general framework for finite mixture models and latent class regression in R. *Journal of Statistical Software*, 11(8):1–18, 2004. doi: 10.18637/jss.v011.i08. [p41]
- E. F. Lock and D. B. Dunson. Bayesian consensus clustering. *Bioinformatics*, 29(20):2610–2616, 2013. doi: 10.1093/bioinformatics/btt425. [p42, 45]
- Z. Lu and N. K. Chandra. A sparse factor model for clustering high-dimensional longitudinal data. *Statistics in Medicine*, 43(19):3633–3648, 2024. URL [doi:10.1002/sim.10151](https://doi.org/10.1002/sim.10151). [p57]
- Z. Lu and W. Lou. Shape invariant mixture model for clustering non-linear longitudinal growth trajectories. *Statistical Methods in Medical Research*, 28(12):3769–3784, 2019. URL <https://doi.org/10.1177/0962280218815301>. [p45]
- Z. Lu and W. Lou. Bayesian approaches to variable selection in mixture models with application to disease clustering. *Journal of Applied Statistics*, pages 1–21, 2021. URL <https://doi.org/10.1080/02664763.2021.1994529>. [p57]
- Z. Lu and W. Lou. Bayesian approaches to variable selection: a comparative study from practical perspectives. *International Journal of Biostatistics*, 18(1):89–108, May 2022a. URL <https://doi.org/10.1515/ijb-2020-0130>. [p57]
- Z. Lu and W. Lou. Bayesian consensus clustering for multivariate longitudinal data. *Statistics in Medicine*, 41(1):108–127, 2022b. doi: 10.1002/sim.9225. [p41, 42, 45]
- Z. Lu, M. Ahmadiankalati, and Z. Tan. Joint clustering multiple longitudinal features: A comparison of methods and software packages with practical guidance. *Statistics in Medicine*, 42(29):5513–5540, 2023. URL <https://doi.org/10.1002/sim.9917>. [p41]
- Z. Lu, P. Subbarao, and W. Lou. A Bayesian latent class model for integrating multi-source longitudinal data: application to the CHILD cohort study. *Journal of the Royal Statistical Society Series C: Applied Statistics*, 73(2):398–419, 2024. URL <https://doi.org/10.1093/rssc/qlad100>. [p57]
- M. Maechler. *Rmpfr: R MPFR - Multiple Precision Floating-Point Reliable*, 2022. URL <https://CRAN.R-project.org/package=Rmpfr>. R package version 0.8-9. [p57]
- M. Maechler, P. Rousseeuw, A. Struyf, M. Hubert, and K. Hornik. *cluster: Cluster Analysis Basics and Extensions*, 2022. URL <https://CRAN.R-project.org/package=cluster>. R package version 2.1.4. [p57]
- A. Magrini. *gbmt: Group-Based Multivariate Trajectory Modeling*, 2022a. URL <https://CRAN.R-project.org/package=gbmt>. R package version 0.1.3. [p41]
- A. Magrini. Assessment of agricultural sustainability in european union countries: A group-based multivariate trajectory approach. *Asta Advances in Statistical Analysis*, 106(4):673–703, 2022b. URL <https://doi.org/10.1007/s10182-022-00437-9>. [p41]
- A. Marson, R. Appleton, G. Baker, D. Chadwick, J. Doughty, B. Eaton, C. Gamble, A. Jacoby, P. Shackley, D. Smith, et al. A randomised controlled trial examining the longer-term outcomes of standard versus new antiepileptic drugs. the sanad trial. *Health Technology Assessment-southampton-*, 11(37), 2007. URL <https://doi.org/10.3310/hta11370>. [p46]
- A. D. Martin, K. M. Quinn, and J. H. Park. MCMCpack: Markov chain Monte Carlo in R. *Journal of Statistical Software*, 42(9):22, 2011. doi: 10.18637/jss.v042.i09. R package version 1.6-3. [p57]
- E. C. Merkle, D. Furr, and S. Rabe-Hesketh. Bayesian comparison of latent variable models: Conditional versus marginal likelihoods. *Psychometrika*, 84(3):802–829, 2019. URL <https://doi.org/10.1007/s11336-019-09679-0>. [p45]
- D. S. Nagin, B. L. Jones, V. L. Passos, and R. E. Tremblay. Group-based multi-trajectory modeling. *Statistical Methods in Medical Research*, 27(7):2015–2023, 2018. URL <https://doi.org/10.1177/0962280216673085>. [p41]
- K. Nasserinejad, J. van Rosmalen, W. de Kort, and E. Lesaffre. Comparison of criteria for choosing the number of classes in Bayesian finite mixture models. *Plos One*, 12(1):e0168838, 2017. URL <https://doi.org/10.1371/journal.pone.0168838>. [p45]

- B. Neelon, G. K. Swamy, L. F. Burgette, and M. L. Miranda. A Bayesian growth mixture model to examine maternal hypertension and birth outcomes. *Statistics in Medicine*, 30(22):2721–2735, 2011. URL <https://doi.org/10.1002/sim.4291>. [p45]
- F. Novomestky, S. Nadarajah, and M. F. Novomestky. *truncdist: Truncated Random Variables*, 2016. URL <https://CRAN.R-project.org/package=truncdist>. R package version 1.0-2. [p57]
- P. Papastamoulis. *label.switching*: An R package for dealing with the label switching problem in MCMC outputs. *Journal of Statistical Software, Code Snippets*, 69(1):1–24, 2016. doi: 10.18637/jss.v069.c01. R package version 1.8. [p43, 57]
- M. Plummer, N. Best, K. Cowles, and K. Vines. *Coda: Convergence diagnosis and output analysis for mcmc*. *R News*, 6(1):7–11, 2006. URL <https://journal.r-project.org/archive/>. R package version 0.19-4. [p45, 57]
- B. Ripley, B. Venables, D. M. Bates, K. Hornik, A. Gebhardt, D. Firth, and Ripley. *MASS: Support Functions and Datasets for Venables and Ripley's MASS*, 2013. URL <https://CRAN.R-project.org/package=crosstalk>. R package version 7.3-58.1. [p57]
- L. Scrucca, C. Fraley, T. B. Murphy, and A. E. Raftery. *Model-Based Clustering, Classification, and Density Estimation Using mclust in R*. Chapman and Hall/CRC, 2023. ISBN 978-1032234953. doi: 10.1201/9781003277965. URL <https://mclust-org.github.io/book/>. R package version 6.0.0. [p57]
- D. J. Spiegelhalter, N. G. Best, B. P. Carlin, and A. Van Der Linde. Bayesian measures of model complexity and fit. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 64(4): 583–639, 2002. URL <https://doi.org/10.1111/1467-9868.00353>. [p45]
- M. Stephens. Dealing with label switching in mixture models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 62(4):795–809, 2000. URL <https://doi.org/10.1111/1467-9868.00265>. [p43]
- Z. Tan, C. Shen, P. Subbarao, W. Lou, and Z. Lu. A joint modeling approach for clustering mixed-type multivariate longitudinal data: Application to the CHILD cohort study. *Arxiv Preprint Arxiv:2210.08385*, 2022. URL <https://doi.org/10.48550/arXiv.2210.08385>. [p41, 42, 45]
- W. N. Venables and B. D. Ripley. *Modern Applied Statistics with S*. Springer, New York, fourth edition, 2002. URL <https://www.stats.ox.ac.uk/pub/MASS4/>. ISBN 0-387-95457-0, R package version 7.3-18. [p57]
- S. Watanabe. Asymptotic equivalence of bayes cross validation and widely applicable information criterion in singular learning theory. *Journal of Machine Learning Research*, 11(Dec):3571–3594, 2010. URL <http://jmlr.org/papers/v11/watanabe10a.html>. [p45]
- H. Wickham. *ggplot2: Elegant Graphics for Data Analysis*. Springer-Verlag New York, 2016. ISBN 978-3-319-24277-4. URL <https://ggplot2.tidyverse.org>. [p57]
- C. O. Wilke. *cowplot: Streamlined Plot Theme and Plot Annotations for 'ggplot2'*, 2024. URL <https://CRAN.R-project.org/package=cowplot>. R package version 1.1.3. [p51]

Zhiwen Tan  
Department of Public Health Sciences  
Queen's University

Chang Shen  
Department of Electrical and Computer Engineering  
Queen's University

Zihang Lu  
Department of Public Health Sciences  
and Department of Mathematics and Statistics  
Queen's University  
[zihang.lu@queensu.ca](mailto:zihang.lu@queensu.ca)

# reslr: An R Package for Relative Sea Level Modelling

by Maeve Upton, Andrew Parnell, and Niamh Cahill

**Abstract** We present reslr, an R package to perform Bayesian modelling of relative sea level data. We include a variety of different statistical models previously proposed in the literature, with a unifying framework for loading data, fitting models, and summarising the results. Relative sea-level data often contain measurement error in multiple dimensions, and so our package allows for these to be included in the statistical models. When plotting the output sea level curves, the focus is often on comparing rates of change, and so our package allows for computation of the derivatives of sea level curves with appropriate consideration of the uncertainty. We provide a large example dataset from the Atlantic coast of North America and show some of the results that might be obtained from our package.

## 1 Introduction

Understanding the rates and spatial patterns of Relative Sea-Level (RSL) change across various timescales, spanning from decades to millennia, poses a significant challenge. The task involves analysing sparse and noisy proxy and/or instrumental data sources that often have large measurement uncertainties. To address these complexities and provide robust assessments, statistical models play a pivotal role and have become indispensable in the task of quantifying RSL changes (as examined by Cahill et al. (2015a) and Khan et al. (2015)) and in the evaluation of temporal and spatial variability (e.g., Kopp, 2013; Kopp et al., 2016; Kemp et al., 2018; Walker et al., 2021). To that end, the paleo sea-level community would benefit from a comprehensive toolset capable of analysing the historical evolution of sea-level changes across different times and locations. This motivated us to create the `reslr` package, which is available on the Comprehensive R Archive Network at <https://cran.r-project.org/web/packages/reslr> or on GitHub at <https://github.com/maeveupton/reslr>. Our package includes a suite of statistical models appropriate for modelling the complexity of sea-level over time and space, while accounting for sea-level data uncertainties and remaining computationally tractable. The output of our package provides insight into temporal and spatial sea-level variability and rates of sea-level change.

The `reslr` package includes a comprehensive dataset of proxy RSL reconstructions for 21 locations along the Atlantic coast of North America (Kemp et al., 2013). These reconstructions rely heavily on dated geological archives obtained from coastal sediments (e.g., Gehrels, 1994) or corals (e.g., Meltzner et al., 2017). Moreover, users of the package have the option to incorporate instrumental sea-level data sourced from the Permanent Service Mean Sea Level (PSMSL) online database, which provides annual RSL measurements for approximately 1,500 tide-gauge stations worldwide (Holgate et al., 2013). By offering this diverse range of data sources, the `reslr` package caters to the needs of researchers seeking to explore and analyse sea-level variations across different locations and time periods.

The `reslr` package offers a range of statistical models which include: linear regression (e.g., Ashe et al. (2019)), change point models (e.g., Cahill et al. (2015b)), integrated Gaussian process (IGP) models (e.g., Cahill et al. (2015a)), temporal splines (e.g., de Boor (1978)), spatio-temporal splines (e.g., Simpson (2018)) and generalised additive models (GAM) (e.g., Upton et al. (2023)). In all cases, a Bayesian framework is employed, facilitating the estimation of unknown parameters based on the RSL data while fully accounting for the associated uncertainties. The `reslr` package enables researchers to gain comprehensive insights into sea-level variations, leveraging the flexibility and robustness of these statistical models.

When it comes to addressing measurement uncertainty in proxy records, the `reslr` package offers two distinct approaches. The first approach involves employing the Errors-in-Variables (EIV) method, which takes into account the inherent uncertainties in the input variables (Dey et al., 2000). This method acknowledges that the input variables are not error-free and incorporates this knowledge into the analysis. The second approach offered by the package is the Noisy Input (NI) uncertainty method. This method tackles uncertainty by inflating the output noise variance with a corrective term that is directly linked to the input noise variance (McHutchon and Rasmussen, 2011). Both the EIV and NI uncertainty methods have their respective advantages, and the `reslr` package recommends the most suitable uncertainty method based on the statistical model being employed. This ensures that researchers can select the appropriate approach to effectively address measurement uncertainties within their specific analysis context.

For each model, the `reslr` package generates informative plots illustrating the model-based estimates of RSL. In the case of more complex models, like the IGP, splines, and GAMs, the resulting

plots not only provide RSL estimates but also offer insights into the rates of RSL change. Of particular significance to the paleo-sea level community, the GAM model provides estimates for separate components that represent potential drivers of RSL change. This feature enables comparisons between different components and contributes to a more comprehensive understanding of the factors influencing RSL fluctuations (Upton et al., 2023). These visual representations enable researchers to gain a clearer understanding of when and where RSL changes occurred, including the magnitude of their temporal variations. Moreover, the package grants users access to the posterior samples used to generate the plots, providing the option to delve deeper into the underlying statistical distributions and uncertainties associated with the estimated RSL changes. The combination of these outputs serves as a valuable resource for researchers, aiding in the investigation and interpretation of RSL dynamics across various spatial and temporal contexts.

The `reslr` package is uniquely tailored to address the challenges inherent in analysing historic RSL changes using proxy records. Its design, characterised by minimal functions and a user-friendly interface, draws inspiration from other packages such as `mgcv`, renowned for its diverse statistical modelling options (Wood, 2015). Crucially, there are currently no competing R packages that match the breadth of capabilities offered by `reslr`. This uniqueness stems from its capability to provide users with a selection of Bayesian statistical models and account for bi-variate uncertainty, crucial in sea-level research.

Our paper has the following structure. First, we introduce the example dataset provided within the package, which serves as the foundation for the examples presented throughout the paper. We also provide insight into additional data sources. Second, we offer an overview of the statistical models available in the package, providing necessary background information. Next, we explore the uncertainty methods employed within these statistical models. Following this, we provide a detailed description of the functionality of the `reslr` package, outlining the diverse outputs and plots accessible to users. Finally, we conclude with important remarks and discuss potential future extensions for the package's advancement. Whilst this paper is just a summary of the features of `reslr`, a more complete vignette containing examples of the full functionality of the package is available at <https://maeveupton.github.io/reslr/>.

## 2 Data and models

### 2.1 Data sources

Proxy sea-level data are vital sources of information for examining historic changes in RSL prior to the instrumental data period. A proxy refers to a characteristic that can be observed and used to estimate a variable of interest, which cannot be measured directly, and can be of physical, biological, or chemical nature (e.g., Gornitz, 2009). In sea-level studies, the proxy data can be sourced from microorganisms such as foraminifera (e.g., Edwards and Wright, 2015), geochemical measurements (e.g., Marshall, 2015), or vegetation that has accumulated in the tidal realm (e.g., Kemp and Telford, 2015). The datasets we use have had their proxy measurements transformed into sea level using various techniques which are beyond the scope of our paper (e.g., Gehrels (1994), Shennan et al. (2015), Kemp et al. (2018)). In the `reslr` package, we provide an example proxy dataset which contains 21 proxy sea-level records (See Appendix) from the Atlantic coast of North America as used in Upton et al. (2023).

Within the context of sea-level analysis, instrumental data plays an important role by providing direct measurements obtained from tide gauges and satellites (although the latter is currently not incorporated into the `reslr` package). To enhance the versatility of the package, we have implemented a feature that allows users to download annual tide-gauge data from the PSMSL Level online database and store it in a temporary file, making it readily available when needed (PSMSL, 2023; Holgate et al. (2013); Woodworth and Player (2003)).

To ensure the comparability of the tide-gauge data with proxy records, we apply two processing steps. First, the tide-gauge data in the PSMSL database is given in millimetres relative to a revised local reference datum (a coordinate system that defines the zero level for sea level measurements Pugh and Woodworth (2014)). Within `reslr`, we transform the data by removing 7,000 mm to revert the tide-gauge data into the observed reference frame and convert the RSL to metres following the guidance from the PSMSL website as described in Aarup et al. (2006). The second processing step involves averaging the tide-gauge data to align with the temporal resolution of the more recent proxy data. However, we provide flexibility for users to adjust this averaging period according to the specific characteristics of their data.

## 2.2 Statistical Models

Within the `reslr` package, a Bayesian hierarchical framework is employed for each statistical modelling technique. Markov Chain Monte Carlo (MCMC) simulations are carried out using the Just Another Gibbs Sampler (JAGS) tool (Plummer, 2003) and implemented using the `rjags` package in R (Plummer et al., 2016). Other tools, such as Stan (Carpenter et al., 2017), are used for MCMC simulations. However, our preference for JAGS is based on its user-friendly interface, computational efficiency, and flexible model design capabilities.

Mathematically, the data level for each statistical model is described as:

$$y = f(\mathbf{x}, t) + \epsilon_y \quad (1)$$

where  $y$  is the response data (RSL in metres).  $f(\mathbf{x}, t)$  is the process mean that depends on location  $\mathbf{x}$  and time  $t$ .  $\epsilon_y$  is the error term given by  $\epsilon_y \sim N(0, \sigma_y^2 + s_y^2)$ , where  $\sigma_y^2$  is the residual variance and  $s_y$  the known measurement error associated with RSL, which occurs during the collection of proxy records (e.g., Kemp and Telford (2015)). In Table 1, we provide a list of all the possible options for  $f$  within the `reslr` package. Since some of the models we fit do not vary over space (they apply to a single site or treat a set of sites as identical), we use  $f(t)$  rather than  $f(\mathbf{x}, t)$  to denote the process model. In Figure 1, we provide a simplified graphic to describe the Bayesian hierarchical framework used for each statistical model in Table 1.

When using proxy RSL data, measurement error is also present in the input variable (time) due to the dating technique used. For the input measurements,  $\tilde{t}$  is assumed to be a noisy estimate of the true time value  $t$ :

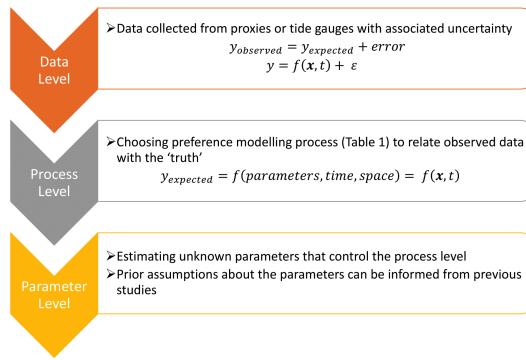
$$\tilde{t} = t + \epsilon_t \quad (2)$$

with the error term given by  $\epsilon_t \sim N(0, s_t^2)$  where  $s_t$  is the known measurement error associated with time.

We use two methods to account for the time measurement uncertainty. The first is the Errors-in-variables (EIV) method, which assumes that the input variable, e.g., time, is measured as an error-prone substitute and models it directly (Dey et al., 2000). The second uncertainty method is the Noisy Input (NI) method. This method fits an initial model and uses the derivative of the mean of  $f$  to calculate a corrective variance term. Then, the model is re-run with this additional corrective variance term, allowing for the input noise variation to be learned from the complete outputs of the model (McHutchon and Rasmussen, 2011). The `reslr` package employs the EIV method for linear regression, change point, and IGP models, while employing the NI uncertainty method for temporal spline, spatio-temporal spline, and GAM models. In general, the EIV method tends to be slower but models the uncertain input process directly, whilst the NI method is faster but requires the model to be fitted twice.

In Table 1, we present a range of statistical modelling techniques for  $f(\mathbf{x}, t)$ , the component of our approach, available in the `reslr` package. Below we discuss each technique and provide insight into the potential uses of these techniques for the paleo-environmental community.

Statistical Model	Model Information	model_type code
Errors in variables simple linear regression	A straight line of best fit taking into account any age and measurement errors in the RSL values using the method of Cahill et al. (2015a)	"eiv_slr_t"
Errors in variables change point model	An extension of the linear regression modelling process. It uses piece-wise linear sections and estimates where/when trend changes occur in the data (Cahill et al., 2015a)	"eiv_cp_t"
Errors in variables integrated Gaussian process	A non-linear fit that uses a Gaussian process prior on the rate of sea-level change that is then integrated (Cahill et al., 2015a).	"eiv_igp_t"
Noisy Input spline in time	A non-linear fit using regression splines (Upton et al., 2023).	"ni_spline_t"
Noisy Input spline in space and time	A non-linear fit for a set of sites across a region using the method of Upton et al. (2023).	"ni_spline_st"
Noisy Input Generalised Additive model for the decomposition of the RSL signal	A non-linear fit for a set of sites across a region and provides a decomposition of the signal into regional, local linear, and non-linear local components. This full model is as described in Upton et al. (2023).	"ni_gam_decomp"



**Figure 1:** A visual representation of the Bayesian hierarchical model structure used for each statistical model in the `reslr` package.

**Table 1:** List of all statistical models available in the `reslr` package. We provide a short description and the relevant literature for each model. The `model_type` code column represents the text input the user should use when implementing their preferred modelling technique.

### 2.3 EIV Linear Regression

The EIV linear regression model is given by:

$$f(t) = \alpha + \beta t \quad (3)$$

where  $\alpha$  is the intercept,  $\beta$  is the slope, and  $t$  is time. Earlier studies, for example, Shennan and Horton (2002) and Engelhart et al. (2009), employed linear regression when evaluating the rate of RSL change over the past 4,000 years. The `reslr` package implements a temporal linear regression as its simplicity is popular for approximate estimates of linear rates of RSL change. However, linearity assumptions for RSL change are often unrealistic when examining long-term historical trends.

### 2.4 EIV Change Point Model

The EIV change point (CP) model, an extension of the linear regression model, assumes the RSL process is piecewise linear and estimates when trend changes occur in the data (Cahill et al., 2015b). Mathematically, the multiple CP model,  $f(t)$  is described as:

$$f(t) = \begin{cases} \alpha_1 + \beta_j(t - \lambda_1) & \text{when } j = 1, 2, \\ \alpha_{j-1} + \beta_j(t - \lambda_{j-1}), & \text{when } j = 3, \dots, m+1 \end{cases} \quad (4)$$

where  $\alpha_j$  is the expected value of the response at the  $j$ th CP.  $\lambda_j$  is the time at which the CP occurs with the prior restriction that  $\lambda_1 < \lambda_2 < \dots < \lambda_m$  and  $m$  is the number of CPs (Cahill et al., 2015b). In the `reslr` package, the user can select  $m$  to be 1, 2, or 3 CPs.  $\beta_1$  and  $\beta_{m+1}$  are the slopes before and after the first and last CP, respectively.  $\beta_j$  for  $j = 2, \dots, m$  are the slopes between the  $(j-1)$ th and  $j$ th CP.

This technique has been used in different aspects of the sea-level literature. For example, Kemp et al. (2009) determined the magnitude and the timing of recent accelerated sea-level rise using change point models in North Carolina, USA. Brain et al. (2012) used the CP method to examine the impacts of sediment compaction on reconstructing recent sea-level rise in the United Kingdom. Hogarth et al. (2020) used CP models to obtain more consistent estimates of sea-level rise since 1958 for the British Isles. The main advantage of the CP model is its ability to identify sudden changes in RSL. However, the number of change points must be specified by the user.

### 2.5 Integrated Gaussian Process

An Integrated Gaussian process (IGP) is a modelling strategy that is extensively used by the sea-level community when examining the temporal evolution of sea level change (e.g. Cahill et al. (2015a); Hawkes et al. (2016); Kemp et al. (2017); Shaw et al. (2018); Dean et al. (2019); Stearns et al. (2023); Kirby et al. (2023)).

The IGP uses Gaussian Process (GP) to directly estimate the rate of change of the response

(Holsclaw et al., 2013). In order to extract the original  $f(t)$  we integrate  $p(t)$ :

$$f(t) = \alpha + \int_0^t p(u)du \quad (5)$$

where  $\alpha$  is the intercept, and  $p(t)$  is the rate of change,  $p(t) = \frac{df}{dt}$ , described as:

$$p(t) \sim GP(\mu(t), k(t, t')) \quad (6)$$

with  $t$  time and  $\mu(t)$  the mean function, and  $k(t, t')$  is the covariance function. The covariance function provides insight into the relationship between the outcome variables, i.e., if the input variables,  $t$  and  $t'$ , are in close proximity, the corresponding outcomes will be more correlated, and vice versa (Rasmussen and Williams, 2006). It is written as (Cahill et al., 2015a):

$$k(t, t') = \nu^2 \rho^{(t-t')^2} \quad (7)$$

where  $\rho$  is the correlation parameter and  $\nu^2$  is the variance of the rate process.

The technique, described by Cahill et al. (2015a), offers insights into examining rates of Relative Sea Level (RSL) change using proxy records from a single location. Apart from the IGP model, the `reslr` package does not rely on GP methods. We acknowledge that the use of GP modelling has gained considerable traction within the sea-level research community, particularly for investigating the spatio-temporal evolution of sea-level changes, as evidenced by notable studies (Kopp et al., 2009; Kopp, 2013; Kopp et al., 2016; Kemp et al., 2018; Walker et al., 2021). Nevertheless, in the context of the `reslr` package, we have intentionally opted for computationally efficient alternatives—splines and GAMs—as detailed in our prior work (Upton et al., 2023). These methods offer practical and effective approaches to analysing sea-level data, accommodating the complexities of spatio-temporal dynamics while ensuring computational tractability.

## 2.6 Temporal Spline

Splines are mathematical tools used in a wide range of settings from interpolation to data smoothing. There are a variety of different splines available, yet in this research we focus on B-splines (de Boor, 1978; Dierckx, 1995) and P-splines (Eilers and Marx, 1996). Mathematically, B-splines are described in the following way:

$$f(t) = \sum_{k=1}^K b_k(t) \beta_k \quad (8)$$

where  $b_k(t)$  is the spline basis function and  $\beta_k$  is the spline coefficient.

Following on from B-splines, Eilers and Marx (1996) describe a method to overcome the difficulty of choosing the correct number of knots by developing penalised spline or P-splines. Penalised differences in the spline coefficients control the smoothness of the spline based on differences (of order  $d$ ) of the spline coefficients. The first order differences are written as:

$$\Delta \beta_k = \beta_k - \beta_{k-1} \quad (9)$$

The spline coefficient will be centered on the previous value with an inverse smoothness parameter  $\sigma_\beta^2$ :

$$\Delta \beta_k \sim \mathcal{N}(0, \sigma_\beta^2) \quad (10)$$

In our package, P-splines are used for the NI spline in time and the extendable nature of these splines allows for different components to be examined within the GAM, which is described below.

## 2.7 Spatio-Temporal Spline

We use a spatio-temporal spline to examine RSL evolving over time at multiple locations. We include a tensor product to capture the variability over time and space (represented with longitude and latitude). For each individual covariate, time ( $t$ ) and longitude ( $x_1$ ) and latitude ( $x_2$ ), we construct a B-spline basis (Wood, 2017a). These basis functions are combined product-wise in the following way (Wood, 2006):

$$f(t, x_1, x_2) = \sum_{h=1}^H \sum_{i=1}^I \sum_{j=1}^J b_h(t) b_i(x_1) b_j(x_2) \beta_{hij} \quad (11)$$

where  $\beta_{hij}$  is the spline coefficient.  $H$  is the number of knots for  $b_h(t)$  the spline basis function in time  $t$ .  $I$  is the number of knots for  $b_i(x_1)$  the spline basis function for longitude.  $J$  is the number of knots for  $b_j(x_2)$  the spline basis functions for latitude values. The prior for the spline coefficient is given as:

$$\beta_{hij} \sim \mathcal{N}(0, \sigma_\beta^2) \quad (12)$$

where  $\sigma_\beta^2$  is the smoothness parameter for the spatio-temporal spline. The `reslr` package uses B-splines for the NI spline in space and time, allowing for multiple sites to be examined. The advantage of the tensor B-spline approach is that the basis functions are simple to construct, each depending on only one input variable. However, the number of parameters to estimate does increase considerably.

## 2.8 Generalised Additive Models

Generalised additive models are an extension of generalised linear models that use a basis expansion and a smoothing penalty to create linear predictors that are dependent on the sum of smooth functions of the predictor variable (GAMs; (Wood, 2017b)). The model developed by Upton et al. (2023) uses splines and random effects to create a spatio-temporal relative sea level surface. It identifies variations of sea-level at different spatial and temporal scales, encompassing multiple underlying processes and avoiding a focus on specific physical processes. The decomposition of this mean relative sea level surface can be written as:

$$f(\mathbf{x}, t) = r(t) + g(z_x) + h(z_x) + l(\mathbf{x}, t) \quad (13)$$

where  $r(t)$  is the regional component at time  $t$  represented with a spline in time.  $g(z_x)$  is the linear local component at location  $x$  represented by a random effect with  $z_x$  representing each data site.  $h(z_x)$  is the spatial vertical offset for each data site.  $l(\mathbf{x}, t)$  is the non-linear local component represented with a spline in space-time.

The regional component ( $r(t)$ ) represents temporal processes that are common to all locations, including barystatic and thermosteric contributions, where the former is caused by the transfer of mass between land-based ice and oceans (Gregory et al., 2019) and the latter is influenced by changes in global temperature creating density variations within the oceans (Grinsted, 2015). It is described using a spline in time:

$$r(t) = \sum_{s=1}^{k_r} b_{rs}(t) \beta_s^r \quad (14)$$

where  $\beta_s^r$  is the  $s^{th}$  spline coefficient,  $k_r$  is the number of knots, and  $b_{rs}(t)$  is the  $s^{th}$  spline basis function at time  $t$ . The prior for the spline coefficients of the regional component  $\beta_s^r$  are:

$$\beta_s^r \sim \mathcal{N}(0, \sigma_r^2) \quad (15)$$

where the smoothness of the model fit is controlled by  $\sigma_r$ , the standard deviation of the spline coefficient.

The linear local component ( $g(z_x)$ ) of the sea level model aims to capture linear trends present in the relative sea level signal. One such cause is glacial isostatic adjustment (GIA), which is a response of the Earth, the gravitational field, and the ocean to changes in the size of ice sheets (Whitehouse, 2018). On relatively short timescales, it is approximated to be linear through time with spatial variability along the Atlantic coastline of North America (Engelhart et al., 2009). Mathematically, we define our linear local component as an unstructured random effect for each site which is formulated as:

$$g(z_{x_j}) = \beta_j^g t \quad (16)$$

where  $\beta_j^g$  is a slope parameter specific for each site  $j$ . The prior for the linear local component is given by:

$$\beta_j^g \sim \mathcal{N}(m_{g_j}, s_{g_j}^2) \quad (17)$$

where  $m_{g_j}$  and  $s_{g_j}^2$  are the empirically estimated rate and associated variance for the dataset (refer to Upton et al., 2023, for a detailed description).

The site-specific vertical offset  $h$  is a random effect used to capture vertical shifts associated with measurement variability between sites and is formulated as:

$$h(z_{x_j}) = \beta_j^h \quad (18)$$

where  $\beta_j^h$  contains the random effect coefficients for site  $j$ . The prior for the site-specific vertical offset  $\beta_j^h$  is given as:

$$\beta_j^h \sim \mathcal{N}(0, \sigma_h^2) \quad (19)$$

where  $\sigma_h^2$  is the variance of the random intercept across all data sites.

The non-linear local component ( $l(\mathbf{x}, t)$ ) captures structured and unstructured RSL variability on century timescales, including dynamic sea-level changes (atmospheric and oceanic circulation patterns (Gregory et al., 2019)) and site-specific processes (e.g., sediment compaction affecting solid Earth's surface (Horton et al., 2018)). It is described using a spatio-temporal spline function formed using a tensor product and is formulated as:

$$l(\mathbf{x}, t) = \sum_{s=1}^{k_l} b_{l_s}(\mathbf{x}, t) \beta_s^l \quad (20)$$

where  $\beta_s^l$  is the  $s^{th}$  spline coefficient,  $k_l$  is the number of knots, and  $b_{l_s}(\mathbf{x}, t)$  is the  $s^{th}$  spline basis function at time  $t$  and location  $\mathbf{x}$ . The prior for the spline coefficient  $\beta_s^l$  is given as:

$$\beta_s^l \sim \mathcal{N}(0, \sigma_l^2) \quad (21)$$

where  $\sigma_l^2$  is the variance of the spline coefficients over space and time.

As described in Upton et al. (2023), B-splines are used for both the regional and local terms as this model structure balances both model usability and computational efficiency for examining proxy-based sea level reconstructions on a regional to local scale. B-splines also allow for easier prior elicitation of the smoothness parameters since they directly control the variability of the spline weights in the model.

### 3 Implementation

Within the package, we keep the number of functions to a minimum to ensure accessibility for users with varying experience in R. We run the statistical models using an MCMC algorithm and include a summary function to obtain a high-level insight into the outputs. We use S3 classes to access the summary, print, and plot commands. The package has functions to plot the input data and resulting model fits using ggplot2 (Wickham, 2016). The user has access to all the underlying information used to create these plots, allowing these visualisations to be re-created. In addition, the functions within the package are extendable, allowing advanced users access to more complex outputs. A detailed description for each function and associated commands for the `reslr` package can be found in the vignettes.

In this section, we provide insight into the example dataset and additional data sources using tide-gauge data within the `reslr` package. A discussion is provided into each function using two separate case studies: a single location and multiple locations. In the first case study, we demonstrate the Noisy Input temporal spline (`model_type = "ni_spline_t"`) which is an example modelling strategy for a single location. In the second case study, we examine multiple locations using the Noisy Input GAM decomposition (`model_type = "ni_gam_decomp"`).

#### 3.1 Example proxy dataset

We include an example dataset called `NAACproxydata`. The full dataset with the names of the locations and associated literature is in the Appendix. The `NAACproxydata` is a data frame with 1715 rows and 8 columns which include:

- *Region*: Region name
- *Site*: Site name
- *Latitude*: Latitude of the site
- *Longitude*: Longitude of the site
- *RSL*: Relative Sea level in metres
- *RSL\_err*: 1 standard deviation error associated with relative sea level measured in metres
- *Age*: Age in years Common Era (CE)
- *Age\_err*: 1 standard deviation error associated with the age in years CE

### 3.2 Including tide-gauge data

The tide-gauge data available to users can be obtained from the [PSMSL online database](#) (([Holgate et al., 2013](#); [Woodworth and Player, 2003](#))) through the `reslr` package. To ensure compatibility with the proxy records, several processing steps are performed within the package, as discussed earlier.

When incorporating tide-gauge data, users have three methods to select their preferred tide gauge(s). The first option is to provide a list of tide-gauge names from the PSMSL database, allowing users the freedom to select any tide gauge available. The second option is to automatically identify the nearest tide gauge to the proxy location that has more than 20 years of observations. This option proves particularly useful when examining proxy records and extending the temporal range to capture recent changes in RSL.

The final option enables the selection of all tide gauges within a 1-degree radius (latitude and longitude) of the proxy location, provided they have more than 20 years of observations. This option grants users access to a wide array of tide gauges within a larger geographic area. Moreover, users can combine the first option with either the second or the third option, allowing for the freedom to choose specific tide gauges while incorporating the nearest tide gauge or multiple tide gauges.

All the values mentioned in this paragraph are arguments that can be adjusted within the function, giving users flexibility in customizing their data selection process according to their specific requirements.

### 3.3 Case Study for 1 location

In the following sections, we use one site, Cedar Island, North Carolina, USA, from the example dataset, `NAACproxydata`:

```
CedarIslandNC <- reslr::NAACproxydata %>%
  dplyr::filter(Site == "Cedar Island")

glimpse(CedarIslandNC)

#> Rows: 104
#> Columns: 8
#> $ Region    <chr> "North Carolina", "North Carolina", "North Carolina", "North~
#> $ Site       <chr> "Cedar Island", "Cedar Island", "Cedar Island", "Cedar Islan~
#> $ Latitude   <dbl> 34.971, 34.971, 34.971, 34.971, 34.971, 34.971, 34.971, 34.9~
#> $ Longitude  <dbl> -76.38, -76.38, -76.38, -76.38, -76.38, -76.38, -76.38, -76.~
#> $ RSL        <dbl> -0.12, -0.14, -0.16, -0.18, -0.19, -0.21, -0.22, -0.23, -0.2~
#> $ Age         <dbl> 2005, 1996, 1988, 1979, 1974, 1963, 1957, 1951, 1941, 1937, ~
#> $ Age_err    <dbl> 2.25, 2.00, 5.00, 5.75, 5.50, 5.50, 7.00, 7.75, 7.75, 8.00, ~
#> $ RSL_err    <dbl> 0.06, 0.06, 0.06, 0.06, 0.06, 0.06, 0.06, 0.06, 0.06, 0.06, ~
```

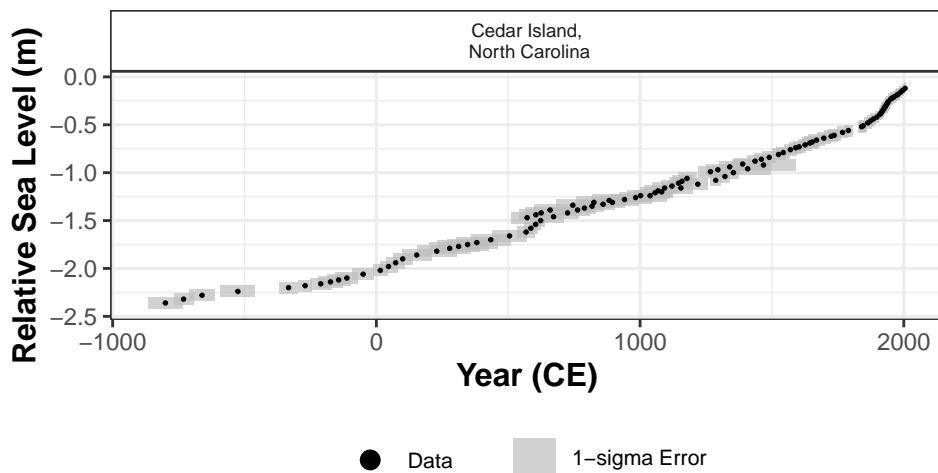
For a single location such as this case study, we recommend using an EIV IGP or a NI spline in time as they estimate RSL changes over time. In this example, tide-gauge data is not included but it is an option available to the user if they require. The next example will demonstrate a more complex analysis with the inclusion of tide gauges.

After selecting the data site from the example dataset, we use the `reslr_load` function to process the data prior to running the statistical model, and it has a number of different settings that the user can alter depending on the model choice. One such setting is the `prediction_grid_res` option. This provides the resolution at which predictions of RSL and RSL rates are made and subsequently plotted. We set the default at 50 years, and if a finer grid is required, the user can alter the setting for `prediction_grid_res`. The `reslr_load` function includes additional settings to include tide-gauge data and linear rates, which will be discussed in the next case study. For the single site case study, we demonstrate the `reslr_load` function:

```
CedarIslandNC_input <- reslr_load(data = CedarIslandNC)
```

The output of this function is a list of two data frames called `data` and `data_grid`. The `data` data frame is the inputted data with an additional column called `data_type_id` which distinguishes proxy records from tide-gauge data. The `data_grid` is a data frame that is evenly spaced in time based on the `prediction_grid_res` value chosen by the user and is used to create the plots. A brief insight into the outputs of the `reslr_input` function can be obtained using the `print` function, which provides the number of observations and the sources of the data, as shown below:

```
print(CedarIslandNC_input)
```



**Figure 2:** A plot of the raw data for our example site Cedar Island North Carolina. The x-axis is time in years in the Common Era (CE) and the y-axis is relative sea level in metres. The grey boxes are 1 standard deviation vertical and horizontal (temporal) uncertainty. The black dots are the midpoints of the uncertainty boxes.

```
#> This is a valid reslr input object with 104 observations and 1 site(s).
#> There are 1 proxy site(s) and 0 tide gauge site(s).
#> The age units are; Common Era.
#> Decadally averaged tide gauge data was not included. It is recommended for the ni_gam_decomp model
#> The linear_rate or linear_rate_err was not included. It is required for the ni_gam_decomp model
```

The next step is using the `plot` function to plot the raw data, shown in Figure 2, using the following:

```
plot(CedarIslandNC_input, plot_caption = FALSE)
```

### Noisy Input Spline in time

The NI spline in time ("ni\_spline\_t") examines how the response variable, RSL, varies in time. While the EIV-IGP method is commonly used in the sea-level community, we demonstrate that the NI spline in time is a faster alternative. Unlike the Gaussian process, which has a computational complexity that grows exponentially with the number of data points, the spline equivalent uses pre-computed basis functions, resulting in a more efficient computation (Wood, 2017b).

For this model type, we use the `reslr_mcmc` function to implement the MCMC simulation using JAGS, and the model type setting is selected to be `model_type = "ni_spline_t"`.

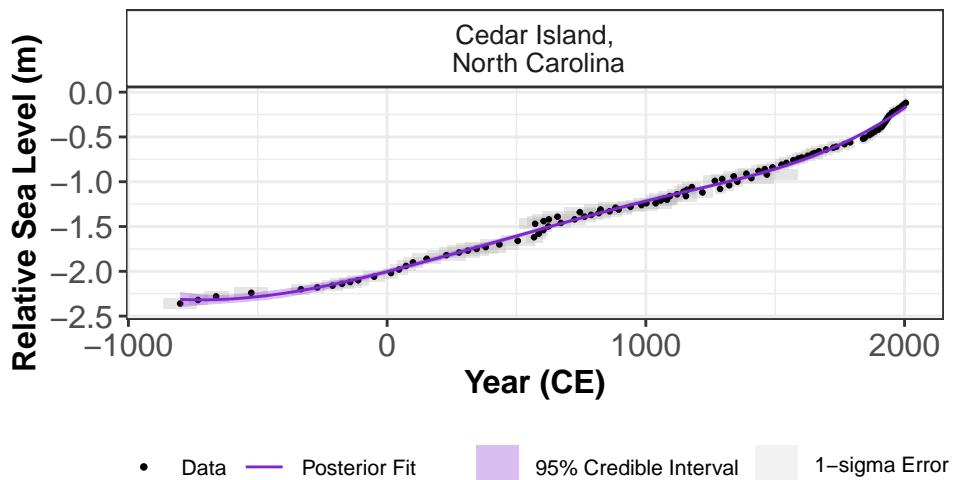
```
res_ni_spline_t <- reslr_mcmc(input_data = CedarIslandNC_input,
                                model_type = "ni_spline_t")
```

The output of the `reslr_mcmc` function is a list that stores the JAGS model run, the input dataframe, and the dataframes for plotting the results. The user can set the size of the credible intervals by changing the `CI` setting in this function; the current default is `CI = 0.95`. In addition, the user can alter the number of iterations which will be required if the model is not converging.

To obtain a brief insight into the outputs of the `reslr_mcmc` function, the user can use the `print` function which provides the number of iterations and the model type:

```
print(res_ni_spline_t)

#> This is a valid reslr output object with 104 observations and 1 site(s).
#> There are 1 proxy site(s) and 0 tide gauge site(s).
#> The age units are; Common Era.
#> The model used was the Noisy Input Spline in time model.
#> The input data has been run via reslr_mcmc and has produced 3000 iterations over 3 MCMC chains.
```



**Figure 3:** The plot of the Noisy Input spline in time model fit for our example site, Cedar Island, North Carolina. The x-axis is time in years in the Common Era (CE) and the y-axis is relative sea level in metres. The grey boxes are 1 standard deviation vertical and horizontal (temporal) uncertainty. The black dots are the midpoints of the uncertainty boxes. The solid purple line represents the posterior model fit with a 95% credible interval denoted by shading.

The convergence of the MCMC algorithm can be examined for the “ni\_spline\_t” model using the summary function and ensures the scale reduction factor ( $R\text{-hat}$ ) is close to 1 (Gelman and Rubin, 1992; Gelman et al., 2013). If the model run has converged, the package will print: “No convergence issues detected”. If the package prints: “Convergence issues detected, a longer run is necessary”. The user is recommended to update the reslr\_mcmc function with additional iterations as described above. The summary function provides insight into the parameter estimates from the model using the following:

```
summary(res_ni_spline_t)

#> No convergence issues detected.

#> # A tibble: 2 x 7
#>   variable     mean      sd      mad      q5      q95 rhat
#>   <chr>     <dbl>    <dbl>    <dbl>    <dbl>    <dbl> <dbl>
#> 1 sigma_beta 2.04    0.617   0.520   1.29    3.24    1.00
#> 2 sigma_y    0.00633 0.00478 0.00456 0.000513 0.0153  1.00
```

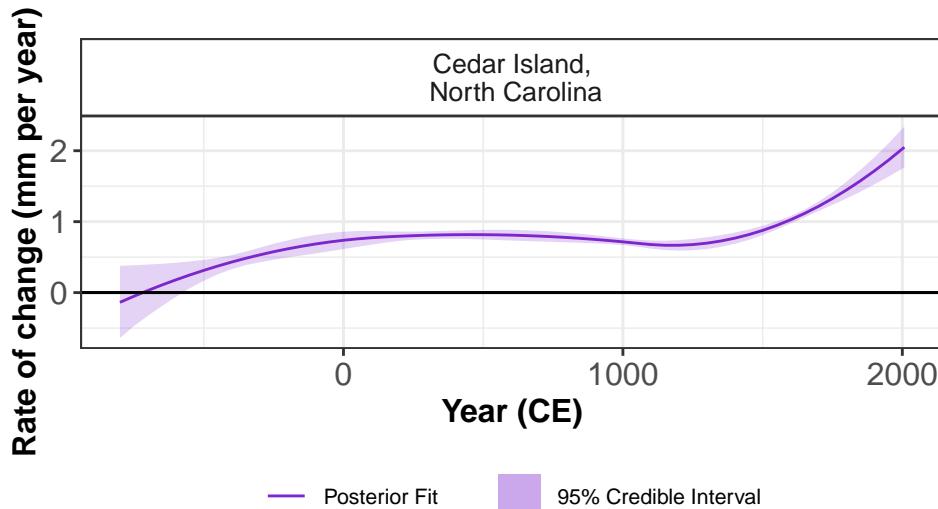
For the parameter estimates, “sigma\_beta” acts as a smoothness parameter controlling the penalisation of the spline coefficients for the spline in time model, and “sigma\_y” represents the data model variation. These are  $\sigma_\beta$  and  $\sigma_y$  as described in Section 2.

The final results from the “ni\_spline\_t” model can be illustrated using the plot function, and the corresponding dataframes are stored in the res\_ni\_spline\_t object called output\_dataframes as a named list element. Figure 3 demonstrates the posterior model fit for our example site using:

```
plot(res_ni_spline_t,
     plot_type = "model_fit_plot",
     plot_caption = FALSE)
```

In Figure 4, the rate of change of this posterior model fit is presented and can be viewed using:

```
plot(res_ni_spline_t,
     plot_type = "rate_plot",
     plot_caption = FALSE)
```



**Figure 4:** The rate of change model fit using the Noisy Input spline in time model for our example site, Cedar Island, North Carolina. The rate is calculated by taking the first derivative of the total model fit. The x-axis is time in years in the Common Era (CE) and the y-axis is the instantaneous rate of change of sea level in mm per year. The solid purple line represents the posterior model fit with a 95% credible interval denoted by shading. There is a black horizontal line which is the zero rate of change for this site.

### 3.4 Case Study for multiple sites

The sea-level research community is commonly interested in temporal and spatial variations in RSL. To cater to this interest, the `reslr` package offers two models for spatio-temporal modeling. The first model is a noisy-input spline that accounts for noise in both time and space, providing a robust representation of RSL dynamics. The second model, a more intricate option, is the Noisy Input GAM. Mathematical details concerning the Noisy Input GAM can be found in Upton et al. (2023). In our upcoming example, we will focus on this model as it empowers users to explore the decomposition of the RSL signal over time and space, unraveling valuable insights into the underlying dynamics.

#### Noisy Input Generalised Additive Model for decomposition of response signal

We demonstrate the functions settings required for the NI GAM. This model requires an adequate number of sites to perform the decomposition and the minimum sites required will depend on the signal in the data. In this example, we use nine sites from the example dataset, `NAACproxydata`, which are selected in the following manner:

```
multi_site <- reslr::NAACproxydata %>%
  dplyr::filter(Site %in% c("Cedar Island", "Nassau",
                           "East River Marsh", "Swan Key",
                           "Placentia",
                           "Pelham Bay", "Fox Hill Marsh",
                           "Snipe Key", "Big River Marsh"))
```

Next, the `reslr_load` function is required for the preparation of input data for the NI GAM, which necessitates additional information not required by earlier models. Firstly, the statistical model relies on an estimate of the “linear local rate” and its associated uncertainty. By setting `include_linear_rate = TRUE`, the package incorporates this rate, which is assumed to stem from physical processes like Glacial Isostatic Adjustment (GIA). Users have the flexibility to include their preferred linear rate values as additional columns (`linear_rate` and `linear_rate_err`) in the input data frame. If these values are not provided, the package automatically calculates them using the available data.

Secondly, users are encouraged to include tide-gauge data by setting `include_tide_gauge = TRUE`. As discussed previously, users need to make a decision regarding the inclusion

of the closest tide gauge (`TG_minimum_dist_proxy = TRUE`), selecting specific tide gauges by providing a list of names (`list_preferred_TGs = c("ARGENTIA")`), or including all tide gauges within a one-degree proximity of the proxy site (`all_TG_1deg = TRUE`). Additionally, the tide-gauge data requires values for the `linear_rate` and `linear_rate_err` columns, which are calculated using the ICE-5G (VM2) Model (Peltier, 2004) with an uncertainty value of 0.3 mm/year (Engelhart et al., 2009), both provided within the `reslr` package.

Thirdly, the tide-gauge data is averaged over a decade to align with the resolution of proxy records. If necessary, users can adjust the size of the averaging window to accommodate varying sediment accumulation rates. For example, a longer sediment accumulation rate would result in a larger average, such as 20 years. The default setting for `sediment_average_TG` is 10 years, which we will use in our example.

The final setting of the `reslr_load` function is `prediction_grid_res`, allowing users to modify the resolution of the output plots. The default setting of 50 years serves as a starting point, but users have the flexibility to explore alternative options. For our example, we will utilize nine proxy sites and select all tide-gauges within a one-degree range of our proxy site, maximizing the number of data points to demonstrate the capabilities of our package. The specific settings employed are described below:

```
multi_site_input <- reslr_load(
  data = multi_site,
  include_tide_gauge = TRUE,
  include_linear_rate = TRUE,
  TG_minimum_dist_proxy = TRUE,
  all_TG_1deg = TRUE)
```

Similar to the previous example, the output of this function is a list of two dataframes called `data` and `data_grid`. The `data` dataframe is the inputted data with additional columns for the `data_type_id` which will contain “ProxyRecord” and “TideGaugeData”. The `data_grid` is a dataframe that is evenly spaced in time based on the `prediction_grid_res` value chosen by the user and is used to create the plots. In this example, we have 9 proxy sites and 26 tide gauges sites and the `print` provides insights into the outputs of the `reslr_input` function yields insights such as the number of observations (1,130).

```
print(multi_site_input)

#> This is a valid reslr input object with 1130 observations and 35 site(s).
#> There are 9 proxy site(s) and 26 tide gauge site(s).
#> The age units are; Common Era.
#> Decadally averaged tide gauge data included by the package.
#> The linear_rate and linear_rate_err has been included.
```

A plot of the raw data can be created using the `plot` function, with an option to plot the tide gauges and the proxy records together or have separate plots for each data source. Figure 5 demonstrates the resulting plot for the proxy records only using the following function:

```
plot(x = multi_site_input,
      plot_proxy_records = TRUE,
      plot_tide_gauges = FALSE,
      plot_caption = FALSE)
```

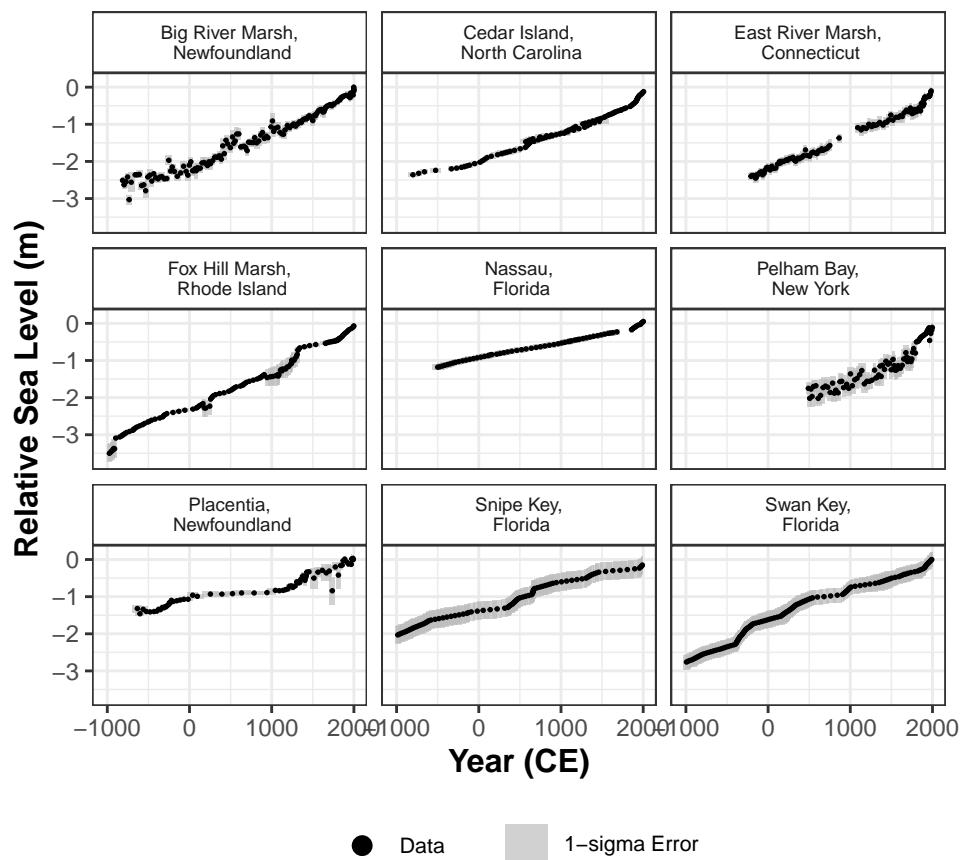
For this model type, the `reslr_mcmc` function should specify the `model_type = "ni_gam_decomp"`, and the MCMC simulation settings can be altered to ensure convergence.

```
res_ni_gam_decomp <- reslr_mcmc(
  input_data = multi_site_input,
  model_type = "ni_gam_decomp"
)
```

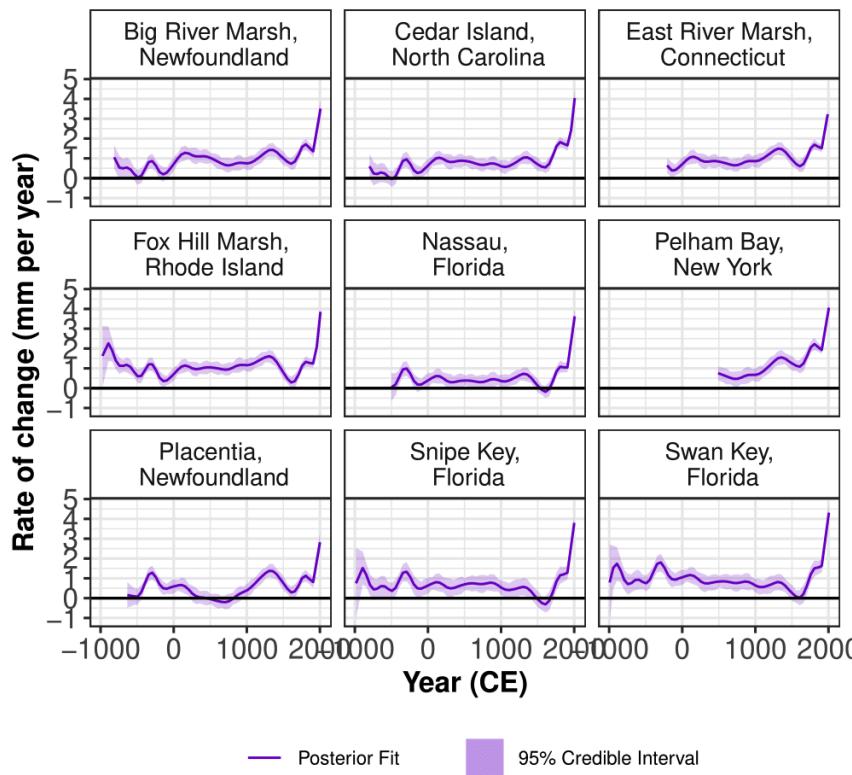
The output of the `reslr_mcmc` function is a list that stores the JAGS model run, the input dataframe, and the dataframes for plotting the results. Identical to the other model processes, the convergence of the MCMC algorithm is examined and the parameter estimates from the model can be investigated using the following:

```
summary(res_ni_gam_decomp)

#> No convergence issues detected.
#> # A tibble: 4 × 7
#>   variable     mean     sd     mad     q5     q95    rhat
```



**Figure 5:** A plot of the raw data for our nine example sites along the Atlantic coast of North America. The x-axis is time in years in the Common Era (CE) and the y-axis is relative sea level in metres. The grey boxes are 1 standard deviation vertical and horizontal (temporal) uncertainty. The black dots are the midpoints of the uncertainty boxes.



**Figure 6:** The rate of change for the total model fit for the Noisy Input generalised additive model for sites along the Atlantic coast of North America. It is calculated by finding the derivative of the total model fit. The solid purple line is the mean rate of change fit and the shading denotes a 95% credible interval for each site along the Atlantic coast of North America. The x-axis is time in years in the Common Era (CE), and the y-axis is the rate of change in mm per year.

```
#>   <chr>      <num>    <num>    <num>    <num>    <num>    <num>
#> 1 sigma_beta_h 1.82    0.257    0.247    1.45    2.29    1.00
#> 2 sigma_beta_r 0.288   0.0547   0.0505   0.213   0.391   1.00
#> 3 sigma_beta_l 0.939   0.148    0.145    0.717   1.21    1.00
#> 4 sigma_y     0.0155  0.00107  0.00108  0.0138  0.0173  1.00
```

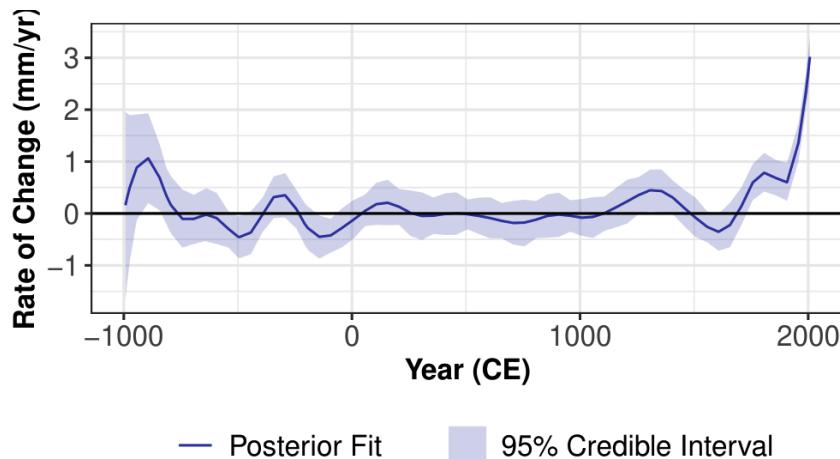
For the parameter estimates, we provide the standard deviation associated with each component of the NI GAM decomposition. Specifically, “sigma\_beta\_r” represents the standard deviation of the spline coefficient for the regional component, “sigma\_beta\_l” represents the standard deviation of the spline coefficient for the non-linear local component, “sigma\_beta\_h” denotes the standard deviation of the site-specific vertical offset component, and “sigma\_y” indicates the data model variation. These names correspond to the algebraic components described in Section 2 above.

One of the key advantages of the NI GAM approach is its ability to decompose regional RSL change into separate components. The results from the `ni_gam_decomp` model can be visualized using the `plot` function, which generates individual plots for each component. Additionally, all components, except for the linear local component, have corresponding rate plots. Users can access the data used to create each plot in the `res_ni_gam_decomp` object as separate dataframes for each component.

In our example, we demonstrate the rate of change for the total model fit in Figure 6. This figure illustrates the rate of change at each site, which is useful for understanding the variations of the relative sea-level signal, i.e.,  $f(\mathbf{x}, t)$ . To plot the rate of change, users can employ the following method:

```
plot(res_ni_gam_decomp,
  plot_type = "rate_plot",
  plot_caption = FALSE)
```

The regional component ( $r(t)$ ) captures the mean of RSL change along the Atlantic coast of North America. The associated rate of change of the regional component, as seen in



**Figure 7:** The rate of change for the regional component of the Noisy Input generalised additive model for the nine proxy sites and the eleven tide gauges along the Atlantic coast of North America. It is calculated by finding the derivative of the regional component fit. The solid blue line is the mean rate of change fit, and the shading denotes a 95% credible interval. The x-axis is time in years in the Common Era (CE), and the y-axis is the rate of change in mm per year.

Figure 7, provides an important visual insight into the rate at which this trend varied over the past 3,000 years. It is accessed by:

```
plot(res_ni_gam_decomp,
     plot_type = "regional_rate_plot",
     plot_caption = FALSE)
```

## 4 Summary

In this paper, we have presented an overview of the `reslr` package and discussed its various features and design decisions. Our goal was to address the specific needs of the paleo sea-level community and provide an efficient and flexible R package that caters to different types of source data, whilst maintaining a simple workflow that does not require the user to learn too many different functions.

Through two case studies, we demonstrated the simplicity and accessibility of the package. The first case study examined a single site using the NI spline in time. Our results showed that the `reslr` package can provide RSL estimates and associated rate of change values over time for a single location. In the second case study, we showcased the capabilities of the `reslr` package when analysing data from multiple locations. We highlighted its flexibility, allowing for the decomposition of the relative sea-level signal into different components. Additionally, we presented a comprehensive method for incorporating tide-gauge data, which can help to provide valuable insights into recent changes in RSL not captured by proxy records.

There are several potential extensions for the `reslr` package. A potential enhancement involves replacing our current MCMC algorithm, employed via JAGS software (Plummer, 2003), with more efficient alternatives like Integrated Nested Laplace Approximations (INLA: Rue et al., 2009). Another improvement could be the integration of other instrumental data sources, such as satellite data, enabling the examination of other variables related to climate change. Overall, the `reslr` package offers a powerful toolkit for the paleo sea-level community, and we anticipate that it will continue to evolve and expand its capability to meet the evolving needs of researchers in this field.

## 5 Acknowledgements

Upton's work is supported by A4 (Aigéin, Aeráid, agus athrú Atlantaigh); the project is funded by the Marine Institute (grant: PBA/CC/18/01). Parnell's work is supported by the

SFI awards 17/CDA/4695; 16/IA/4520; 12/RC/2289P2. Cahill's research is conducted with the financial support of Science Foundation Ireland and co-funded by Geological Survey Ireland under Grant number 20/FFP-P/8610.

## 6 Appendix

### 6.1 Example dataset

The `reslr` package contains a dataset used as an example called `NAACproxydata`. This dataset contains proxy records from the Atlantic coast of North America as used in Upton et al. (2023). The 21 different proxy data sites and the references for each data source can be found in Table 3. The `reslr` package contains a dataset used as an example called `NAACproxydata`. This dataset contains proxy records from the Atlantic coast of North America as used in Upton et al. (2023). The 21 different proxy data sites and the references for each data source can be found in Table 3.

Site Name	Reference
Barn Island, Connecticut	(Donnelly et al., 2004; Gehrels et al., 2020)
Big River Marsh, Newfoundland	(Kemp et al., 2018)
Cape May Courthouse, New Jersey	(Kemp et al., 2013; Cahill et al., 2016)
Cedar Island, North Carolina	(Kemp et al., 2011, 2017)
Cheesequake, New Jersey	(Walker et al., 2021)
Chezzetcook Inlet, Nova Scotia	(Gehrels et al., 2020)
East River Marsh, Connecticut	(Kemp et al., 2015; Stearns et al., 2023)
Fox Hill Marsh, Rhode Island	(Stearns et al., 2023)
Leeds Point, New Jersey	(Kemp et al., 2013; Cahill et al., 2016)
Les Sillons, Magdalen Islands	(Barnett et al., 2017)
Little Manatee River, Florida	(Gerlach et al., 2017)
Nassau, Florida	(Kemp et al., 2014)
Pelham Bay, New York	(Kemp et al., 2017; Stearns et al., 2023)
Placentia, Newfoundland	(Kemp et al., 2018)
Revere, Massachusetts	(Donnelly, 2006)
Saint Simeon, Quebec	(Barnett et al., 2017)
Sanborn Cove, Maine	(Gehrels et al., 2020)
Sand Point, North Carolina	(Kemp et al., 2011, 2017)
Snipe Key, Florida	(Khan et al., 2022)
Swan Key, Florida	(Khan et al., 2022)
Wood Island, Massachusetts	(Kemp et al., 2011)

**Table 3:** Presents the names of all the sites available in the example dataset within the `reslr` package. For each site, we include the reference in the literature to the source of the data.

## References

- T. Aarup, M. Merrifield, B. Pérez Gómez, I. Vassie, and P. Woodworth. Manual on Sea-level Measurements and Interpretation, Volume IV : An update to 2006. *Intergovernmental Oceanographic Commission of UNESCO*, 4, 2006. URL [https://psmsl.org/train\\_and\\_info/training/manuals](https://psmsl.org/train_and_info/training/manuals). [p62]
- E. L. Ashe, N. Cahill, C. Hay, N. S. Khan, A. Kemp, S. E. Engelhart, B. P. Horton, A. C. Parnell, and R. E. Kopp. Statistical modeling of rates and trends in Holocene relative sea level. *Quaternary Science Reviews*, 204:58–77, 2019. doi: <https://doi.org/10.1016/j.quascirev.2018.10.032>. [p61]
- R. L. Barnett, P. Bernatchez, M. Garneau, and M.-N. Juneau. Reconstructing late Holocene relative sea-level changes at the Magdalen Islands (Gulf of St. Lawrence, Canada) using

- multi-proxy analyses. *Journal of Quaternary Science*, 32(3):380–395, 2017. doi: <https://doi.org/10.1002/jqs.2931>. [p76]
- M. J. Brain, A. J. Long, S. A. Woodroffe, D. N. Petley, D. G. Milledge, and A. C. Parnell. Modelling the effects of sediment compaction on salt marsh reconstructions of recent sea-level rise. *Earth and Planetary Science Letters*, 345–348:180–193, 2012. ISSN 0012-821X. doi: <https://doi.org/10.1016/j.epsl.2012.06.045>. [p64]
- N. Cahill, A. C. Kemp, B. P. Horton, and A. C. Parnell. Modeling sea-level change using Errors-in-Variables integrated Gaussian Process 1. *The Annals of Applied Statistics*, 9(2): 547–571, 2015a. doi: <https://doi.org/10.1214/15-AOAS824>. [p61, 63, 64, 65]
- N. Cahill, S. Rahmstorf, and A. C. Parnell. Change points of global temperature. *Environmental research letters*, 10(8):84002, 2015b. doi: <https://doi.org/10.1088/1748-9326/10/8/084002>. [p61, 64]
- N. Cahill, A. C. Kemp, B. P. Horton, and A. C. Parnell. A Bayesian hierarchical model for reconstructing relative sea level: from raw data to rates of change. *Climate of the Past*, 12 (2):525–542, 2016. doi: <https://doi.org/10.5194/cp-12-525-2016>. [p76]
- B. Carpenter, A. Gelman, M. D. Hoffman, D. Lee, B. Goodrich, M. Betancourt, M. A. Brubaker, J. Guo, P. Li, and A. Riddell. Stan: A probabilistic programming language. *Journal of statistical software*, 76, 2017. doi: 10.18637/jss.v076.i01. URL <https://www.jstatsoft.org/index.php/jss/article/view/v076i01>. [p63]
- C. de Boor. A Practical Guide to Spline. *Applied Mathematical Sciences*, New York: Springer, 1978, 27, 01 1978. doi: 10.2307/2006241. [p61, 65]
- S. Dean, B. P. Horton, N. Evelydou, N. Cahill, G. Spada, and D. Sivan. Can we detect centennial sea-level variations over the last three thousand years in Israeli archaeological records? *Quaternary Science Reviews*, 210:125–135, 2019. ISSN 0277-3791. doi: <https://doi.org/10.1016/j.quascirev.2019.02.021>. URL <https://www.sciencedirect.com/science/article/pii/S0277379117310831>. [p64]
- D. K. Dey, S. K. Ghosh, and B. K. Mallick. *Generalized linear models: a Bayesian perspective*. CRC Press, 2000. doi: <https://doi.org/10.1201/9781482293456>. [p61, 63]
- P. Dierckx. *Curve and surface fitting with splines*. Oxford University Press, 1995. URL <https://books.google.ie/books?id=-RIQ3SR0sZMC>. [p65]
- J. P. Donnelly. A Revised Late Holocene Sea-Level Record for Northern Massachusetts, USA. *Journal of Coastal Research*, 22(5):1051–1061, 2006. doi: <https://doi.org/10.2112/04-0207.1>. [p76]
- J. P. Donnelly, P. Cleary, P. Newby, and R. Ettinger. Coupling instrumental and geological records of sea-level change: Evidence from southern New England of an increase in the rate of sea-level rise in the late 19th century. *Geophysical Research Letters*, 31(5), 2004. doi: <https://doi.org/10.1029/2003gl018933>. [p76]
- R. Edwards and A. Wright. *Foraminifera: Handbook of Sea-Level Research*, chapter 13, pages 191–217. John Wiley & Sons, Ltd, 2015. ISBN 9781118452547. doi: <https://doi.org/10.1002/9781118452547.ch13>. [p62]
- P. Eilers and B. Marx. Flexible Smoothing with B-splines and Penalties. *Statistical Science*, 11, 1996. doi: <https://doi.org/10.1214/ss/1038425655>. [p65]
- S. E. Engelhart, B. P. Horton, B. C. Douglas, W. R. Peltier, and T. E. Törnqvist. Spatial variability of late Holocene and 20th century sea-level rise along the Atlantic coast of the United States. *Geology*, 37(12):1115–1118, 12 2009. ISSN 0091-7613. doi: 10.1130/G30360A.1. [p64, 66, 72]
- R. W. Gehrels. Determining Relative Sea-Level Change from Salt-Marsh Foraminifera and Plant Zones on the Coast of Maine, U.S.A. *Journal of Coastal Research*, 10(4):990–1009, 1994. URL <https://www.jstor.org/stable/4298291>. [p61, 62]

- W. Gehrels, S. Dangendorf, N. L. M. Barlow, M. H. Saher, A. J. Long, P. L. Woodworth, C. G. Piecuch, and K. Berk. A Preindustrial Sea Level Rise Hotspot Along the Atlantic Coast of North America. *Geophysical Research Letters*, 47(4), 2020. doi: <https://doi.org/10.1029/2019GL085814>. [p76]
- A. Gelman and D. B. Rubin. Inference from iterative simulation using multiple sequences. *Statistical Science*, 7(4):457–472, 1992. ISSN 08834237. URL <http://www.jstor.org/stable/2246093>. [p70]
- A. Gelman, J. B. Carlin, H. S. Stern, D. B. Dunson, A. Vehtari, and D. B. Rubin. *Bayesian data analysis*. CRC press, 2013. doi: <https://doi.org/10.1201/b16018>. [p70]
- M. J. Gerlach, S. E. Engelhart, A. C. Kemp, R. P. Moyer, J. M. Smoak, C. E. Bernhardt, and N. Cahill. Reconstructing Common Era relative sea-level change on the Gulf Coast of Florida. *Marine Geology*, 390:254–269, 2017. ISSN 0025-3227. doi: <https://doi.org/10.1016/j.margeo.2017.07.001>. [p76]
- V. Gornitz. *Paleoclimate Proxies, An Introduction*, pages 716–721. Springer Netherlands, 2009. doi: <https://doi.org/10.1007/978-1-4020-4411-3>. [p62]
- J. M. Gregory, S. M. Griffies, C. W. Hughes, J. A. Lowe, J. A. Church, I. Fukimori, N. Gomez, R. E. Kopp, F. Landerer, G. L. Cozannet, R. M. Ponte, D. Stammer, M. E. Tamisiea, and R. S. van de Wal. Concepts and Terminology for Sea Level: Mean, Variability and Change, Both Local and Global. *Surveys in Geophysics*, 40(6):1251–1289, 11 2019. ISSN 15730956. URL <https://link.springer.com/article/10.1007/s10712-019-09525-z>. [p66, 67]
- A. Grinsted. Projected Change—Sea Level. *Second assessment of climate change for Baltic Sea basin*, pages 253–263, 2015. URL [https://link.springer.com/chapter/10.1007/978-3-319-16006-1\\_14](https://link.springer.com/chapter/10.1007/978-3-319-16006-1_14). [p66]
- A. D. Hawkes, A. C. Kemp, J. P. Donnelly, B. P. Horton, W. R. Peltier, N. Cahill, D. F. Hill, E. Ashe, and C. R. Alexander. Relative sea-level change in northeastern Florida (USA) during the last 8.0ka. *Quaternary Science Reviews*, 142:90–101, 2016. ISSN 0277-3791. doi: <https://doi.org/10.1016/j.quascirev.2016.04.016>. URL <https://www.sciencedirect.com/science/article/pii/S0277379116301275>. [p64]
- P. Hogarth, C. Hughes, S. Williams, and C. Wilson. Improved and extended tide gauge records for the british isles leading to more consistent estimates of sea level rise and acceleration since 1958. *Progress in Oceanography*, 184:102333, 2020. ISSN 0079-6611. doi: <https://doi.org/10.1016/j.pocean.2020.102333>. URL <https://www.sciencedirect.com/science/article/pii/S0079661120300720>. [p64]
- S. J. Holgate, A. Matthews, P. L. Woodworth, L. J. Rickards, M. E. Tamisiea, E. Bradshaw, P. R. Foden, K. M. Gordon, S. Jevrejeva, and J. Pugh. New Data Systems and Products at the Permanent Service for Mean Sea Level. *Journal of Coastal Research*, 29(3):493 – 504, 2013. doi: 10.2112/JCOASTRES-D-12-00175.1. [p61, 62, 68]
- T. Holsclaw, B. Sansó, H. K. H. L. Lee, K. Heitmann, S. Habib, D. Higdon, and U. Alam. Gaussian process modeling of derivative curves. *Technometrics*, 55(1):57–67, 2013. ISSN 00401706, 15372723. URL <http://www.jstor.org/stable/24587388>. [p65]
- B. P. Horton, R. E. Kopp, A. J. Garner, C. C. Hay, N. S. Khan, K. Roy, and T. A. Shaw. Mapping Sea-Level Change in Time, Space, and Probability. *Annual Review of Environment and Resources*, 43(1):481–521, 2018. ISSN 1543-5938. doi: <https://doi.org/10.1146/annurev-environ-102017-025826>. [p67]
- A. C. Kemp and R. J. Telford. *Transfer functions: Handbook of Sea-Level Research*, chapter 31, pages 470–499. John Wiley & Sons, Ltd, 2015. ISBN 9781118452547. doi: <https://doi.org/10.1002/9781118452547.ch31>. [p62, 63]
- A. C. Kemp, B. P. Horton, S. J. Culver, D. R. Corbett, O. van de Plassche, W. R. Gehrels, B. C. Douglas, and A. C. Parnell. Timing and magnitude of recent accelerated sea-level rise (North Carolina, United States). *Geology*, 37(11):1035–1038, 11 2009. ISSN 0091-7613. doi: <https://doi.org/10.1130/G30352A.1>. [p64]

- A. C. Kemp, B. P. Horton, J. P. Donnelly, M. E. Mann, M. Vermeer, and S. Rahmstorf. Climate related sea-level variations over the past two millennia. *Proceedings of the National Academy of Sciences*, 108(27):11017–11022, 2011. doi: 10.1073/pnas.1015619108. [p76]
- A. C. Kemp, B. P. Horton, C. H. Vane, C. E. Bernhardt, D. R. Corbett, S. E. Engelhart, S. C. Anisfeld, A. C. Parnell, and N. Cahill. Sea level change during the last 2500 years in New Jersey, USA. *Quaternary Science Reviews*, 81:90–104, 2013. ISSN 02773791. doi: 10.1016/j.quascirev.2013.09.024. [p61, 76]
- A. C. Kemp, C. E. Bernhardt, B. P. Horton, R. E. Kopp, C. H. Vane, W. R. Peltier, A. D. Hawkes, J. P. Donnelly, A. C. Parnell, and N. Cahill. Late Holocene sea- and land-level change on the U.S. southeastern Atlantic coast. *Marine Geology*, 357:90–100, 2014. ISSN 0025-3227. doi: <https://doi.org/10.1016/j.margeo.2014.07.010>. [p76]
- A. C. Kemp, A. D. Hawkes, J. P. Donnelly, C. H. Vane, B. P. Horton, T. D. Hill, S. C. Anisfeld, A. C. Parnell, and N. Cahill. Relative sea level change in Connecticut (USA) during the last 2200 yrs. *Earth and Planetary Science Letters*, 428:217–229, 2015. ISSN 0012-821X. doi: <https://doi.org/10.1016/j.epsl.2015.07.034>. [p76]
- A. C. Kemp, T. D. Hill, C. H. Vane, N. Cahill, P. M. Orton, S. A. Talke, A. C. Parnell, K. Sanborn, and E. K. Hartig. Relative sea-level trends in New York City during the past 1500 years. *The Holocene*, 27(8):1169–1186, 2017. doi: 10.1177/0959683616683263. [p64, 76]
- A. C. Kemp, A. J. Wright, R. J. Edwards, R. L. Barnett, M. J. Brain, R. E. Kopp, N. Cahill, B. P. Horton, D. J. Charman, A. D. Hawkes, T. D. Hill, and O. van de Plassche. Relative sea-level change in Newfoundland, Canada during the past 3000 years. *Quaternary Science Reviews*, 201:89–110, 2018. doi: <https://doi.org/10.1016/j.quascirev.2018.10.012>. [p61, 62, 65, 76]
- N. S. Khan, E. Ashe, T. A. Shaw, M. Vacchi, J. Walker, W. R. Peltier, R. E. Kopp, and B. P. Horton. Holocene relative sea-level changes from near-, intermediate-, and far-field locations. *Current Climate Change Reports*, 1(4):247–262, 2015. doi: <https://doi.org/10.1007/s40641-015-0029-z>. [p61]
- N. S. Khan, E. Ashe, R. P. Moyer, A. C. Kemp, S. E. Engelhart, M. J. Brain, L. T. Toth, A. Chappel, M. Christie, R. E. Kopp, and B. P. Horton. Relative sea-level change in South Florida during the past 5000 years. *Global and Planetary Change*, 216:103902, 2022. ISSN 0921-8181. doi: <https://doi.org/10.1016/j.gloplacha.2022.103902>. [p76]
- J. R. Kirby, E. Garrett, and W. R. Gehrels. Holocene relative sea-level changes in northwest ireland: An empirical test for glacial isostatic adjustment models. *The Holocene*, 2023. doi: <https://doi.org/10.1177/09596836231169992>. [p64]
- R. E. Kopp. Does the mid-Atlantic United States sea level acceleration hot spot reflect ocean dynamic variability? *Geophysical Research Letters*, 40(15):3981–3985, 8 2013. doi: 10.1002/GRL.50781. [p61, 65]
- R. E. Kopp, F. J. Simons, J. X. Mitrovica, A. C. Maloof, and M. Oppenheimer. Probabilistic assessment of sea level during the last interglacial stage. *Nature*, 462(7275):863–867, 2009. ISSN 00280836. doi: <https://doi.org/10.1038/nature08686>. [p65]
- R. E. Kopp, A. C. Kemp, K. Bittermann, B. P. Horton, J. P. Donnelly, R. W. Gehrels, C. C. Hay, J. X. Mitrovica, E. D. Morrow, and S. Rahmstorf. Temperature-driven global sea-level variability in the Common Era. *Proceedings of the National Academy of Sciences of the United States of America*, 113(11):E1434–E1441, 2016. ISSN 10916490. doi: 10.1073/pnas.1517056113. [p61, 65]
- W. Marshall. *Chronohorizons: Handbook of Sea-Level Research*, chapter 25, pages 373–385. John Wiley & Sons, Ltd, 2015. ISBN 9781118452547. doi: <https://doi.org/10.1002/9781118452547.ch25>. [p62]

- A. McHutchon and C. E. Rasmussen. Gaussian Process training with input noise. *Advances in Neural Information Processing Systems 24: 25th Annual Conference on Neural Information Processing Systems 2011, NIPS 2011*, pages 1–9, 2011. URL [https://proceedings.neurips.cc/paper\\_files/paper/2011/file/a8e864d04c95572d1aece099af852d0a-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2011/file/a8e864d04c95572d1aece099af852d0a-Paper.pdf). [p61, 63]
- A. J. Meltzner, A. D. Switzer, B. P. Horton, E. Ashe, Q. Qiu, D. F. Hill, S. L. Bradley, R. E. Kopp, E. M. Hill, J. M. Majewski, D. H. Natawidjaja, and B. W. Suwargadi. Half-metre sea-level fluctuations on centennial timescales from mid-Holocene corals of Southeast Asia. *Nature Communications*, 8(1):14387, 2017. doi: 10.1038/ncomms14387. [p61]
- W. Peltier. Global Glacial Isostasy and the Surface of the Ice-Age Earth: The ICE-5G (VM2) Model and GRACE. *Annual Review of Earth and Planetary Sciences*, 32:111–149, 2004. doi: <https://doi.org/10.1146/annurev.earth.32.082503.144359>. [p72]
- M. Plummer. JAGS: A program for analysis of Bayesian graphical models using Gibbs sampling. *Proceedings of the 3rd International Workshop on Distributed Statistical Computing, TUWien*, 124, 2003. URL <https://api.semanticscholar.org/CorpusID:265976949>. [p63, 75]
- M. Plummer, A. Stukalov, and M. Denwood. rjags: Bayesian graphical models using MCMC. *R package version*, 4(6), 2016. URL <https://search.r-project.org/CRAN/refmans/rjags/html/rjags-package.html>. [p63]
- D. Pugh and P. Woodworth. *Tidal forces: Sea-Level Science: Understanding Tides, Surges, Tsunamis and Mean Sea-Level Changes*, page 36–59. Cambridge University Press, 2014. doi: 10.1017/CBO9781139235778.006. [p62]
- C. E. Rasmussen and C. K. I. Williams. Gaussian Processes for Machine Learning. Technical report, MIT Press, 2006. [p65]
- H. Rue, S. Martino, and N. Chopin. Approximate Bayesian inference for latent Gaussian models by using integrated nested Laplace approximations. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 71(2):319–392, 2009. doi: <https://doi.org/10.1111/j.1467-9868.2008.00700.x>. URL <https://rss.onlinelibrary.wiley.com/doi/abs/10.1111/j.1467-9868.2008.00700.x>. [p75]
- T. A. Shaw, A. J. Plater, J. R. Kirby, K. Roy, S. Holgate, P. Tutman, N. Cahill, and B. P. Horton. Tectonic influences on late Holocene relative sea levels from the central-eastern Adriatic coast of Croatia. *Quaternary Science Reviews*, 200:262–275, 2018. ISSN 0277-3791. doi: <https://doi.org/10.1016/j.quascirev.2018.09.015>. URL <https://www.sciencedirect.com/science/article/pii/S0277379118300659>. [p64]
- I. Shennan and B. Horton. Holocene land- and sea-level changes in Great Britain. *Journal of Quaternary Science*, 17(5-6):511–526, 2002. doi: <https://doi.org/10.1002/jqs.710>. [p64]
- I. Shennan, A. J. Long, and B. P. Horton. Handbook of sea-level research. *John Wiley & Sons, Ltd*, 2015. doi: <https://doi.org/10.1002/9781118452547>. [p62]
- G. L. Simpson. Modelling palaeoecological time series using generalised additive models. *Frontiers in Ecology and Evolution*, 6:149, 2018. doi: <https://doi.org/10.3389/fevo.2018.00149>. [p61]
- R. B. Stearns, S. E. Engelhart, A. C. Kemp, T. D. Hill, M. J. Brain, and D. R. Corbett. Within-region replication of late Holocene relative sea-level change: An example from southern New England, United States. *Quaternary Science Reviews*, 300:107868, 2023. ISSN 0277-3791. doi: <https://doi.org/10.1016/j.quascirev.2022.107868>. URL <https://www.sciencedirect.com/science/article/pii/S0277379122004991>. [p64, 76]
- M. Upton, A. Parnell, A. Kemp, E. Ashe, G. McCarthy, and N. Cahill. A noisy-input generalised additive model for relative sea-level change along the atlantic coast of north america. *arXiv preprint arXiv:2301.09556*, 2023. doi: <https://doi.org/10.48550/arXiv.2301.09556>. [p61, 62, 63, 65, 66, 67, 71, 76]

- J. S. Walker, R. E. Kopp, T. A. Shaw, N. Cahill, N. S. Khan, D. C. Barber, E. L. Ashe, M. J. Brain, J. L. Clear, D. R. Corbett, and B. P. Horton. Common Era sea-level budgets along the U.S. Atlantic coast. *Nature Communications*, 12(1):1841, 2021. doi: 10.1038/s41467-021-22079-2. [p61, 65, 76]
- P. L. Whitehouse. Glacial isostatic adjustment modelling: historical perspectives, recent advances, and future directions. *Earth Surf. Dynam.*, 6:401–429, 2018. doi: <https://doi.org/10.5194/esurf-6-401-2018>. [p66]
- H. Wickham. *ggplot2: Elegant Graphics for Data Analysis*. Springer-Verlag New York, 2016. ISBN 978-3-319-24277-4. URL <https://ggplot2.tidyverse.org>. [p67]
- S. Wood. Package ‘mgcv’. *R package version*, 1(29):729, 2015. URL <https://cran.r-project.org/web/packages/mgcv/index.html>. [p62]
- S. N. Wood. Low-Rank Scale-Invariant Tensor Product Smooths for Generalized Additive Mixed Models. *Biometrics*, 62(4):1025–1036, 2006. doi: <https://doi.org/10.1111/j.1541-0420.2006.00574.x>. [p65]
- S. N. Wood. P-splines with derivative based penalties and tensor product smoothing of unevenly distributed data. *Statistics and Computing*, 27(4):985–989, 2017a. doi: <https://doi.org/10.1007/s11222-016-9666-x>. [p65]
- S. N. Wood. *Generalized additive models: an introduction with R*. CRC Press, 2017b. doi: <https://doi.org/10.1201/9781315370279>. [p66, 69]
- P. Woodworth and R. Player. The Permanent Service for Mean Sea Level: An update to the 21st century. *Journal of Coastal Research*, 19:287–295, 03 2003. doi: <https://www.jstor.org/stable/4299170>. [p62, 68]

*Maeve Upton*  
Maynooth University  
Hamilton Institute, ICARUS and Department of Mathematics and Statistics  
Maynooth, Ireland  
ORCID: 0000-0002-1151-6657  
[uptonmaeve010@gmail.com](mailto:uptonmaeve010@gmail.com)

*Andrew Parnell*  
Maynooth University  
Hamilton Institute, ICARUS and Department of Mathematics and Statistics  
Maynooth, Ireland  
ORCID: 0000-0001-7956-7939  
[andrew.parnell@mu.ie](mailto:andrew.parnell@mu.ie)

*Niamh Cahill*  
Maynooth University  
Department of Mathematics and Statistics and ICARUS  
Maynooth, Ireland  
ORCID: 0000-0003-3086-550X  
[niamh.cahill@mu.ie](mailto:niamh.cahill@mu.ie)

# Rfssa: An R Package for Functional Singular Spectrum Analysis

by Hossein Haghbin, Jordan Trinka and Mehdi Maadooliat

**Abstract** Functional Singular Spectrum Analysis (FSSA) is a non-parametric approach for analyzing Functional Time Series (FTS) and Multivariate FTS (MFTS) data. This paper introduces Rfssa, an R package that addresses implementing FSSA for FTS and MFTS data types. Rfssa provides a flexible container, the funts class, for FTS/MFTS data observed on one-dimensional or multi-dimensional domains. It accepts arbitrary basis systems and offers powerful graphical tools for visualizing time-varying features and pattern changes. The package incorporates two forecasting algorithms for FTS data. Developed using object-oriented programming and Rcpp/RcppArmadillo, Rfssa ensures computational efficiency. The paper covers theoretical background, technical details, usage examples, and highlights potential applications of Rfssa.

## 1 Introduction

In recent times, advancements in data acquisition techniques have made it possible to collect data in high-resolution formats. Due to the presence of temporal-spatial dependence, one may consider this type of data as *functional data*. Functional Data Analysis (FDA) focuses on developing statistical methodologies for analyzing data represented as functions or curves. While FDA methods are particularly well-suited for handling smooth continuum data, they can also be adapted and extended to effectively analyze functional data that may not exhibit perfect smoothness, including high-resolution data and data with inherent variability. The widely-used R package for FDA is [fda](#) (Ramsay et al., 2023), which is designed to support analysis of functional data, as described in the textbook by Ramsay and Silverman (2005). Additionally, there are over 40 other R packages available on CRAN that incorporate functional data analysis, such as [funFEM](#) (Bouveyron, 2021), [fda.usc](#) (Frbrero-Band and de la Fuente, 2012), [refund](#) (Goldsmith et al., 2023), [fdapace](#) (Gajardo et al., 2022), [funData](#) (Happ-Kurz, 2020), [ftsspec](#) (Tavakoli, 2015), [rainbow](#) (Shang and Hyndman, 2022), and [ftsa](#) (Hyndman and Shang, 2023). One crucial initial requirement for any of these packages is to establish a framework for representing and storing infinite-dimensional functional observations. The [fda](#) package, for instance, employs the fd class as a container for functional data defined on a one-dimensional (1D) domain. An fd object represents functional data as a finite linear combination of known basis functions (e.g., Fourier, B-splines, etc.), storing both the basis functions and their respective coefficients for each curve. This representation aligns with the practical implementation found in many papers within the field of FDA. Conversely, several other R packages store functional data in a discrete form evaluated on grid points (e.g., [fda.usc](#), [refund](#), [funData](#), [rainbow](#), and [fdapace](#)). These packages also provide the capability to analyze functions beyond the one-dimensional case, such as image data treated as two-dimensional (2D) functions (e.g., [refund](#), [fdasrvf](#), and [funData](#)). To the best of our knowledge, packages that support representation beyond 1D functions utilize the grid point representation for execution and storage. Moreover, recent packages have been developed to handle multivariate functional data, which consist of more than one function per observation unit. Examples of such packages include [roahd](#), [fda.usc](#), and [funData](#). While some recent FDA packages have focused on analyzing and implementing techniques for Functional Time Series (FTS), where sequences of functions are observed over time, none of them handle Multivariate FTS (MFTS) or multidimensional MFTS. For example, see the packages [ftsspec](#), [rainbow](#), and [ftsa](#). In summary, there is still a need for a unified and flexible container for FTS/MFTS data, defined on either one or multidimensional domains. The funts class in [Rfssa](#) (Haghbin et al., 2023), the package discussed in this article, aims to address this gap. One of the primary contributions of the package is its capacity to handle and visualize 2-dimensional FTS, including image data. Furthermore, the package accommodates MFTS, especially when observed on distinct domains. This flexibility empowers users to analyze and visualize FTS with multiple variables, even when they do not share the same domain. Notably, the [Rfssa](#) package introduces novel visualization tools (as exemplified in Figure 5). These tools include heatmaps and 3D plots, thoughtfully designed to provide a deeper understanding of functional patterns over time. They enhance the ability to discern trends and variations that might remain inconspicuous in conventional plots. An additional feature of the funts class is its ability to accept any arbitrary basis system as input for the class constructor, including FDA basis functions or even empirical basis represented as matrices evaluated at grid points. The classes in the [Rfssa](#) package are developed using the S3 object-oriented programming system, and for computational efficiency, significant portions of the package are implemented using the [Rcpp](#)/[RcppArmadillo](#) packages. Notably, the package includes a shiny web application that provides a user-friendly GUI for implementing Functional Singular

Spectrum Analysis (FSSA) on real or simulated FTS/MFTS data. The **Rfssa** package was initially developed to implement FSSA for FTS, as discussed in the work of [Haghbin et al. \(2021\)](#). FSSA extends Singular Spectrum Analysis (SSA), a model-free procedure commonly used to analyze time series data. The primary goal of SSA is to decompose the original series into a collection of interpretable components, such as slowly varying trends, oscillatory patterns, and structureless noise. Notably, SSA does not rely on restrictive assumptions like stationarity, linearity, or normality ([Golyandina and Zhitljavskiy, 2013](#)). It's worth noting that SSA finds applications beyond the functional framework, including smoothing and forecasting purposes ([Hassani and Mahmoudvand, 2013; de Carvalho and Rua, 2017](#)). The non-functional version of FSSA, known as SSA, has previously been implemented in the **Rssa** package ([Golyandina et al., 2015](#)) and the **ASSA** package ([de Carvalho and Martos, 2020](#)). The **Rssa** package provides various visualization tools to facilitate the grouping stage, and the **Rfssa** package includes equivalent functional versions of those tools ([Golyandina et al., 2018](#)). While the foundational theory of FSSA was originally designed for univariate FTS, it has since been extended to handle multidimensional FTS data, referred to as Multivariate FSSA (MFSSA) ([Trinka et al., 2022](#)). Furthermore, in line with the developments in SSA for forecasting, two distinct algorithms known as Recurrent Forecasting (FSSA R-forecasting) and Vector Forecasting (FSSA V-forecasting) were introduced for FSSA by [Trinka et al. \(2023\)](#). Both of these forecasting algorithms, along with the capabilities for handling MFSSA, have been seamlessly integrated into the most recent version of the **Rfssa** package. The remainder of this manuscript is organized as follows. Section 2 introduces the FTS/MFTS data preparation theory used in the `funts` class. Section 3 discusses the FSSA methodology, including the basic schema of FSSA, FSSA R-forecasting, and FSSA V-forecasting. Technical details of the **Rfssa** package are provided in Section 4, where we describe the available classes in the package and illustrate their practical usage with examples of real data. Section 5 focuses on the reconstruction stage and FSSA/MFSSA forecasting. In Section 6, we provide a summary of the embedded shiny app. Finally, we conclude the paper in Section 7.

## 2 Data preparation in FTS

Define  $\mathbf{y}_N = (y_1, \dots, y_N)$  to be a collection of observations from an FTS. In the theory of FTS,  $y_i$ 's are considered as functions in the space  $\mathbb{H} = L^2(\mathcal{T})$  where  $\mathcal{T}$  is a compact subset of  $\mathbb{R}$ . Let  $s \in \mathcal{T}$  and consider  $y_i(s) \in \mathbb{R}^p$ , the sequence of  $\mathbf{y}_N$  is called (univariate) FTS if  $p = 1$ , and multivariate FTS (or MFTS) if  $p > 1$ . In the realm of functional data analysis, we operate under the assumption that the underlying sample functions, denoted as  $y_i(\cdot)$ , exhibit smoothness for each sample  $i$ , where  $i = 1, \dots, N$ . Nevertheless, in practical scenarios, observations are typically acquired discretely at a set of grid points and are susceptible to contamination by random noise. This phenomenon can be represented as follows:

$$Y_{i,k} = y_i(t_k) + \varepsilon_{i,k}, \quad k = 1, \dots, K. \quad (1)$$

In this expression,  $t_k \in \mathcal{T}$ , and  $K$  denotes the count of discrete grid points across all samples. The  $\varepsilon_{i,k}$  terms represent i.i.d. random noise. To preprocess the raw data, it is customary to employ smoothing techniques, converting the discrete observations  $Y_{i,k}$  into a continuous form,  $y_i(\cdot)$ . This is typically performed individually for each variable and sample. One widely used approach is finite basis function expansion ([Ramsay and Silverman, 2005](#)). In this method, a set of basis functions  $\{\nu_i\}_{i \in \mathbb{N}}$  is considered (not necessarily orthogonal) for the function space  $\mathbb{H}$ . Each sample function  $y_i(\cdot)$  in (1) is then considered as a finite linear combination of the first  $d$  basis functions:

$$y_i(s) = \sum_{j=1}^d c_{ij} \nu_j(s). \quad (2)$$

Subsequently, the coefficients  $c_{ij}$  can be estimated using least square techniques. By adopting the linear representation form for the functional data in (2), we establish a correspondence between each function  $y_i(\cdot)$  and its coefficient vector  $\mathbf{c}_i = (c_{ij})_{j=1}^d$ . As a result, the coefficient vectors  $\mathbf{c}_i$  can serve to store and retrieve the original functions,  $y_i(\cdot)$ 's. This arises from the inherent isomorphism between two finite vector spaces of the same dimension (in this case,  $d$ ). Consequently,  $\mathbf{c}_i$ 's are stored as the primary attribute of `funts` objects within the **Rfssa** package. Take two elements  $x, y \in \mathbb{H}$  with corresponding coefficient vectors  $\mathbf{c}_x$  and  $\mathbf{c}_y$ . Then, the inner product of  $x, y$  can be computed in matrix form as  $\langle x, y \rangle = \mathbf{c}_x^\top \mathbf{G} \mathbf{c}_y$ , where  $\mathbf{G} = [\langle \nu_i, \nu_j \rangle]_{i,j=1}^d$  is the Gram matrix. It is important to note that  $\mathbf{G}$  is Hermitian. Furthermore, because the basis functions  $\{\nu_i\}_{i=1}^d$  are linearly independent,  $\mathbf{G}$  is positive definite, making it invertible ([Horn and Johnson, 2012, Thm. 7.2.10](#)). Moreover, let  $A : \mathbb{H} \rightarrow \mathbb{H}$  be a linear operator and  $y = A(x)$ . Then,  $\mathbf{c}_y = \mathbf{G}^{-1} \mathbf{A} \mathbf{c}_x$ , where  $\mathbf{A} = [\langle A(\nu_j), \nu_i \rangle]_{i,j=1}^d$  is called the corresponding matrix of the operator  $A$ . It is worth noting that while the FSSA theory extends to arbitrary dimensions, practical implementation for dimensions greater than 2 introduces considerable computational complexity. Moreover, high-dimensional FTS data are relatively rare in real-world

applications. Therefore, within the **Rfssa** package, we have chosen to confine the funts object to support functions observed over domains that are one or two-dimensional. In the **Rfssa** package, the task of preprocessing the raw discrete observations and converting those to the funts object is assigned to the funts(.) constructor.

### 3 An overview of the FSSA methodology

FSSA is a nonparametric technique to decompose FTS and MFTS, and the methodology can also be used to forecast such data (Haghbin et al., 2021; Trinka et al., 2022, 2023); it can also be used as a visualization tool to illustrate the concept of seasonality and periodicity in the functional space over time.

#### 3.1 Basic schema of FSSA

Basic FSSA consists of two stages where each stage includes two steps. We outline the four steps of the FSSA algorithm here.

##### I) First stage: decomposition

###### 1. Embedding

For a positive integer,  $L < N/2$ , let  $\mathbb{H}^L$  be the Cartesian product of  $L$  copies of  $\mathbb{H}$  and define the *trajectory* operator  $\mathcal{X} : \mathbb{R}^K \rightarrow \mathbb{H}^L$  with

$$\mathcal{X}\mathbf{a} := \sum_{j=1}^K a_j \mathbf{x}_j, \quad \mathbf{a} = (a_1, \dots, a_K)^\top \in \mathbb{R}^K. \quad (3)$$

where  $K = N - L + 1$  and

$$\mathbf{x}_j(s) := (y_j(s), y_{j+1}(s), \dots, y_{j+L-1}(s))^\top, \quad j = 1, \dots, K, \quad (4)$$

are called *lag-vectors*. One may consider the trajectory operator  $\mathcal{X}$  in (3) as an  $L \times K$  matrix with functional entries in the form of

$$\mathcal{X}(s) = \begin{bmatrix} y_1(s) & y_2(s) & y_3(s) & \cdots & y_K(s) \\ y_2(s) & y_3(s) & y_4(s) & \cdots & y_{K+1}(s) \\ y_3(s) & y_4(s) & y_5(s) & \cdots & y_{K+2}(s) \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ y_L(s) & y_{L+1}(s) & y_{L+2}(s) & \cdots & y_N(s) \end{bmatrix}. \quad (5)$$

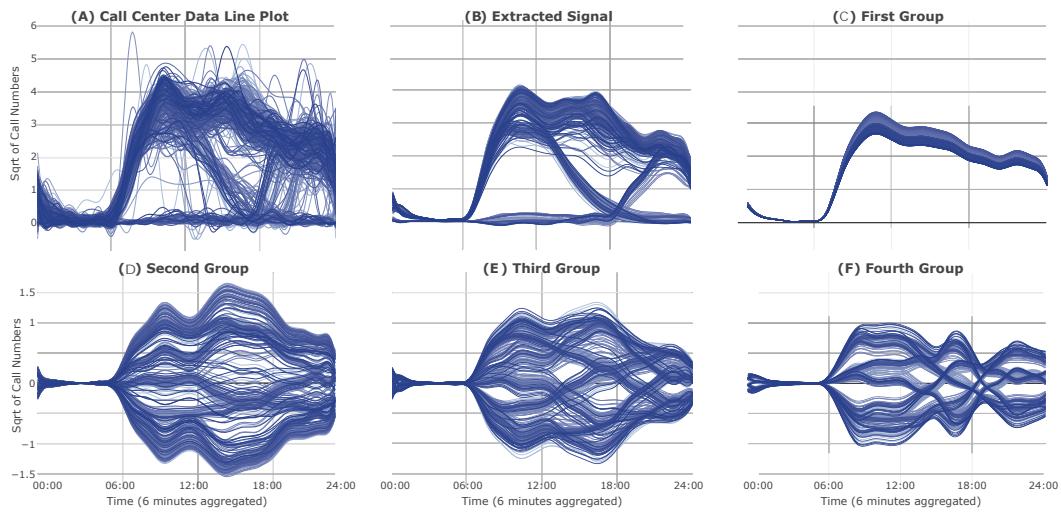
Note that the antidiagonal elements of the matrix in (5) are all equal. Such matrices are called Hankel, and since  $\mathcal{X}(s)$  is a Hankel matrix for any  $s$  in the domain, we call  $\mathcal{X}$  a Hankel operator. In practice, a main challenge is how to use the original basis of  $\mathbb{H}$  to represent the lag-vectors in  $\mathbb{H}^L$ . To do this, one may define a quotient sequence  $q_k$ , and a remainder sequence  $r_k$ , by  $k = (q_k - 1)L + r_k$ , where  $1 \leq r_k \leq L$ , and  $1 \leq q_k \leq d$ . Now, consider  $\boldsymbol{\phi}_k$  as a functional vector of length  $L$  with all zero functions, except the  $r_k$ -th element, which is  $v_{q_k}$ . Using this definition, the lag-vector  $\mathbf{x}_j$  given in (4) can be represented as a linear combination of  $\{\boldsymbol{\phi}_k\}_{k=1}^{Ld}$  with the corresponding coefficient vector  $\mathbf{b}_j = (c_{1j}, \dots, c_{1,j+L-1}, c_{2j}, \dots, c_{2,j+L-1}, \dots, c_{d,j+L-1})^\top$ .

###### 2. FSVD

We apply the functional singular value decomposition (FSVD) to  $\mathcal{X}$  and obtain a collection of singular values,  $\{\sqrt{\lambda_i}\}_{i=1}^r$ , orthonormal right singular vectors  $\{\mathbf{v}_i\}_{i=1}^r$  (that are elements of  $\mathbb{R}^K$ ), and orthonormal left singular functions  $\{\boldsymbol{\psi}_i\}_{i=1}^r$  (that are elements of  $\mathbb{H}^L$ ). The collection  $(\sqrt{\lambda_i}, \boldsymbol{\psi}_i, \mathbf{v}_i)$  will be called the  $i^{th}$  eigentriple of the FSVD, and  $r$  is the rank of  $\mathcal{X}$ :

$$\mathcal{X} = \sum_{i=1}^r \mathcal{X}_i = \sum_{i=1}^r \sqrt{\lambda_i} \mathbf{v}_i \otimes \boldsymbol{\psi}_i, \quad (6)$$

where  $\mathcal{X}_i : \mathbb{R}^K \rightarrow \mathbb{H}^L$  is a rank one elementary operator. To implement the FSVD of  $\mathcal{X}$  given in (6), let  $\mathbf{B} := [\mathbf{b}_1, \dots, \mathbf{b}_K]$ ,  $\mathbf{G} := [\langle \boldsymbol{\phi}_i, \boldsymbol{\phi}_j \rangle_{\mathbb{H}^L}]_{i,j=1}^{Ld}$ ,  $\mathbf{X} := \mathbf{G}^{1/2} \mathbf{B}$ , and  $(\sqrt{\lambda_i}, \mathbf{u}_i, \mathbf{v}_i)$ 's be the eigentriples of the SVD of the matrix  $\mathbf{X}$ . It can be shown that  $(\sqrt{\lambda_i}, \boldsymbol{\psi}_i, \mathbf{v}_i)$  is the  $i^{th}$  eigentriple of the FSVD of  $\mathcal{X}$ , where the left singular function,  $\boldsymbol{\psi}_i$ , is corresponding to the coefficient vector  $\mathbf{G}^{-1/2} \mathbf{u}_i$ . See (Haghbin et al., 2021) for more details.



**Figure 1:** (A): Line plot of Callcenter data; (B): FTS reconstructed using  $I_5 = \{1, \dots, 7\}$ ; (C): FTS reconstructed by setting  $I_1 = \{1\}$ ; (D): FTS reconstructed using  $I_2 = \{2, 3\}$ ; (E): FTS reconstructed from  $I_3 = \{4, 5\}$ ; (F): FTS reconstructed from  $I_4 = \{6, 7\}$ .

These steps can also be extended to the multivariate case, i.e. MFSSA. See Trinka et al. (2022) for more details. In the **Rfssa** package, the results of the decomposition stage are held in an object from the **fssa** class. The constructor, **fssa(·)**, performs the decomposition for both FSSA and MFSSA algorithms and returns an object of class **fssa**. Further discussion about the attributes and methods of the **fssa** class is given in the technical details section.

## II) Second stage: reconstruction

### 3. Grouping

We partition the set of indices  $\{1, \dots, r\}$  into disjoint sets  $I_q$ , where  $q \in \{1, \dots, m\}$  and  $m \leq r$ . From here, we obtain the group  $\mathcal{X}_{I_q}$  by combining the respective elementary operators accordingly:

$$\mathcal{X}_{I_q} = \sum_{i \in I_q} \mathcal{X}_i.$$

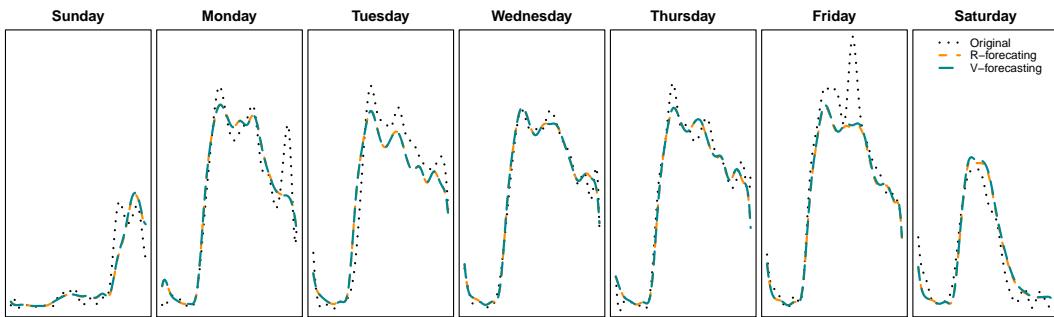
Exploratory plots of singular values, right singular vectors, and left singular functions that investigate the different modes of variation extracted in the decomposition stage are used to decide how to form the sets,  $I_q$ , and we discuss such plots in further detail in section five.

### 4. Hankelization

Since each  $\mathcal{X}_{I_q}$  is not necessarily Hankel, we perform diagonal averaging of the entries to Hankelize each operator. From each Hankelized  $\mathcal{X}_{I_q}$ , we obtain an FTS,  $\mathbf{y}_N^q$ , that describes main characteristics of  $\mathbf{y}_N$  such as mean, seasonal, trend, and noise behaviors.

For the reconstruction stage, the **Rfssa** package provides the function **freconstruct(·)** which returns a list of objects of class **funts** associated with the groups specified by the user. If the supplied input to the **fssa(·)** function is an MFTS, the signal extraction process is almost identical as compared to the univariate case with the exception that now, we have that each element of the time series is a tuple of functions comprised of elements observed over one or two-dimensional domains.

As a motivating example, consider the Callcenter dataset in Figure 1(A), which has been previously discussed in Maadoliat et al. (2015) and is included in the package. This dataset records the number of calls received by a call center in 1999. Each functional observation corresponds to the square root of the daily call count, and the entire FTS represents all the days in 1999. When dealing with time series data exhibiting known periodic patterns, it is a common practice in SSA to choose the window length,  $L$ , as a multiple of this underlying periodicity. This choice ensures that the method effectively extracts the corresponding periodic components (Golyandina et al., 2001). For our illustrative example using the Callcenter dataset, we specifically opt for a window length of  $L = 28$ . This selection allows us to effectively capture the weekly periodic patterns that inherently exist within this functional time series (Haghbin et al., 2021). After applying FSSA, we obtain reconstructed FTS representations under different grouping considerations, as illustrated in Figures 1(B-F). Particularly noteworthy is the reconstructed FTS obtained using the leading 7 eigentriples of the FSVD, as shown in Figure 1(B). It is important to mention that the selection of groups in FTS is typically guided by



**Figure 2:** The prediction of the last 7 days of the year 1999 for the Callcenter data, based on FSSA R-forecasting and V-forecasting, and comparing the results with the observed functions.

various SSA-type plotting tools, which rely on the similarity in the harmonic structure of extracted elementary components. Such SSA-type plots have been available for non-functional time series in the **Rssa** package, and analogous tools for the functional case have been developed in **Rfssa** (Figure 7). However, it's worth noting that this extension presented its own set of challenges, both in theory and implementation. Detailed code for generating the reconstructed functions shown in Figure 1 can be found in the subsequent sections of this paper.

### 3.2 FSSA Forecasting

After the decomposition stage, one may perform forecasting instead of reconstruction (Figure 3). R-forecasting and V-forecasting are two main defined approaches in FSSA/MFSSA (Trinka et al., 2023). In both forecasting methods, the goal is to obtain the FTS,  $\mathbf{g}_{N+M}^q = (g_1^q, \dots, g_{N+M}^q)$ , where the first  $N$  terms are close to  $\mathbf{y}_N^q$  and the goal is to forecast the remaining elements,  $\{g_i^q\}_{i=N+1}^{N+M}$ . Forecasts in the R-forecasting method are obtained using a linear combination of the previous  $L - 1$  elements of  $\mathbf{g}_{N+M}^q$ . Consider a positive integer  $k < r$  and define  $\mathbf{V} = \sum_{i=1}^k \pi_i \otimes \pi_i$ , where  $\pi_i \in \mathbb{H}$  is the last element of the left singular function  $\psi_i$ . Define  $\mathcal{A}_j : \mathbb{H} \rightarrow \mathbb{H}$  such that  $\mathcal{A}_j = \sum_{i=1}^k \psi_{j,i} \otimes (\mathcal{I} - \mathbf{V})^{-1} \pi_i$ , where  $\psi_{j,i}$  is the  $j^{\text{th}}$  component of  $\psi_i$ , and  $\mathcal{I} : \mathbb{H} \rightarrow \mathbb{H}$  is the identity operator. Then the R-forecasting can be obtained using the following equation

$$g_i^q = \begin{cases} y_i^q, & i = 1, \dots, N \\ \sum_{j=1}^{L-1} \mathcal{A}_j g_{i+j-L}^q, & i = N + 1, \dots, N + M. \end{cases} \quad (7)$$

Forecasts in the V-forecasting method are obtained by predicting functional vectors. One may define an orthogonal projection,  $\Pi : \mathbb{H}^{L-1} \rightarrow \mathbb{H}^{L-1}$ , that projects onto the space spanned by  $\{\psi_i^\nabla\}_{i=1}^k$ , where  $\psi_i^\nabla \in \mathbb{H}^{L-1}$  is formed from the first  $L - 1$  elements of  $\psi_i$ . Now let  $\mathcal{Q} : \mathbb{H}^L \rightarrow \mathbb{H}^L$  be given by

$$\mathcal{Q}(\mathbf{x}) = \left( \frac{\Pi(\mathbf{x}^\Delta)}{\sum_{j=1}^{L-1} \mathcal{A}_j x_j^\Delta} \right), \quad \mathbf{x} \in \mathbb{H}^L \quad (8)$$

where  $\mathbf{x}^\Delta$  contains the last  $L - 1$  components of  $\mathbf{x}$ , and  $x_j^\Delta$  is the  $j^{\text{th}}$  component of  $\mathbf{x}^\Delta$ . The V-forecasting algorithm is given in the following steps:

1. Define  $\mathbf{w}_j^q$  as
$$\mathbf{w}_j^q = \begin{cases} \mathbf{x}_j^\Delta, & j = 1, \dots, K \\ \mathcal{Q}\mathbf{w}_{j-1}^q, & j = K + 1, \dots, K + M, \end{cases}$$
where  $\{\mathbf{x}_j^q\}_{j=1}^K$  spans the range of operator  $\mathcal{X}_{I_q}$ .
2. Form the operator  $\mathcal{W}^q : \mathbb{R}^{K+M} \rightarrow \mathbb{H}^L$  whose range is linearly spanned by the set  $\{\mathbf{w}_i^q\}_{i=1}^{K+M}$ .
3. Hankelize  $\mathcal{W}^q$  in order to extract the FTS  $\mathbf{g}_{N+M}^q$ .
4. The functions,  $g_{N+1}^q, \dots, g_{N+M}^q$ , form the  $M$  terms of the FSSA vector forecast.

Continuing from the previous example, we utilized the first 358 days of the year to train the FSSA model on the Callcenter dataset. Subsequently, we employed the FSSA R-forecasting and V-forecasting methods to make predictions for the final week ( $len = 7$ ). Figure 2 illustrates both the actual and forecasted curves for each day of the week, employing each respective method. Detailed code for

Function	Descriptions	Main arguments	Returns
<code>funts(·)</code>	Create FTS/MFTS objects	Discretely sampled data or coefficients, basis system specifications, a set of argument values corresponding to the observations in $X$ , the time specifications arguments.	An object of class <code>funts</code> .
<code>fssa(·)</code>	Performs the decomposition (including embedding and FSVD steps) stage for FTS/MFTS data.	An object of class <code>funts</code> and window length $L$ .	An object of class <code>fssa</code> .
<code>freconstruct(·)</code>	Performs the reconstruction (including grouping and Hankelization steps) stage.	An object of class <code>fssa</code> and a list of numeric vectors includes indices of elementary components of a group.	A list of <code>funts</code> objects reconstructed according to the specified groups.
<code>fforecast(·)</code>	Performs FSSA R-forecasting or FSSA V-forecasting.	An object of class <code>fssa</code> , a list of numeric vectors includes indices of elementary components of a group used for reconstruction and forecasting, and forecast horizon $h$ .	An object of class <code>fforecast</code> .

**Table 1:** A summary of FSSA written functions in the `Rfssa` package.

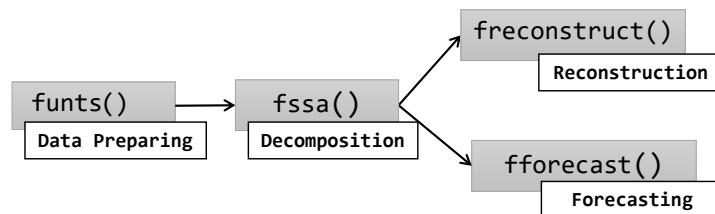
this process can be found in the subsequent sections. In the `Rfssa` package, we offer the `fforecast(·)` function designed for the execution of R-forecasting or V-forecasting algorithms. This function expects an input argument of class `fssa` and yields an output of class `fforecast`. The latter comprises a list of objects of class `funts`, with each `funts` representing a forecasted group.

## 4 Technical details of the `Rfssa` package

The roadmap of the main functions used in the `Rfssa` package is given in Figure 3. The inputs and outputs of these functions are described in Table 1. As it can be seen from Table 1, three classes (`funts`, `fssa` and `fforecast`) are used to support the return objects of these functions. The `funts(·)`, `fssa(·)` and `fforecast(·)` functions are the constructors of the `funts`, `fssa` and `fforecast` classes, respectively. In the rest of this section we present these classes with illustrative examples, and later we describe the reconstruction and forecasting functions in detail.

### 4.1 The `funts` class

The `funts(·)` constructor is used to create an S3 object of class `funts`. This object is designed to encapsulate various forms of FTS, including both univariate and multivariate types. It offers a versatile framework for the creation and manipulation of `funts` objects, accommodating different basis systems and dimensions. It accepts the following arguments:

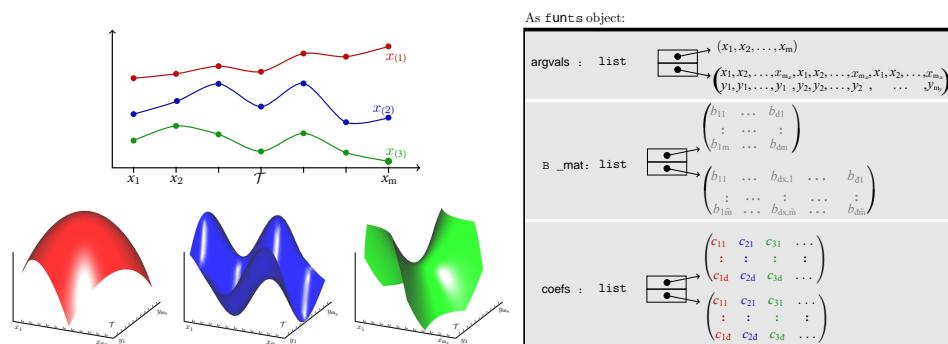


**Figure 3:** The roadmap of the `Rfssa` package.

- **X:** A matrix, three-dimensional array, or a list of matrix or array objects. When `method="data"`, it represents the observed curve values at discrete sampling points or argument values. When `method="coefs"`, **X** specifies the coefficients corresponding to the basis system defined in `basisobj`. If **X** is a list, it defines a multivariate FTS, with each element being a matrix or three-dimensional array object. In matrix objects, rows correspond to argument values, and columns correspond to the length of the FTS. In three-dimensional array objects, the first and second dimensions correspond to argument values, and the third dimension to the length of the FTS.
- **basisobj:** This argument should be an object of class `basisfd`, a matrix of empirical basis, or a list of `basisfd` or empirical basis objects. In the case of empirical basis, rows correspond to basis functions, and columns correspond to grid points.
- **argval:** A vector list of length  $p$ , representing a set of argument values corresponding to the observations in **X**. Each entry in this list should either be a numeric value or a list of numeric elements, depending on the dimension of the domain over which the variable is observed. It's worth noting that these values can vary from one variable to another. If `argval` is set to `NULL`, the default values are the integers from 1 to  $n$ , where  $n$  is the size of the first dimension in the **X** argument.
- **method:** This parameter determines the type of the **X** matrix, and it can take one of two values: `coefs` or `data`.
- **start** and **end:** Specify the time of the first and last observations. They can be a single positive integer or an object of classes `Date`, `POSIXct`, or `POSIXt`, representing a natural time unit.
- **tname**, **vnames** and **dnames:** These parameters accept strings or lists of strings to specify the names of time, variables, and domains.

The `funts()` constructor offers flexibility to users. Users can either provide their custom basis or request `Rfssa` to generate the basis for them, leveraging the capabilities of the `fda` package. It is assumed that each variable is observed over a regular and equidistant grid. Furthermore, each variable in the `funts` object is assumed to be observed over either a one or two-dimensional domain, as illustrated in Figure 4. To enhance the representation of time, the `funts` function introduces two parameters, namely `start` and `end`, which capture the time series duration. This design allows users to specify time information in a structured and standardized manner. Users have the flexibility to set `start` and `end` using various time and date classes, such as `Date`, `POSIXct`, or `POSIXt`. If users do not provide the `start` and `end` arguments, default values are used, with `start=1` and `end=N`, where  $N$  represents the length of the time series. This default approach aligns with common practices, as seen in classes like `ts` in the `stats` package, as well as `fts` classes in `rainbow` or `ftsa`. An object of class `funts` is a list encompassing the following elements:

- **N:** Represents the length of the time series.
- **dimSupp:** A list specifying the dimensions of the support domain of the variables.
- **time:** The time object.
- **coefs:** A list containing basis coefficients.
- **basis:** A list of basis systems.
- **B\_mat:** Evaluated basis functions on initial arguments.
- **argval:** Initial arguments of the observed functions.



**Figure 4:** The main roadmap of `funts` objects.

The `funts` class provides essential functionalities for managing FTS objects, ensuring users have well-defined basic operations. It supports arithmetic operations like addition and multiplication, along

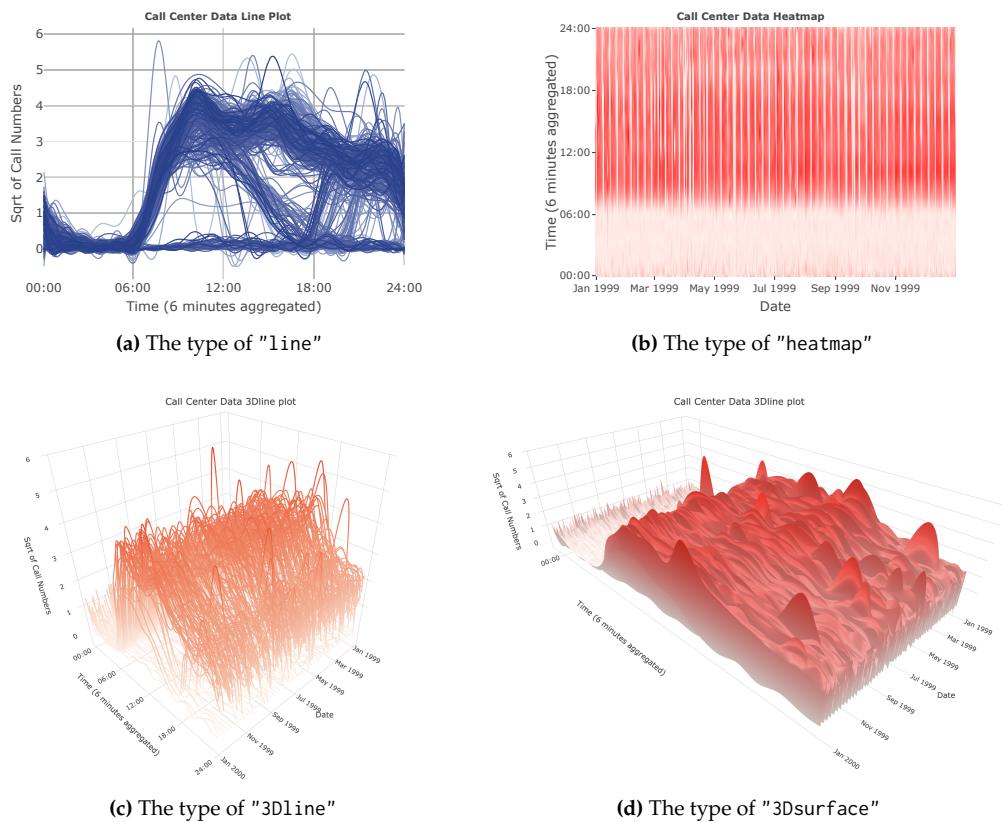
with indexing methods. Additionally, three generic methods, `length()`, `print()`, and `plot()`, are available. The `eval.funts()` method allows users to evaluate a `funts` object on a specified grid. To determine if an object belongs to the `funts` class, the `is.funts()` method is provided. Converting objects from the `fda` package's `fd` class or the `rainbow` package's `fts` class to a `funts` object is made easy using the `as.funts()` function. The strength of `funts` objects lies in their powerful visualization capabilities within the field of FTS. Users can create two types of plots using the graphical commands `plot()` and `plotly_funts()`. The `plot()` method utilizes base graphics objects to generate FTS plots. On the other hand, the `plotly_funts()` function offers a versatile `plotly` platform for visualizing FTS data, providing several plot types: `line`, `3Dline`, `3Dsurface`, and `heatmap`. These plots make it easier to detect trends or patterns within each curve's behavior over time. Furthermore, users can directly apply the `plotly_funts()` function to objects from the `fd`, `fds`, or `fts` classes in packages like `rainbow`, `fds`, `ftsa`, or `fda`, without the need for conversion to `funts` objects. Additionally, when converting objects from packages like `fds` or `fts`, the `xname` and `yname` arguments are automatically captured and used as the `xlab` and `zlab` arguments, ensuring that resulting plots are informative and intuitive. To demonstrate the capabilities of the `Rfssa` package for handling FTS, two illustrative examples are provided. The first example showcases the `Callcenter` dataset consisting of curves observed over a one-dimensional domain. In contrast, the second example involves a bivariate FTS dataset, which includes a sequence of two types of remote sensing images (two-dimensional functional data domain).

### Example: Creating `funts` objects for FTS

The first example uses the raw `Callcenter` dataset which was discussed before. The `funts` object can be made and plotted using the following codes. The generated plots are shown in Figure 5.

```
# Load necessary libraries
require(Rfssa)
require(fda)
# Load Callcenter data
Call_data <- loadCallcenterData()
# Prepare the data
D <- matrix(sqrt(Call_data$calls)), nrow = 240)
bs1 <- create.bspline.basis(c(0, 23), 22)
# Create a 'funts' object
Call_funts <- funts(D, bs1, start = as.Date("1999-1-1"),
                     vnames = "Sqrt of Call Numbers",
                     dnames = "Time (6 minutes aggregated)",
                     tname = "Date")
xtlab <- list(c("00:00", "06:00", "12:00", "18:00", "24:00"))
xtloc <- list(c(1, 60, 120, 180, 240))
# Generate a line plot using Plotly
plotly_funts(Call_funts, main = "Call Center Data Line Plot",
              xticklabels = xtblab, xticklocs = xtloc)
# Generate a heatmap plot using Plotly
plotly_funts(Call_funts, type = "heatmap", main = "Call Center Data Heatmap",
              xticklabels = xtblab, xticklocs = xtloc)
# Generate a 3D line plot using Plotly
plotly_funts(Call_funts, type = "3Dline", main = "Call Center Data 3Dline plot",
              xticklabels = xtblab, xticklocs = xtloc);
# Generate a 3D surface plot using Plotly
plotly_funts(Call_funts, type = "3Dsurface", main = "Call Center Data 3Dsurface plot",
              xticklabels = xtblab, xticklocs = xtloc);
```

As one can see from Figure 5(A), the overlapping line plot is a common way to view FTS data observed over a one-dimensional domain, where the curves that are recorded on dates closer to the start of 1999 are given in light blue while functions obtained on later dates are plotted in a darker blue. The heatmap plot in Figure 5(B) is a newer technique used to visualize FTS data observed over a one-dimensional domain, allowing the user to see the evolution of the FTS over time instead of relying on different colorings of the curves to specify date. Such one-dimensional FTS can also be represented in a interactive 3D view. To obtain such plots the user simply needs to specify `type="3Dsurface"` or `type="3Dline"`.



**Figure 5:** The plot types of the generated plots by `plotly_funts()`

### Example: Creating funts objects for MFTS

The second example considers two collections of images where each image is drawn from a region southeast of the city of Jambi, Indonesia, located between latitudes of  $1.666792^{\circ}$  S -  $1.598042^{\circ}$  S and longitudes of  $103.608963^{\circ}$  E -  $103.677742^{\circ}$  E. The images were recorded using the MODerate Resolution Imaging Spectroradiometer (MODIS) Terra satellite with a resolution of 250 meters every 16 days starting from February 18, 2000 and ending July 28, 2019. The output of each image is the normalized difference vegetation index (NDVI) which is used to quantify how much vegetation is present and enhanced vegetation index (EVI). The NDVI values closer to one being indicative of more vegetation and values closer to zero indicate less vegetation (see Haghbin et al., 2021, for more details). The following code can be used to load the data, define a funts object for a MFTS, slice the funts object to select a specific variable (here NDVI), and plot the smoothed images over time in an animation, that a snapshot is given in Figure 6.

```
# Load the Jambi dataset
Jambi <- loadJambiData()
# Extract NDVI and EVI array data
NDVI <- Jambi$NDVI
EVI <- (Jambi$EVI)
# Create a list containing NDVI and EVI array data
Jambi_Data <- list(NDVI, EVI)
# Create a multivariate B-spline basis
require(fda)
bs2 <- create.bspline.basis(c(0, 1), 13)
bs2d <- list(bs2, bs2)
bsmv <- list(bs2d, bs2d)
# Create a funts object
Y_J <- funts(X = Jambi_Data,
              basisobj = bsmv,
              start = as.Date("2000-02-18"), end = as.Date("2019-07-28"),
              vnames = c("NDVI", "EVI"), tname = "Date",
              dnames = list(c("Latitude", "Longitude"), c("Latitude", "Longitude")))
# Create a Plotly-based visualization of the NDVI Image
```

```
plotly_funts(Y_J[, 1],
  main = "NDVI Image (Jambi)",
  xticklabels = list(c("103.61\u00B0 E", "103.68\u00B0 E")),
  yticklabels = list(c("1.67\u00B0 S", "1.60\u00B0 S")),
  xticklocs = list(c(0, 1)),
  yticklocs = list(c(0, 1)),
  color_palette = "RdYlGn");
```

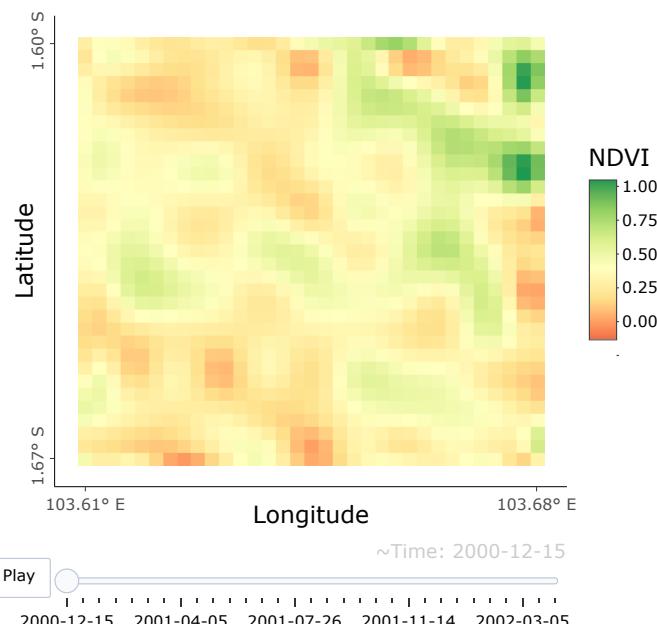
## 4.2 The fssa Class

Once the funts object is created and  $L$  is chosen, one can apply the `fssa(·)` constructor to obtain an S3 object of class `fssa` that contains our singular values, left singular functions, and right singular vectors. An object of class `fssa` is a list of right singular functions, which is packed in an object of class `funts` and the following components:

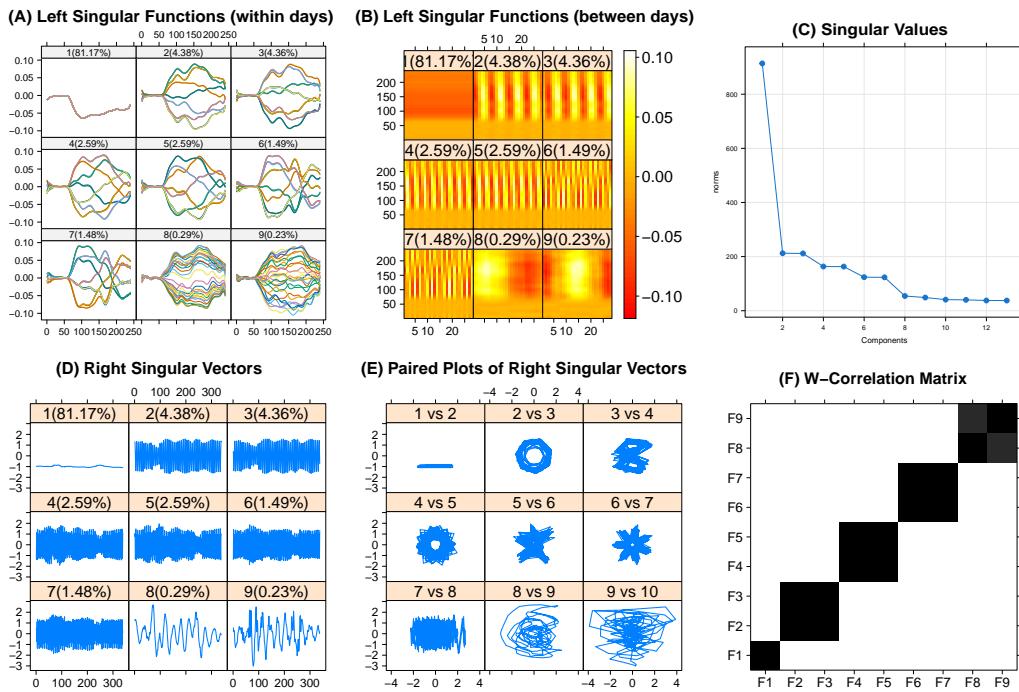
- `values`: A numeric vector of singular values.
- `L`: The specified window length.
- `N`: The length of the functional time series.
- `Y`: The original `funts` object.

The generic `plot(·)` is developed for the `fssa` class to help the user make decisions on how to do the grouping stage of FSSA/MFSSA. This method provides a complete set of visualization tools for the user to check the quality of the decomposition stage. These SSA-type plots encompass various types, each providing unique insights into the data:

- "`values`": Plots the singular values (default).
- "`paired
- "wcor": Generates a plot of the W-correlation matrix for the reconstructed objects.
- "vectors": Displays the right singular vectors, aiding in the detection of period length.
- "lcurves": Showcases the left singular functions, assisting in period length detection.
- "lheats": Heatmap plots of left singular functions are available, designed for funts variables observed over one or two-dimensional domains. These plots are valuable for identifying meaningful patterns.
- "periodogram": Periodogram plots of the right singular vectors can be generated, which help detect the frequencies of oscillations in functional data.`



**Figure 6:** An NDVI snapshot from the city of Jambi, Indonesia.



**Figure 7:** (A): Line plot of the left singular functions, used to identify periodic and trend components; (B): Heatmap of the left singular functions; (C): Scree plot of the singular values often used for grouping; (D): Right singular vectors, used to identify periodic and trend components; (E): Paired plots of the right singular vectors, used for grouping and identifying periodicity in the FTS; (F): Weighted correlation (W-correlation) matrix used for grouping.

While efforts have been made to align these plots in `Rfssa` with the non-functional versions available in the `Rssa` package, there are some fundamental differences. Notably, "lheats" and "lcurves" are novel plot types introduced in the `Rfssa` package for functional data. Additionally, "paired" and "vectors" types are developed based on the right singular vectors. All these plot types utilize the `lattice` graphics engine. The following example codes generate these plots for the Callcenter dataset.

### Example: performing decomposition stage of FSSA

In the rest of the paper, we will use the pre-generated `fnts` class object datasets such as `Callcenter` which are included within the package. These ready-to-use datasets serve as practical examples and templates, allowing users to test the FSSA procedure without the need to start from scratch with data preprocessing. As mentioned previously, our approach involves performing FSSA with a lag window of  $L = 28$ . We then generate a variety of SSA-type plots, as illustrated in Figure 7, using the following code:

```
# Load the Callcenter dataset
data("Callcenter")
# Perform FSSA:
fssa_results <- fssa(Callcenter, L = 28)
# FSSA plots:
plot(fssa_results, d = 9, type = "lcurves",
      main = "(A) Left Singular Functions (within days)")
plot(fssa_results, d = 9, type = "lheats",
      main = "(B) Left Singular Functions (between days)")
plot(fssa_results, d = 13, main = "(C) Singular Values")
plot(fssa_results, d = 9, type = "vectors",
      main = "(D) Right Singular Vectors")
plot(fssa_results, d = 10, type = "paired",
      main = "(E) Paired Plots of Right Singular Vectors")
plot(fssa_results, d = 9, type = "wcor",
      main = "(F) W-Correlation Matrix")
```

The W-correlation matrix in Figure 7(F) is built by measuring the correlation between FTS that are

reconstructed using the grouping of indices by setting  $m = r$  (see the discussion on grouping in section 2). We also have that Figure 7(E) is a plot of successive right singular vectors against one another. Using Figures 7(C, E-F), we identify four groups in the Callcenter data such that  $I_1 = \{1\}$ ,  $I_2 = \{2, 3\}$ ,  $I_3 = \{4, 5\}$ , and  $I_4 = \{6, 7\}$ . Specifically, as discussed in Golyandina et al. (2001), periodic components in FTS typically exhibit a rank of 2, consisting of pairs of harmonic elementary components (sine and cosine functions with the same frequency). To identify these pairs, we can examine pairwise scatterplots of the right singular vectors, as illustrated in Figure 7(E). In such scatterplots, components with identical frequencies, amplitudes, and phases form points that lie along a circular path. Additionally, the number of vertices in the resulting regular T-vertex polygon corresponds to the periodicity of the component. In Figure 7(E), subplots 2 vs. 3, 4 vs. 5, and 6 vs. 7 reveal harmonic factors with frequencies of  $1/7$ ,  $8/14$ , and  $4/7$ , respectively. For a more detailed exploration of extracting meaningful components from the extracted signal, additional references in the SSA literature, such as Golyandina and Zhitlavsky (2013), can be consulted. In addition, using Figures 7(A-B, D-E), we see clear weekly periodic patterns captured in the decomposition, for example, the seven distinct curves found in various subplots of Figure 7(A) and the seven corners seen in subplot 2 vs. 3 of Figure 7(E). For interested readers, we provide further results for the decomposition stage and the fssa object in the GitHub repository of the Rfssa package (<https://github.com/haghbinh/Rfssa>).

## 5 Reconstruction and forecasting

After obtaining an object of class fssa, the user may then choose to perform reconstruction using the freconstruct(.) function or perform forecasting using fforecast(.). The reconstruction and forecasting functions both return a list of funts objects with length  $m$  (number of groups). We note that even though it is common to perform forecasting using a combination of groups that best reconstruct the original signal, the user may try forecasting using several different combinations of groups.

### 5.1 Reconstruction

We start by reconstructing the Callcenter data using the grouping suggested from the FSSA decomposition. The following code implements the reconstruction methodology and gives the plots of the reconstruction in Figure 1.

```
# Define groups and their labels
groups <- list(1, 2:3, 4:5, 6:7, 1:7)
group_labels <- c("(B) First Group",
                 "(C) Second Group",
                 "(D) Third Group",
                 "(E) Fourth Group",
                 "(F) Extracted Signal")
# Perform FSSA reconstruction
reconstructed_data <- freconstruct(fssa_results, groups)
# Create and visualize plots
for (i in 1:length(groups)) {
  print(plotly_funts(reconstructed_data[[i]], main = group_labels[i],
                     xticklocs = xlab, xticklabels = xtloc))
}
```

One may observe the mean behavior (from group  $I_1$ ) in Figure 1(C), and the weekly behaviors (from groups  $I_2, I_3$  and  $I_4$ ) in Figures 1(D-F). Note that the weekly trajectories are more well-separated in  $I_4$  as opposed to the groups  $I_2$  and  $I_3$ . We consider the last group,  $I_5 = \{1, \dots, 7\}$ , as the set of indices corresponding to the leading eigentriples that capture more than 98% of the variation in the signal, to reconstruct the original FTS in Figure 1(B).

### 5.2 The fforecast Class

As previously mentioned, in addition to performing reconstruction of the FTS, the package offers the capability to perform forecasting using an object of class fssa. The constructor for this class, i.e., the fforecast(.) function, accepts the following arguments:

- U: An object of class fssa holding the decomposition.
- groups: A list of numeric vectors where each vector is used for reconstruction and forecasting.
- len: An integer representing the desired length of the forecasted FTS.

- **method:** A character string specifying the type of forecasting to perform, with options including:
  - "recurrent" for FSSA R-forecasting.
  - "vector" for FSSA V-forecasting.
- **only.new:** A logical argument, when set to TRUE, returns only the forecasted FTS, otherwise the entire FTS is returned.

The `fforecast()` function returns an S3 class named `fforecast`, which has been introduced to encapsulate the output of the function. This class is designed to provide a more organized and intuitive structure for handling FTS data. The `fforecast` class includes the following attributes:

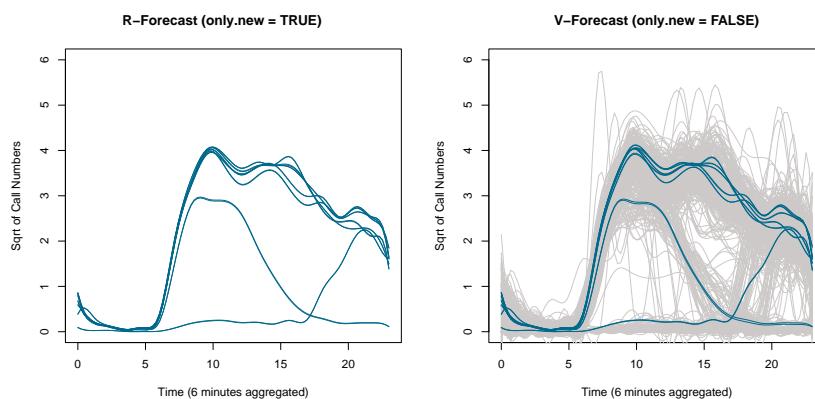
- **original\_funts:** This attribute stores the original FTS object, allowing users to maintain a clear reference to the original data.
- **predicted\_time:** Stores the forecast time index.
- **groups:** Contains a list of numeric vectors, where each vector includes indices of elementary components of a group used for reconstruction and forecasting.
- **method:** A character string specifying the type of forecasting performed.

To streamline user interactions further, we have developed a `print()` method for the `fforecast` class, making it more convenient to view and assess forecasted FTS data. To illustrate these enhancements, consider the continuation of the Callcenter data example in the following:

```
# Perform FSSA R-forecasting
pr_R <- fforecast(U = fssa_results, groups = c(1:3), len = 14, method = "recurrent")
# Perform FSSA V-forecasting
pr_V <- fforecast(U = fssa_results, groups = list(1,1:7), len = 14, method = "vector",
                    only.new = FALSE)
plot(pr_R, main = 'R-Forecast (only.new = TRUE)')
plot(pr_V, main = 'V-Forecast (only.new = FALSE)')
print(pr_V)

# FSSA Forecast (fforecast) class:
# Groups: List of 2
# : num 1
# : int [1:7] 1 2 3 4 5 6 7
# Prediction method: vector
# Predicted series length: 14
# Predicted time: Date[1:14], format: "2000-01-01" "2000-01-02" ...
# -----The original series-----
# Functional time series (funts) object:
# Number of variables: 1
# Length: 365
# Start: 10592
# End: 10956
# Time: Date[1:365], format: "1999-01-01" "1999-01-02" "1999-01-03" ...
```

The resulted figures are shown in Figure 8. The next example will forecast the Callcenter FTS one



**Figure 8:** `plot()` method of `fforecast` class.

week into the future, based on the first 358 days of the year, by leveraging the FSSA R-forecasting and V-forecasting methods as the results that had been given in Figure 2.

```

# Define the data length
N <- Callcenter$N
U1 <- fssa(Callcenter[1:(N-7)], 28)
# Perform recurrent forecasting using FSSA
fore_R = fforecast(U1, groups = list(1:7), method = "recurrent", len = 7)[[1]]
# Perform vector forecasting using FSSA
fore_V = fforecast(U1, groups = list(1:7), method = "vector", len = 7)[[1]]
# Extract the true call data
true_call <- Callcenter[(N-7+1):N]
# Define weekdays and colors
wd <- c('Sunday', 'Monday', 'Tuesday', 'Wednesday', 'Thursday', 'Friday', 'Saturday')
clrs <- c("black", "turquoise4", "darkorange")
argvals <- seq(0, 23, length.out = 100);
par(mfrow = c(1,7), mar = c(0.2, 0.2, 2, 0.2));
# Iterate over the days of the week
for(i in 1:7) {
  plot(true_call[i], col = clrs[1], ylim = c(0, 5.3),
        lty = 3, yaxt = "n", xaxt = "n", main = wd[i]);
  plot(fore_R[i], col = clrs[2], lty = 2, add = TRUE);
  plot(fore_V[i], col = clrs[3], lty = 5, add = TRUE);
}
legend("top", c("Original", "R-forecasting", "V-forecasting"), col = clrs,
       lty = c(3, 2, 5));

```

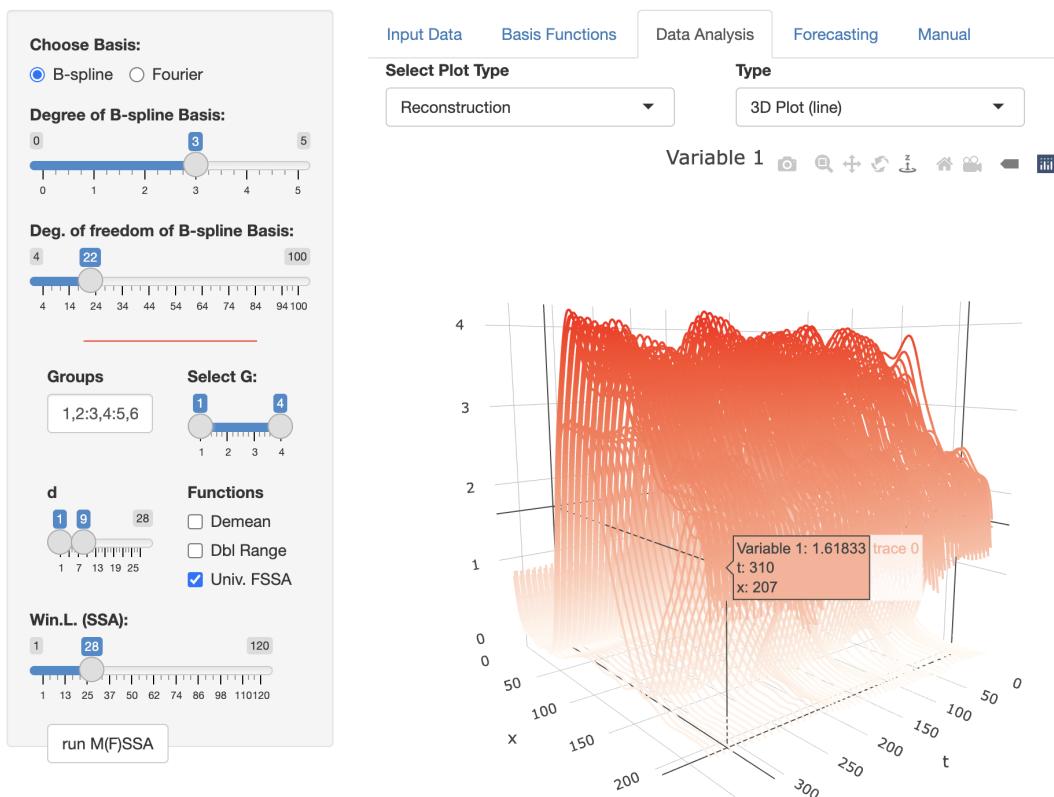
More information on these results may be found in the associated literature of [Haghbin et al. \(2021\)](#) and [Trinka et al. \(2023\)](#).

## 6 A shiny application

The [Rfssa](#) package contains shiny apps for both FSSA and MFSSA methods. Those shiny apps can be called using the `launchApp()` function, and are also available at <http://sctc.mscs.mu.edu/fssa.htm> and <http://sctc.mscs.mu.edu/mfssa.htm>. Here we present the features of the MFSSA app since the FSSA app would be a special case of that. The MFSSA shiny app was developed to visualize and extract the information related to MFTS in a non-parametric framework. The proposed shiny app provides a friendly GUI for users to implement the [Rfssa](#) functionalities and even compare the results with the non-functional version ([Rssa](#)). Figure 9 provides a snapshot of the features available in the MFSSA shiny app. At the top of the sidebar panel, users may specify the basis functions (B-spline or Fourier) and the associated degrees of freedom to represent the funts object. Those basis functions can be visualized in the sub-panel 'Basis Functions' under the main panel. The remaining inputs in the sidebar will be used in other sub-panels. Specifically, 'Groups' is an input box for the third step of the MFSSA algorithm. Each group is specified via a vector (e.g. '`c(1,2,4)`' or '`'1:3'`') and separated from other groups by a comma ('`,`'). The slider 'd' is used to specify the dimensions used in MFSSA (scree, W-correlation, paired, singular vectors & functions, and periodogram plots). The checkbox (a) 'Demean' is used to subtract the mean to obtain mean-zero functions; (b) 'Dbl Range' is used to extend the y-axis to cover all potential mirror functions (e.g. sometimes FPCs may get multiplied by a negative sign); and (c) 'Univ. FSSA' to compare the MFSSA results with marginal FSSA ones respectively. The 'Win.L.' slider specifies the window lengths for the MSSA and the MFSSA. The 'run M(F)SSA' button runs MSSA and MFSSA using the specified parameters for the given dataset. In general, for the sidebar, the top inputs (above the red line) are mostly to describe the basis functions. The bottom inputs (below the red line) are used to specify SSA and FSSA parameters. The main panel includes five sub-panels. Here we briefly describe the features in each of these sub-panels:

- 'Input Data': In this sub-panel, the user can either (a) use the functional datasets available in the [Rfssa](#) package (e.g., `Callcenter` data or remote sensing datasets); (b) simulate MFTS (see [Haghbin et al. 2021](#), for details on simulation setup); or (c) upload any arbitrary FTS matrix (where FTS are given in common grid points and are represented in the columns of the data matrix) to provide the dataset and then analyze it.
- 'Basis Functions': As described before, we illustrate the basis functions selected by the user in this sub-panel.
- 'Data Analysis': In this sub-panel, the user can call various tools to
  - visualize the MFTS.
  - obtain the optimal number of basis functions based on the GCV criteria.

## MFSSA Illustration



**Figure 9:** Snapshot of the MFSSA shiny app

- select a variety of outputs under MSSA and MFSSA, that includes scree plot, W-correlation plot, paired plots, singular vectors plots, periodogram plots, singular functions (heat or regular plots), and reconstruction of FTS using different types of plots (heat, regular, 3Dline and 3Dsurface). An example of the 3Dline reconstruction plot for the Callcenter data is given in Figure 9.
- ‘Forecasting’: This sub-panel would be accessible after the user runs the MFSSA procedure, and it includes the functionalities of R-forecasting and V-forecasting algorithms.
- ‘Manual’: This sub-panel provides a brief instruction manual to use the MFSSA shiny app.

## 7 Summary and conclusion

In summary, **Rfssa** is a pioneering package that brings the power of SSA to the realm of functional time series, offering novel techniques for decomposition, reconstruction, multivariate analysis, and functional forecasting. Its flexible data representation and functional context for SSA make it a valuable addition to the CRAN ecosystem, providing unique capabilities not readily available in other packages. Notably, the package offers extensive capabilities for analyzing FTS/MFTS data, allowing joint analysis of smoothed curves and image data across different dimensional domains. The implementations of the methodologies in the package have been optimized for speed by leveraging the functionalities of **RcppEigen** and **RSpectra** R packages, along with custom C++ code. By utilizing the **Rfssa** package, researchers and practitioners can easily apply advanced FSSA-based techniques to their data, yielding informative results that can significantly enhance decision-making across various applied domains. The intuitive nature and computational efficiency of the package make it a valuable asset for the FTS analysis toolkit.

## Acknowledgments

The authors would like to express their sincere gratitude to the anonymous reviewers for their valuable feedback and constructive comments, which greatly contributed to the improvement of this work.

Additionally, we would like to acknowledge the significant contributions of Dr. S. Morteza Najibi during the development stages of the first version of the `Rfssa` package.

## References

- C. Bouveyron. *funFEM: Clustering in the Discriminative Functional Subspace*, 2021. URL <https://CRAN.R-project.org/package=funFEM>. R package version 1.2. [p82]
- M. de Carvalho and G. Martos. *ASSA: Applied Singular Spectrum Analysis*, 2020. URL <https://CRAN.R-project.org/package=ASSA>. R package version 2.0. [p83]
- M. de Carvalho and A. Rua. Real-time nowcasting the us output gap: Singular spectrum analysis at work. *International Journal of Forecasting*, 33(1):185–198, 2017. [p83]
- M. Febrero-Band and M. O. de la Fuente. Statistical computing in functional data analysis: The R package `fda.usc`. *Journal of Statistical Software*, 51(4):1–28, 2012. doi: 10.18637/jss.v051.i04. [p82]
- A. Gajardo, S. Bhattacharjee, C. Carroll, Y. Chen, X. Dai, J. Fan, P. Z. Hadjipantelis, K. Han, H. Ji, C. Zhu, H.-G. Müller, and J.-L. Wang. *fdatapace: Functional Data Analysis and Empirical Dynamics*, 2022. URL <https://CRAN.R-project.org/package=fdatapace>. R package version 0.5.9. [p82]
- J. Goldsmith, F. Scheipl, L. Huang, J. Wrobel, C. Di, J. Gellar, J. Harezlak, M. W. McLean, B. Swihart, L. Xiao, C. Crainiceanu, P. T. Reiss, Y. Chen, S. Greven, L. Huo, M. G. Kundu, S. Y. Park, D. L. Miller, A.-M. Staicu, E. Cui, and R. Li. *refund: Regression with Functional Data*, 2023. URL <https://CRAN.R-project.org/package=refund>. R package version 0.1.32. [p82]
- N. Golyandina and A. Zhigljavsky. *Singular Spectrum Analysis for Time Series*. Springer Science & Business Media, Berlin, Heidelberg, 2013. [p83, 93]
- N. Golyandina, V. Nekrutkin, and A. A. Zhigljavsky. *Analysis of Time Series Structure: SSA and Related Techniques*. Chapman and Hall/CRC, Boca Raton, FL, 2001. [p85, 93]
- N. Golyandina, A. Korobeynikov, A. Shlemov, and K. Usevich. Multivariate and 2D extensions of the Rssa package. *Journal of Statistical Software*, 67(2):1–78, 2015. doi: 10.18637/jss.v067.i02. URL <http://dx.doi.org/10.18637/jss.v067.i02>. [p83]
- N. Golyandina, A. Korobeynikov, and A. Zhigljavsky. *Singular Spectrum Analysis with R*. Springer Berlin Heidelberg, 2018. [p83]
- H. Haghbin, S. Morteza Najibi, R. Mahmoudvand, J. Trinka, and M. Maadooliat. Functional singular spectrum analysis. *Stat*, page e330, 2021. doi: <https://doi.org/10.1002/sta4.330>. e330 STAT-20-0240.R1. [p83, 84, 85, 90, 95]
- H. Haghbin, J. Trinka, S. M. Najibi, and M. Maadooliat. *Rfssa: Functional Singular Spectrum Analysis*, 2023. URL <https://CRAN.R-project.org/package=Rfssa>. R package version 3.0.2. [p82]
- C. Happ-Kurz. Object-oriented software for functional data. *Journal of Statistical Software*, 93(5):1–38, 2020. doi: 10.18637/jss.v093.i05. [p82]
- H. Hassani and R. Mahmoudvand. Multivariate singular spectrum analysis: A general view and new vector forecasting approach. *International Journal of Energy and Statistics*, 01(01):55–83, 2013. doi: 10.1142/S2335680413500051. URL <https://doi.org/10.1142/S2335680413500051>. [p83]
- R. A. Horn and C. R. Johnson. *Matrix Analysis*. Cambridge University Press, 2012. [p83]
- R. Hyndman and H. L. Shang. *ftsa: Functional Time Series Analysis*, 2023. URL <https://CRAN.R-project.org/package=ftsa>. R package version 6.3.0. [p82]
- M. Maadooliat, J. Z. Huang, and J. Hu. Integrating data transformation in principal components analysis. *Journal of Computational and Graphical Statistics*, 24(1):84–103, 2015. [p85]
- J. O. Ramsay and B. W. Silverman. *Functional Data Analysis*. Springer Series in Statistics. New York, NY, 2005. ISBN 038740080X. URL <http://0-search.ebscohost.com.libus.csd.mu.edu/login.aspx?direct=true&db=cat06952a&AN=mul.b2395232&site=eds-live>. [p82, 83]
- J. O. Ramsay, H. Wickham, S. Graves, and G. Hooker. *fda: Functional Data Analysis*, 2023. URL <https://CRAN.R-project.org/package=fda>. R package version 6.1.4. [p82]

- H. L. Shang and R. Hyndman. *rainbow: Rainbow Plots, Bagplots and Boxplots for Functional Data*, 2022. URL <https://CRAN.R-project.org/package=rainbow>. R package version 3.7. [p82]
- S. Tavakoli. *ftsspec: Spectral Density Estimation and Comparison for Functional Time Series*, 2015. URL <https://CRAN.R-project.org/package=ftsspec>. R package version 1.0.0. [p82]
- J. Trinka, H. Haghbin, and M. Maadooliat. Multivariate functional singular spectrum analysis: A nonparametric approach for analyzing multivariate functional time series. In *Innovations in Multivariate Statistical Modeling: Navigating Theoretical and Multidisciplinary Domains*, pages 187–221. Springer, 2022. [p83, 84, 85]
- J. Trinka, H. Haghbin, H. L. Shang, and M. Maadooliat. Functional time series forecasting: Functional singular spectrum analysis approaches. *Stat*, 12(1):e621, 2023. doi: 10.1002/sta4.621. URL <https://doi.org/10.1002/sta4.621>. [p83, 84, 86, 95]

*Hossein Haghbin*

*Artificial Intelligence and Data Mining Research Group, ICT Research Institute  
Faculty of Intelligent Systems Engineering and Data Science, Persian Gulf University  
Boushehr, Iran  
(ORCID 0000-0001-8416-2354)  
haghbin@pgu.ac.ir*

*Jordan Trinka*

*Department of Mathematical and Statistical Sciences, Marquette University,  
Wisconsin, USA  
(ORCID 0000-0001-9118-5781)  
jordantrinka4@hotmail.com*

*Mehdi Maadooliat*

*Department of Mathematical and Statistical Sciences, Marquette University,  
Wisconsin, USA  
(ORCID: 0000-0002-5408-2676)  
mehdi.maadooliat@mu.edu*

# mcmsupply: An R Package for Estimating Contraceptive Method Market Supply Shares

by Hannah Comiskey and Niamh Cahill

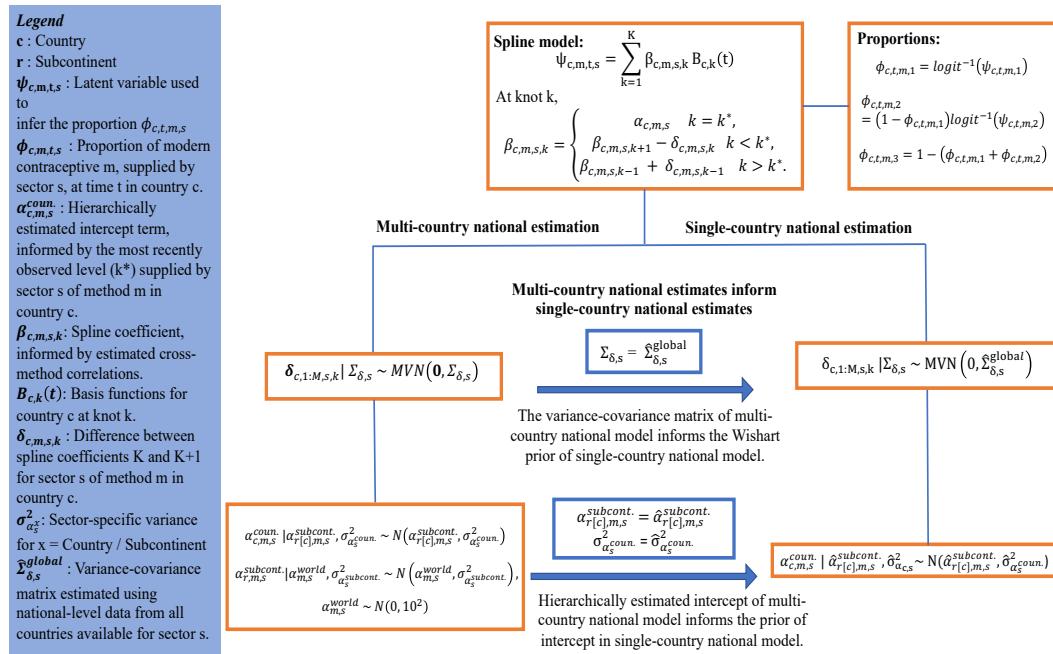
**Abstract** In this paper, we introduce the R package `mcmsupply` which implements Bayesian hierarchical models for estimating and projecting modern contraceptive market supply shares over time. The package implements four model types. These models vary by the administration level of their outcome estimates (national or subnational estimates) and dataset type utilised in the estimation (multi-country or single-country contraceptive market supply datasets). The `mcmsupply` package contains a compilation of national and subnational level contraceptive source datasets, generated by Integrated Public Use Microdata Series (IPUMS) and Demographic and Health Survey (DHS) microdata. We describe the functions that implement the models through practical examples. The annual estimates and projections with uncertainty of the contraceptive market supply, produced by `mcmsupply` at a national and subnational level, are the first of their kind. These estimates and projections have diverse applications, including acting as an indicator of family planning market stability over time and being utilised in the calculation of estimates of modern contraceptive use.

## 1 Introduction

Family Planning 2030 (FP2030) is a ‘global movement dedicated to advancing the rights of people everywhere to access reproductive health services safely and on their own terms’ (FP2030, 2021). One step towards achieving this goal is to quantify *how* people are accessing their modern contraceptive supplies. To date, obtaining estimates of modern contraceptive supply shares in low- and middle-income countries has relied on large-scale national surveys like the Demographic and Health Surveys (DHS). However, these DHS are not annually available and in practice, most countries carry out DHS every 3 to 5 years approximately, with some countries having fewer surveys than this (DHS, 2023). In previous work, we described a model that provides probabilistic estimates of the contraceptive supply share over time with uncertainty and examined the model performance at the national administration division for countries that are participating in FP2030 and have varying amounts of DHS data available (Comiskey et al., 2024). The original modern contraceptive supply share model (`mcmsupply` model) relies on splines, informed by cross-method correlations, to capture temporal variation combined with a hierarchical modelling approach to estimate country-level parameters. Using a multi-country dataset, the original `mcmsupply` model produces estimates of contraceptive supply market shares at the national level for all countries simultaneously. In this paper, we extend the model to estimate supply shares using input data from a single-country and to also include estimation at the subnational administrative division. At the time of writing, this package and models are the first of their kind to estimate and project modern contraceptive method supply shares at the national and subnational administration levels using Demographic and Health survey data.

For the remainder of the paper, we will refer to the original modern contraceptive supply share model as the multi-country national `mcmsupply` model. The single-country `mcmsupply` model uses a scaled-down version of the multi-country approach. It borrows strength from the multi-country model using modular model runs with informative priors placed on key parameters to provide precise outcome estimation, even in the absence of data for a particular contraceptive method. Modularization in Bayesian analysis describes the process within a statistical model where information is restricted to flow only from the prior to the likelihood. Thus, preventing the ‘contamination’ of key parameters from suspect data (Lunn et al. (2009) as cited in Plummer (2015)). In the context of our problem, parameters estimated within the multi-country model are used to inform the priors of the single-country (either national or subnational administrative division) models. This approach prevents spurious parameter estimates due to a lack of data for some locations. To summarise how the single-country and multi-country models are connected to each other, Figure 1 depicts this modelling relationship at the national administration level. The main differences between the multi-country and single-country approaches is that for a single-country model, we only have data for one country (at the national or subnational administration division) and the country-level (in the national model) or subnational-level (in the subnational model) population parameters are informed by estimates from the corresponding multi-country model. In contrast to this, the multi-country model uses national or subnational level data (depending on your administrative level of interest) from many countries simultaneously to estimate model parameters and the country-level (in the national model)

or subnational-level (in the subnational model) population parameters are estimated hierarchically. The **mcmsupply** package contains vignettes that consider case studies of both single-country and multi-country model estimation approaches at the national and subnational administration levels.



**Figure 1:** A flow chart to illustrate the relationship between the multi-country and single-country national estimation approaches. The median estimates of the multi-country national model parameters are used as informative priors in the single-country national model. A legend describing each element of the schematic is located in the blue box.

The need for a single-country version of the **mcmsupply** model arose for two primary reasons. Firstly, there was a demand for improved computational efficiency while ensuring that model accuracy remains uncompromised when generating projections of modern contraceptive method supply shares for a single-country. Secondly, the option to accommodate custom data, like incorporating a new survey dataset, was sought to be included in the estimation process for users who require a more tailored analysis.

The inclusion of the functionality to estimate supply shares at the subnational level in the **mcmsupply** package was spurred by the growing interest in subnational estimation among the family planning community (New et al., 2017; Mercer et al., 2019; Li et al., 2019). The decentralization of family planning services produces a more equitable and efficient service; however, it also shifts the responsibility of service delivery to lower-level organisations that may not have the capacity to carry-out the role (Williamson et al., 2014). Providing subnational-level estimates can lead to a clearer understanding of localised user preferences, localised user access to family planning commodities and a measure of the true stability contraceptive supply market at a smaller geographic scale (Bossert and Beauvais, 2002). These subnational estimates may also be used as part of a localised temperature check for progress towards the FP2030 goals, which include increasing access to contraceptive methods (Munoz and Ali, 2022).

This paper introduces the R package **mcmsupply** for estimating and projecting the contraceptive supply share at the national and subnational administration divisions using multi-country or single-country datasets. The package uses **tidyverse** throughout to carry out data manipulation and visualisation tasks. The graphics are produced using **ggplot2** and **dplyr** is used for any data manipulation procedures. The workflow of the package aims to follow a Bayesian workflow (Gelman et al., 2020). The users are encouraged to explore the data, the model inputs, and the model outputs using diagnostics and visualisations. The individual elements of the workflow are made as accessible as possible through the built-in functions, such as the ‘get\_data’ and ‘plot\_estimates’ functions, and the code provided in the vignettes to test the convergence of the model parameters. To assess convergence, we consider the R-hat values of the model parameters using the plot function of **rjags**, as well as the individual parameter trace plots (Vehtari et al., 2021a). Figure 2 shows a summary of the different ways users can use the R functions and input data in **mcmsupply** to carry out model fitting and estimation. The national data contained in the package is derived from the DHS microdata (ICF, 2004) while the subnational data is derived from the IPUMS DHS datasets (Herger Boyle et al., 2022). In the

'Implementation and operation' section, we review the operation and implementation requirements of the **mcmSupply** R package. In the 'Data' section, a detailed description of the data used within the **mcmSupply** R package is provided, as well as an explanation of the data pre-processing functions `get_data` and `get_modelinputs`. 'The estimation process' section describes the model fitting and visualisation functions within the **mcmSupply** R package. A basic overview of the process model is provided in 'Overview of the mcmSupply process models' with an explanation of some key modelling parameters. The 'Model fitting' section explains the `run_jagsmodel` function. The `run_jagsmodel` function fits models in a Bayesian framework using the JAGS (Just Another Gibbs Sampler) software and produces estimates with uncertainty for the administration level and dataset type of choice (Hornik et al., 2003). The 'Model output' section describes the `plot_estimates` function that takes the model estimates and visualises them using the R package `ggplot2` (Wickham, 2016). The 'Use cases' section discusses five use cases for modelling modern contraceptive supply shares using the **mcmSupply** R package. Finally, the conclusions are presented.

## 2 Implementation and operation

**mcmSupply** contains pre-processing functions to clean and prepare raw input data for model fitting at the administrative level of choice (national or subnational) using the dataset type of choice (multi-country or single-country). This R package includes functions to fit Bayesian hierarchical models using the model inputs. Functionality for post-processing and visualisation of the model estimates are also included. The model fitting process described uses Just Another Gibbs Sampler (JAGS). JAGS uses Markov Chain Monte Carlo (MCMC) sampling to produce model estimates for Bayesian hierarchical models (Hornik et al., 2003). For installation, both R (>=3.5.0) and JAGS (>=4.0.0) are required. JAGS can be downloaded at <https://sourceforge.net/projects/mcmc-JAGS/files/JAGS/>. **mcmSupply** interacts with JAGS using the wrapper functions supplied by the R package **R2jags** (Su and Yajima, 2021). The **mcmSupply** package dependencies are listed in the package `DESCRIPTION` file and will be automatically installed upon installing the main package. There are no minimum RAM, CPU, or HARDDRIVE requirements apart from what is necessary to store model runs, which varies case-by-case. This software was run on a MacBook Air using macOS 13.1 with a 1.6 GHz Dual-Core Intel Core i5 processor and 8GB of memory.

## 3 Data

### 3.1 Input data

The **mcmSupply** package contains the following input data sources:

- Survey data for the national and subnational modern contraceptive supply shares between 1990 and 2022 for a selection of countries participating in the FP2030 initiative.
- Variance-covariance data on the logit scale corresponding to the observed public and commercial medical modern contraceptive supply shares at the national level between 1990 and 2022 for a selection of countries participating in the FP2030 initiative.
- Estimated correlations between the rates of change in method supply shares in the public and private sectors at both the national and subnational administrative divisions. These are derived based on the method outlined in Comiskey et al. (2023).
- Global and subcontinental parameter estimates obtained from the multi-country model run at the national and subnational level.

The `inst/data-raw` folder of the **mcmSupply** package contains sample R code that was used to create the model inputs in the case of the covariance arrays, correlations and parameters. This enables users to recalculate their own single-country model parameters and correlations should they wish to do so. The code for the creation of the main input contraceptive supply source data is provided on the **mcmSupply** [github](#) page but the raw data for these datasets cannot be provided. The raw data may be accessed by users through an application to the DHS program for national level data, or the IPUMS program, for subnational level data. Help files for the contraceptive supply source datasets can be accessed by using the command `?mcmSupply::national_FPsource_data` and `?mcmSupply::subnat_FPsource_data`. The national contraceptive supply source data has data for 33 countries (including participants and non-participants of FP2030) between 1990 and 2022. The subnational contraceptive supply source data has data for 256 provinces across 24 countries (all participating in FP2030) between 1990 and 2020. Table 1 is a sample of 6 rows of the subnational contraceptive supply source data.

**Table 1:** The subnational contraceptive supply source data used in the subnational estimation models. Country and Region list the name of the country and province the observation relates to. The Method column lists the type of the contraceptive method supplied. The mid-year when the survey was collected is listed in average\_year. The supply sector is found in sector\_categories. The observed proportion and standard error are found in proportion and SE.proportion. The number of respondent making up each observation are listed in n.

	Country	Region	Method	average_year	sector_categories	proportion	SE.proportion	n
9512	Zimbabwe	Midlands	OC Pills	2010.5	Commercial_medical	0.1601738	0.0355722	40
9513	Zimbabwe	Midlands	OC Pills	2010.5	Other	0.1018543	0.0259602	29
9514	Zimbabwe	Midlands	OC Pills	2010.5	Public	0.7379720	0.0449851	194
9515	Zimbabwe	Midlands	OC Pills	2015.5	Commercial_medical	0.1929466	0.0395410	60
9516	Zimbabwe	Midlands	OC Pills	2015.5	Other	0.0428227	0.0166039	13
9517	Zimbabwe	Midlands	OC Pills	2015.5	Public	0.7642307	0.0387566	209

Lastly, data on country classification, ISO codes and area groupings is provided in the dataset `Country_and_area_classification` (Table 2). The help file for this dataset can be accessed via the R command `?mcmsupply::Country_and_area_classification`.

**Table 2:** The `Country_and_area_classification` dataset is the Track20 project country and area classification data according to the United Nations Statistical Division, standard country or area codes for statistical use (M49). This data set is how we classify each country in subcontinental regions. The name of the country, the International Organization for Standardization (ISO) code for each country, the continent, sub-continent are listed. Details on whether or not a country is defined as a developing, located in Sub-Saharan Africa and the status of its participation in FP2020 (now FP2030) are also provided.

Country or area	ISO Code	Major area	Region	Developed region	Least developed country	Sub-Saharan Africa	FP2020
Afghanistan	4	Asia	Southern Asia	No	Yes	No	Yes
Albania	8	Europe	Southern Europe	Yes	No	No	No
Algeria	12	Africa	Northern Africa	No	No	No	No
American Samoa	16	Oceania	Polynesia	No	No	No	No
Andorra	20	Europe	Southern Europe	Yes	No	No	No
Angola	24	Africa	Middle Africa	No	Yes	Yes	No

### 3.2 Data pre-processing

In `mcmsupply`, the data processing occurs in two steps: First, the raw input data is retrieved and preliminary cleaning to the dataset is completed. Secondly, the cleaned data is processed to provide the model inputs for the Bayesian hierarchical model. This two-step process removes any black-box element to the model fitting process and allows the user to review the data at both stages and refer to it later when considering model outputs. The first step involves the `get_data` function. This function retrieves the raw data from the stored contraceptive supply source dataset, does data cleaning and processing to address any issues with missing data and regional naming inconsistencies. Its arguments are summarized in Table 3. The `get_data` function carries out several data pre-processing tasks. For example, it will remove observations where two or more sectors are missing. It also checks to see if observations across all three sectors sum to one and will replace missing observations with 0 when the total does sum to one. The function transforms the raw data away from the (0,1) boundary as parameter estimates that lie on the distribution limits cause issues with convergence during the model estimation process. Finally, the `get_data` function imputes any missing standard errors using a binomial approximation. To use the `get_data` function, the user first defines whether they wish to use national or subnational level administrative data via the `national` argument. National level data is accessed when the `national` argument is set to TRUE. When subnational administrative data is required, the user sets `national` to FALSE. Similarly, the user defines whether they want to use multi-country estimation with data from multiple countries or single-country estimation with data from a single-country via the `local` argument. The default setting for the `local` argument is FALSE. This induces a multi-country estimation, where the outcomes for all the countries in the contraceptive supply source dataset will be estimated simultaneously. In the event of single-country estimation, the user sets `local` to TRUE and indicates their country of interest via the `mycountry` argument. The names of the countries listed in the package data can be found in the `country_names` dataset. The help file for this dataset can be accessed via the R command `?mcmsupply::country_names`. The `fp2030` argument controls whether to include countries that are participating in the FP2030 initiative or not. The default includes only the named FP2030 countries (see `country_names`) in the dataset. There is the optional functionality to include a custom dataset. This allows the user to run the model on data outside of that stored within the package. The `surveydata_filepath` is a character string that denotes

the location of the custom dataset. The file must meet a series of internal checks on file type, column names, suitable data ranges and missing data. When a custom dataset is supplied to `get_data`, the function carries out the checks and alerts the user to any differences between what is expected and what has been supplied. If `surveydata_filepath` is left as `NULL`, by default the function uses the stored `national_FPsouce_data` or `subnat_FPsouce_data`, depending on what administrative level the user has specified via the `national` argument (Table 3). The `get_data` function returns a list containing the cleaned data and a list of arguments supplied to the function. Storing the arguments of the `get_data` function allows the set-up information to flow without requiring the user to repeatedly supply the same arguments for each step of the modelling process.

**Table 3:** The arguments of the `get_data` function. The purpose of this function is to retrieve and clean the Demographic and Health Survey (DHS) data or custom user supplied data for use in supply share estimation. The ‘Argument’ column names the function component. Data type describes the argument. Description explains the purpose of the argument and any default entries.

Argument	Data type	Description
<code>national</code>	Character	<p>It indicates whether the user is interested in using data at the national or subnational administration level.            This is a binary TRUE or FALSE argument.            Default is TRUE which retrieves national level data, while FALSE retrieves subnational data.</p>
<code>local</code>	Character	<p>It indicates whether the user is interested in using data for a single population or not.            This is a binary TRUE or FALSE argument.            Default is FALSE.</p>
<code>mycountry</code>	Character	<p>This is the name of the country you wish to do single-country estimation for. The data will only be returned for this country.            Default is <code>NULL</code>.</p>
<code>fp2030</code>	Character	<p>It indicates whether the user is interested in using only countries participating in FP2030.            This is a binary TRUE or FALSE argument.            Default is TRUE.</p>
<code>surveydata_filepath</code>	Character string	<p>Pathway to the location of the custom dataset.            When left as <code>NULL</code>, the function automatically uses the stored datasets.            Default is <code>NULL</code>.</p>

Step two of the data pre-processing is the `get_modelinputs` function. This function takes the cleaned data from the previous step and repackages it into suitable inputs for the model implementation. The arguments of the function are summarised in Table 4. This function uses the arguments set in the `get_data` function as well as additional parameters for the model. These parameters include the year the user wishes to begin their estimation at and the year they finish on. In the `mcmcsupply` package, the models use basis splines (B-splines), to capture the complexities in variation of the contraceptive supply source data over time. The use of splines for the estimation of demographic indicators is growing in popularity. Previous studies used small-area estimation models with splines to capture complex shapes of demographic spatio-temporal data (Ugarte et al., 2010, 2009). In recent years, many international health organisations have also used splines for the estimation of key demographic indicators. These include the estimation and projection of under-5 mortality for United Nations Children’s Fund (UNICEF) (Sharrow et al., 2019) and the estimation of excess morality due to Covid-19 for the World Health Organisation (WHO) (Knutson et al., 2023). We build on these previous studies to use penalised regression splines in the `mcmcsupply` R package. B-splines use basis-functions to create piecewise cubic polynomials. The number of basis functions that are fit to the data is determined by the number of knots. Knots are the locations along the x-axis where the piecewise polynomials of the B-splines join. As you increase the number of knots in the basis functions, the B-splines give a tighter fit to the data. Similarly, if you decrease the number of knots in the basis, you will get a smoother fit to your data. In the `mcmcsupply` package, the user may alter the number of knots (`nsegments`) used in the basis functions. The default number of knots is 12, as was used in Comiskey et al. (2023). This equates to a knot approximately every 3.5 years. Through validation, it was found that having fewer knots dis-improved model fit, while more knots did not significantly improve model fit. Like the `get_data` function, this function returns a list containing the model inputs and the function arguments.

**Table 4:** The arguments of the `get_modelinputs` function. The purpose of this function is to get the model inputs for the JAGS model used for supply share estimation. In the table, the argument name and data type of the argument is stated, a description of the argument and any default values is then provided.

Argument	Data type	Description
<code>startyear</code>	Numeric	The year you wish to start your estimation at.
<code>endyear</code>	Numeric	The year you wish to finish your estimation at.
<code>nsegments</code>	Numeric	The number of knots you wish to include in your basis functions. Default is 12.
<code>raw_data</code>	List	The output of the <code>get_data</code> function, which includes a list of the function arguments used and the cleaned contraceptive supply source data.

## 4 The estimation process

The `mcmcsupply` R package contains four Bayesian models, each of which aims to estimate and project contraceptive method supply shares over time with uncertainty. These models vary by the administration level of their outcome estimates (national or subnational estimates) and dataset type utilised in the estimation (multi-country or single-country contraceptive market supply datasets). A full mathematical description of all the models contained within `mcmcsupply` is described in the supplementary material.<sup>1</sup>. A summary of each of the parameters and their role within each model can be found in Table 5 while a visual summary of the national model, using both multi-country and single-country inputs, can be found in Figure 1.

### 4.1 Brief model overview

The outcome of interest is the components of a compositional vector  $\phi_{q,t,m}$ , which captures the proportion of contraceptive method m, at time t, in population q supplied across the public and private sectors.

Let,

$$\phi_{q,t,m} = (\phi_{q,t,m,1}, \phi_{q,t,m,2}, \phi_{q,t,m,3}), \quad (1)$$

where,

$\phi_{q,t,m,s}$  is the proportion supplied by the public sector ( $s=1$ ), the private commercial medical sector ( $s=2$ ) and the other private sector ( $s=3$ ) of modern contraceptive method m, at time t, in population q (national or subnational).

Figure 3 shows the model set up for the national level models. A similar approach is taken when estimating modern contraceptive method supply at the subnational administration level. For each model within the `mcmcsupply` package, the latent variable  $\psi_{q,m,t,s}$  relies on a spline to capture the underlying process that generates the data, on the logit scale, for sector s, in year t, for method m and population q (depending on the administration level of interest).

$$\psi_{q,t,m,1} = \sum_{k=1}^K \beta_{q,m,1,k} B_{q,k}(t), \quad (2)$$

where,

$\beta_{q,m,1,k}$  is the  $k^{th}$  spline coefficient for sector s, method m in population q.

$B_{q,k}(t)$  is the  $k^{th}$  basis function fit to the data for population q.

We assume that in population q, for method m and sector s, the value of spline coefficient at knot index  $k^*$ , aligning with the year  $t^*$ , the most recent survey available, is  $\alpha_{q,m,s}$ . By doing this, we are assuming that the  $\alpha_{q,m,s}$  parameter will act as the spline coefficient for the reference spline at  $k^*$ . We are then able to calculate the remaining spline coefficients using a random walk model of order 1 on spline coefficients from the reference index ( $k^*$ ) using the penalised  $\delta_{q,m,s}$ .

<sup>1</sup>At the time of publication it is expected that this material will be made available via another open access venue.

Let,

$$\beta_{q,m,s,k} = \begin{cases} \alpha_{q,m,s} & k = k^*, \\ \beta_{q,m,s,k+1} - \delta_{q,m,s,k} & k < k^*, \\ \beta_{q,m,s,k-1} + \delta_{q,m,s,k-1} & k > k^*, \end{cases} \quad (3)$$

where,

$\alpha_{q,m,s}$  is the most recently observed supply share, on the logit scale, for sector s, method m, in population q. This parameter is estimated hierarchically. The geographical set-up of this estimation process adapts to match the administrative level of interest. For example, the subnational multi-country models contain an additional layer of geography (world > subcontinent > country > province) in the hierarchical set-up that the national models don't have, which accounts for the subnational administration levels.

$k$  is the knot index along the set of basis splines  $B_{q,k}(t)$

$k^*$  is the index of the knot that corresponds with  $t^*$ , the year index where the most recent survey occurred in population q.

$\delta_{q,m,s,k-1}$  is the first order difference between spline coefficients  $\beta_{q,m,s,K}$  and  $\beta_{q,m,s,K-1}$ . These reflect the changes in method supply shares over time. Within each sector, the first-order differences are assumed to be correlated between methods. These correlations were estimated using a maximum *a posteriori* estimator for the correlation matrix first described in Azose and Raftery, (2018) and adapted for method supply shares in Comiskey et al. (2023). These estimated correlations are available as data for both the national and subnational models. Please see the Data section of this paper for more details.

**Table 5:** This table summarises the purpose of each parameter within each of the four models described in the `mcmsupply` R package. The four types of models are listed in the 'Model type' column. For each model, the parameters listed in the 'Parameter name' column are the default parameters monitored when the argument `jagsparams=NULL` is used within the `run_jags_model` function. 'Parameter purpose' explains the role each parameter plays within the estimation process and the notation can be linked directly to the 'Model overview' section of this paper.

Model type	Parameter name	Parameter purpose
Multi-country national	P	The method supply share proportions across all countries, methods and sectors.
	beta.k	The set of spline coefficients for $(\beta_{c,m,s,k})$ for all countries, methods and sectors at each knot.
	alpha_cms	The intercept term $(\alpha_{c,m,s})$ for all countries, methods and sectors.
	delta.k	The first order differences between spline coefficients $(\delta_{c,m,s,k})$ across all knots for all countries, methods and sectors
Single-country national	P	The method supply share proportions for a given country, across all methods and sectors.
	alpha_cms	The intercept term $(\alpha_{c,m,s})$ for the country of interest c, across all methods and sectors.
	inv.sigmas_delta	The precision matrix used in the multivariate normal prior of $\delta_{c,1:M,s,k}$
	beta.k	The set of spline coefficients $(\beta_{c,m,s,k})$ for the country of interest c, across all methods and sectors at each knot.
Multi-country subnational	alpha_pms	The intercept term $(\alpha_{p,m,s})$ for all subnational provinces, methods and sectors.
	alpha_cms	The intercept term $(\alpha_{c,m,s})$ for all countries, methods and sectors. $\alpha_{c[p],m,s}$ is the expected value of $\alpha_{p,m,s}$ .
	inv.sigmas_delta	The precision matrix used in the multivariate normal prior of $\delta_{p,1:M,s,k}$
	tau_alpha_pms	The sector-specific precision associated with $\alpha_{p,m,s}$
	beta.k	The set of spline coefficients $(\beta_{p,m,s,k})$ across all subnational provinces, methods and sectors at each knot.
	delta.k	The first order differences between spline coefficients across all knots $(\delta_{p,m,s,k})$ for all provinces, methods and sectors.
Single-country subnational	P	The method supply share proportions for the subnational provinces in the country of interest, for all methods and sectors.
	alpha_pms	The intercept term $(\alpha_{p,m,s})$ for the subnational provinces in the country of interest, for all methods and sectors.
	beta.k	The set of spline coefficients $(\beta_{p,m,s,k})$ for the subnational provinces in the country of interest, across all methods and sectors at each knot

**Table 6:** The arguments of the `run_jags_model` function. The purpose of this function is to run the Bayesian hierarchical models stored within the `mcmcsply` package for either single- or multi-country datasets at the administration level of interest. The argument name and data type of the argument is stated, a description and any default values of the argument are then provided.

Argument	Data type	Description
<code>jagsdata</code>	List	The output of the <code>get_modelinputs</code> function. A list of the initial set-up arguments and the JAGS inputs required using the data.
<code>jagsparams</code>	Vector	A string vector of model parameters to be monitored. NULL invokes a standard vector to be used Default is NULL.
<code>n_iter</code>	Numeric	Number of iterations you wish to run your JAGS model for. Default is 80000.
<code>n_burnin</code>	Numeric	Number of burn-in samples you wish to run your JAGS model for. Default is 10000
<code>n_thin</code>	Numeric	Number of samples you wish to thin your JAGS sample by. Default is 35.

**Table 7:** The arguments of the `plot_estimates` function. The purpose of this function is to plot the data alongside the model estimates so that users can visualise their estimated method supply shares. The argument name and data type of the argument is stated, a description of the argument is then provided.

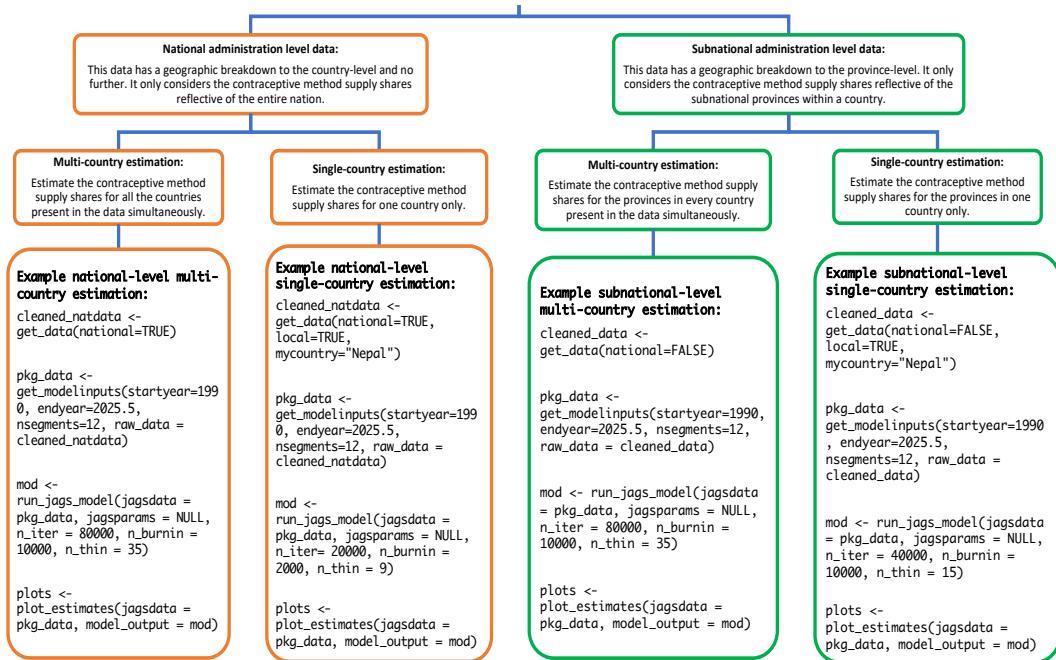
Argument	Data type	Description
<code>jagsdata</code>	List	A list of the initial set-up arguments and the JAGS inputs required using the data retrieved in the <code>get_modelinputs</code> function.
<code>model_output</code>	List	The object assigned to store the list of MCMC results and estimate summary output from the <code>run_jags_model</code> function.

## 4.2 Model fitting

The `run_jags_model` function fits the selected JAGS model to the supplied data and returns a list of MCMC samples and point summaries for the time-period and locations of interest. Initial set-up arguments (`national`, `local`, `mycountry`) are inherited from the specification of the previous `get_modelinputs` function. The additional inputs of this function are summarised in Table 6. The `jagsdata` argument of this function is a list of initial set-up arguments and JAGS model inputs gathered from the `get_modelinputs` function. `jagsparams` is a vector of strings that name the model parameters the user wishes to monitor within the JAGS model. The default is NULL. When `jagsparams` = NULL, the function will refer to a stored vector of parameters to monitor (Table 5). The JAGS parameters of `n_iter` = 80000, `n_burnin` = 10000 and `n_thin` = 35 ensure that the model has converged and that the final posterior sample size is 2000 samples. The function `get_point_estimates` takes the chains produced by the JAGS model and estimates the median, 95% credible intervals and 80% credible intervals. The `get_point_estimates` function runs automatically inside the `run_jags_model` function and returns the point summaries as part of the `run_jags_model` function output. The `run_jags_model` function returns a list containing the JAGS output of the model and the point summaries for the estimates.

## 4.3 Model output

The user then runs the function `plot_estimates`. This function visualises the point estimates with uncertainty alongside the data using the initial set up inputs of the `get_modelinputs` function and the output of the `run_jags_model` function (Table 7). The `plot_estimates` function returns a list of `ggplot2` objects, one for each country (when using the national model) or subnational region (when using the subnational model).



**Figure 2:** A flow chart to illustrate the decision processes that lead to the different estimation types within the mcmsupply package. The first decision is with respect to the administrative level of the estimates you wish to create – they may be either national level or subnational level. An explanation of each division is found on the first row of the figure. The second decision is with respect to the number of countries you wish to estimate – the user may estimate the proportions for all the countries at once or only one country. An explanation of each in the context of the specific administrative division is found on the second row of the figure. A set of sample functions used to estimate and plot the estimates for each modelling option are located on the third row of the figure

## 5 Use cases

### 5.1 Case 1: Estimating contraceptive method supply shares at the national administration level for multiple countries simultaneously

The first use case describes how the user can estimate modern contraceptive method supply shares at the national administrative level over time for multiple countries at once (i.e., using the multi-country national model). This use case is described in `national_multicountry_mod` found in the vignettes folder. This vignette takes approximately 12 hours to run when using the recommended JAGS arguments, on a machine with 1.6 GHz Dual-Core Intel Core i5 processor and 8GB of RAM. The `national_FPsOURCE_data` dataset contains observations for 30 countries. The user begins by accessing the `national_FPsOURCE_data` dataset through the `get_data` function with the argument `national=TRUE`, and the remaining arguments sets to their default values, to indicate that they are interested in national-level data for the FP2030 countries present in the data.

```
cleaned_natdata <- get_data(national = TRUE)

dplyr::glimpse(cleaned_natdata$mydata)
```

Next, this data is supplied to the `get_modelinputs` function. This function reshapes the data into a list of inputs for the JAGS model. At this point, the user must indicate the start and end years they wish to estimate between. The `n_segments` argument controls how many knots will be used in the basis functions. The default number of segments is 12. Lastly, the cleaned national data from the `get_data` function is provided to the `get_modelinputs` function via the `raw_data` argument.

```
pkg_data <- get_modelinputs(startyear = 1990,
                             endyear = 2025.5,
                             nsegments = 12,
                             raw_data = cleaned_natdata)
```

This list of data and model inputs is then fed into the JAGS model via the `run_jags_model` function. In this instance, the user wishes to monitor the default set of parameters within the JAGS model. Therefore, they set the `jagsparams` argument to `NULL`, which invokes the function to use the default list. The parameters for running the JAGS model are set via the `n_iter`, `n_burnin` and `n_thin` arguments. For demonstration purposes in the code below, these arguments are set to low values. However, for optimal model convergence we recommend increasing the iterations to 80,000, the burn-in period to 10,000, and the thinning to 35. Please note that for these increased values of arguments, the model run time is approximately 12 hours. As part of the `run_jags_model` function, the median and 80% and 95% credible intervals for the estimates are calculated. To assess convergence, we considered the convergence diagnostic  $\hat{R}$ , as well as the individual parameter trace plots (Vehtari et al., 2021b).

```
mod <- run_jags_model(jagsdata = pkg_data,
                       jagsparams = NULL,
                       n_iter = 80,
                       n_burnin = 10,
                       n_thin = 2)

plot(mod$JAGS)

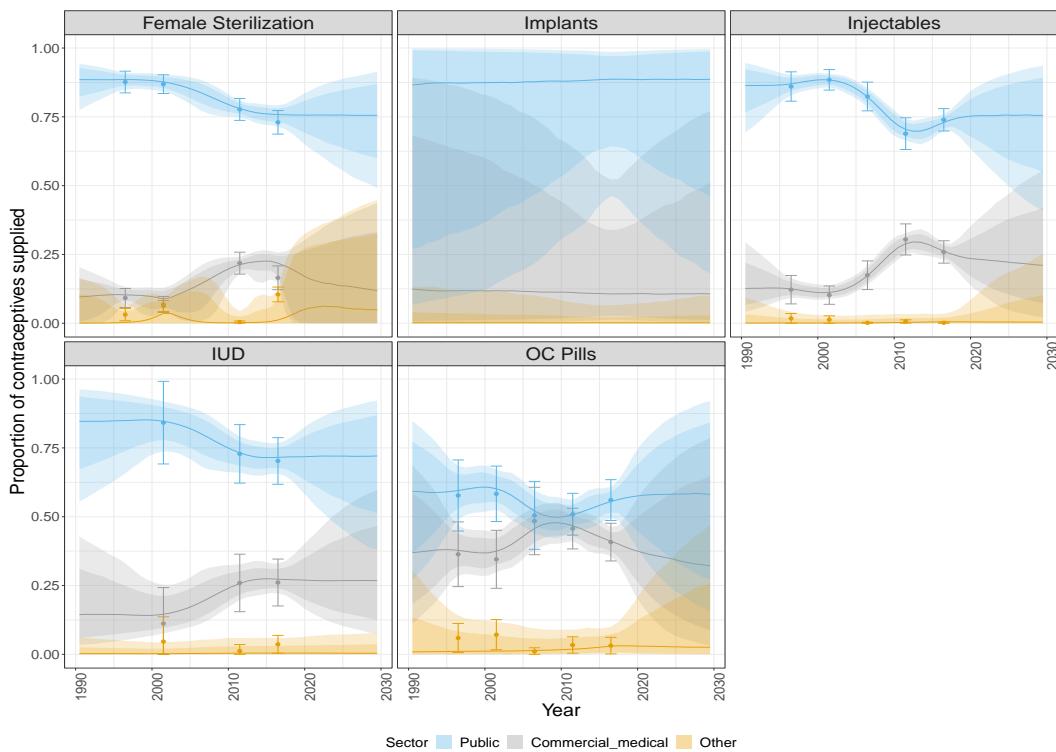
sample_draws <- tidybayes::tidy_draws(mod$JAGS$BUGSoutput$sims.matrix)

var <- sample_draws %>% dplyr::select(.chain, .iteration, .draw, `P[1,2,1,1]`) %>%
  dplyr::mutate(chain = rep(1:2, each=mod$JAGS$BUGSoutput$n.keep)) %>%
  dplyr::mutate(iteration = rep(1:mod$JAGS$BUGSoutput$n.keep, 2))

ggplot2::ggplot(data=var) +
  ggplot2::geom_line(ggplot2::aes(x=iteration, y=`P[1,2,1,1]`, color=as.factor(chain)))
```

The final JAGS model output and the summary estimates are returned as a list. Finally, these summary estimates are visualised via the `plot_estimates` function using the R package `ggplot2` (Figure 3).

```
plots <- plot_estimates(jagsdata = pkg_data,
                        model_output = mod)
```



**Figure 3:** The plotted posterior point estimates for each of the three sectors (public in blue, private commercial medical in grey, and private other in gold) for Nepal at the national administrative level over time with the 80% and 95% uncertainty interval denoted as shaded regions. Time in years is on the x-axis and the proportion of each contraceptive method supply share is on the y-axis. The survey observations are plotted as points with their associated standard error, plotted as vertical lines. This plot was produced using a multi-country set up of the `mcmcSupply` functions.

If the user wishes to pull out specific estimates for a given country and year, they may do so using the `pull_estimates` function. The model object, country name and the year of interest are supplied, the function then returns a tibble of the corresponding estimates.

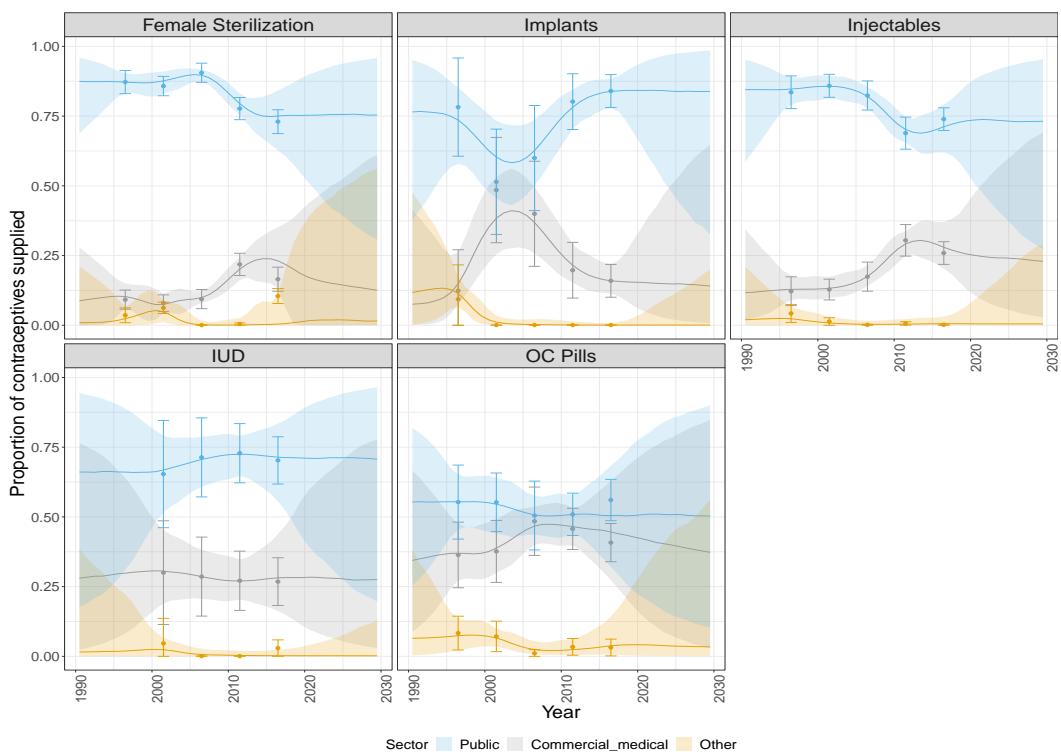
```
estimates_2018 <- pull_estimates(model_output = mod, country = 'Nepal', year=2018)
```

## 5.2 Case 2: Estimating contraceptive method supply shares at the national administration level for a single-country

This case considers when the estimates at the national administration level are required for only one country. Rather than running a multi-country model, which takes several hours, a quicker alternative is the single-country approach, which takes only a few minutes. The main difference between the multi-country and single-country model outputs is that the model estimates of the single-country models have slightly larger uncertainty. This is especially evident where data is absent for a particular method. For example in Figure 3, the width of the 95% credible intervals over time for implants estimated by the multi-country national model are smaller than those estimated in the single-country model (Figure 4). The arguments `local` and `mycountry` control the single-country estimation models in `mcmcSupply`. The user begins as by retrieving the data for Nepal only using the `get_data` function. They set `local=TRUE` and specify which country they are interested in by setting `mycountry=Nepal`.

```
cleaned_natdata <- get_data(national = TRUE,
                             local = TRUE,
                             mycountry = "Nepal")
```

These arguments are the only discernible differences in the commands for users, the rest of the workflow is as described above in Case 1. A complete workflow for this use case can be found in `vignettes/national_singlecountry_mod`. The single-country model produces model estimates that align with those estimated by the multi-country estimation model.



**Figure 4:** The plotted posterior point estimates for each of the three sectors (public in blue, private commercial medical in grey, and private other in gold) for Nepal at the national administrative level over time with the 80% and 95% uncertainty interval denoted as shaded regions. Time in years is on the x-axis and the proportion of each contraceptive method supply share is on the y-axis. The survey observations are plotted as points with their associated standard error, plotted as vertical lines. This plot was produced using a single-country set up of the mcmcsupply functions.

### 5.3 Case 3: Estimating contraceptive method supply shares at the subnational administration level for multiple countries simultaneously

The use case for estimating the contraceptive supply shares via a multi-country model for the sub-national administration division is given by the vignette subnational\_multicountry\_models. This vignette takes approximately 24 hours to run on a machine with 1.6 GHz Dual-Core Intel Core i5 processor and 8GB of RAM. The dataset contains observations for 225 subnational divisions, across 23 countries. The user begins by calling the multi-country dataset at the subnational administration level via the `national` argument.

```
cleaned_subnatdata <- get_data(national = FALSE)

#> Using preloaded dataset!

#> Joining with `by = join_by(Country, average_year)`
#> Joining with `by = join_by(Country)`
```

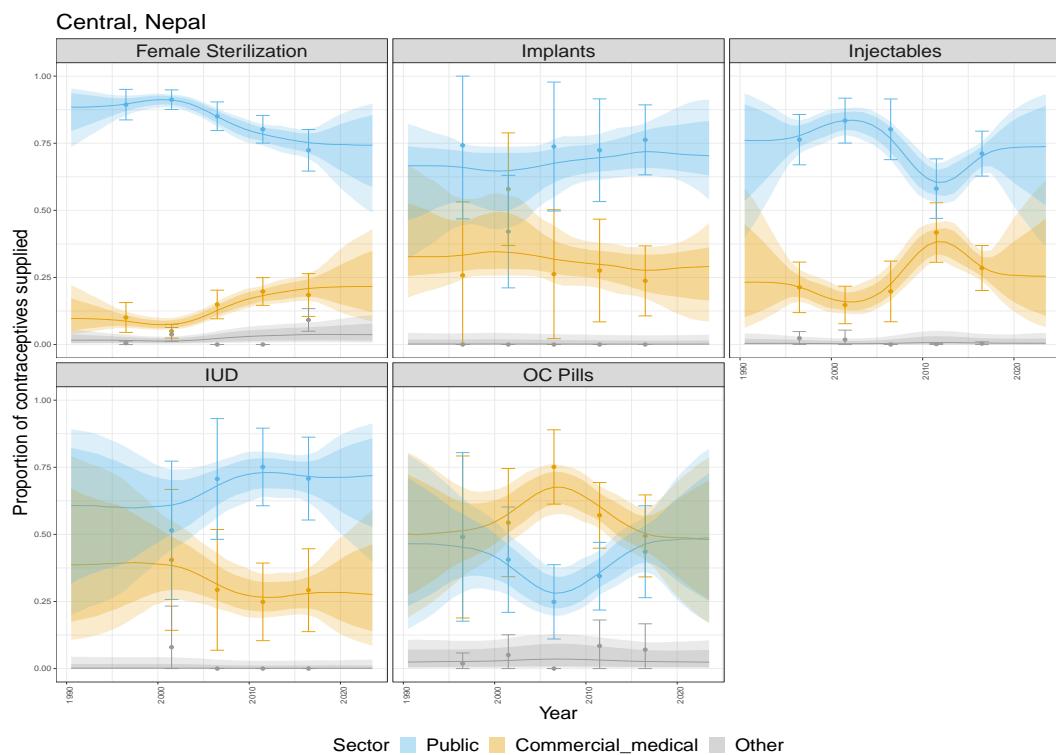
The remaining workflow is the same as described above in Case 1 and is not shown here. A complete workflow for this use case can be found in vignettes/subnational\_multicountry\_models.

### 5.4 Case 4: Estimating contraceptive method supply shares at the subnational administration level for a single-country

This use case is for considering use of the single-country model at the subnational administrative division. The user begins by retrieving the data for Nepal using the `get_data` function by setting the arguments `national=FALSE`, `local=TRUE` and `mycountry='Nepal'`.

```
cleaned_data <- get_data(national = FALSE,
                          local = TRUE,
                          mycountry = "Nepal")
```

As in the previous use cases, the JAGS model is run and point summaries are calculated via the `run_jags_model` function. The visualisations for the subnational regions of Nepal are returns as a list via the `plot_estimates` function. An example of these visualisations is given in Figure 5, where the estimated median method supply shares for Central region of Nepal are plotted with 80% and 95% credible intervals over time in each of the methods.



**Figure 5:** The plotted posterior point estimates for each of the three sectors (public in blue, private commercial medical in grey, and private other in gold) taken from the subnational single-country population for the Central subnational region of Nepal over time with both 80% and 95% uncertainty denoted as corresponding shaded regions. Time in years is on the x-axis and the proportion of each contraceptive method supply share is on the y-axis. The survey observations are plotted as points with their associated standard error, plotted as vertical lines.

### 5.5 Case 5: Estimating contraceptive method supply shares at the national/subnational administration level for a single-country using custom data

It is possible to include custom datasets when estimating contraceptive method supply shares at either the national or subnational administration level. The set-up for using custom data is very similar to the above processes, with a small difference in the data retrieval step using the `get_data` function. When using a custom dataset, the user defines the location of the `.xlsx` file containing the custom data.

```
cleaned_data <- get_data(national = FALSE,
                          local = TRUE,
                          surveydata_filepath =
                            "inst/data-raw/my_custom_data_good.xlsx",
                          mycountry = "Ethiopia")
```

The file must be in `.xlsx` format and match the layout of either the `national_FPsouce_data` or `subnat_FPsouce_data` datasets, depending on your desired administration level. The `get_data` function will carry out a series of internal checks to ensure the custom data matches the stored data layout. If the custom data is not suitable, the `get_data` function will return an error and a message to the user describing the issue with the custom data. A description of each column is given in Table 8. Once the custom data checks are complete and passed, the regular workflow for fitting an `mcmcSupply` model continues. The JAGS model inputs are retrieved using the `get_modelinputs` and the JAGS model is fit using `run_jags_model`. The summary estimates are plotted using the `plot_estimates`

**Table 8:** The required columns when supplying custom data to the `get_data`. The column name is given, the type of data expected in each column along with a description of the column. The ‘Required’ column indicates whether the column belongs in national or subnational custom data files. The `get_data` function will carry out a series of internal checks to ensure the custom data matches the stored data layout.

Column name	Data type	Description	Required
Country	Character	The name of the country you have the custom data for.	National and subnational
Region	Character	The name of the subnational area you have the custom data for	Subnational
Method	Character	The name of the contraceptive method you have the custom data for	National and subnational
average_year	Numeric	The year that the observation was collected in.	National and subnational
sector_categories	Character	The supply sector of the contraceptive method you have the custom data for.	Subnational
sector_category	Character	The supply sector of the contraceptive method you have the custom data for.	National
proportion	Numeric	The proportion of a given contraceptive method supplied by the observation.	National and subnational
SE.proportion	Numeric	The associated standard error of the proportion supplied by the observation.	National and subnational
n	Numeric	The sample size used to estimate the observation.	National and subnational

function without any further changes to the workflow. A vignette of how to run the subnational model using a custom dataset can be found [vignettes/subnational\\_singlecountry\\_customdata\\_models](#).

## 6 Discussion

In this paper we have introduced the package `mcmSupply`. The primary purpose of this package is to estimate and project modern contraceptive method supply shares from the public, private commercial medical and private other sectors with uncertainty, for a selection of countries participating in the FP2030 initiative. The `mcmSupply` package produces estimates within the period of available Demographic and Health Survey data for a given country, as well as projections beyond the most recent data point. The package uses either stored DHS survey data or custom user-supplied survey data as inputs to the modelling process. The package implements four model types. These models vary by the administration level of their outcome estimates (national or subnational administration level) and dataset type utilised in the estimation (multi-country or single-country contraceptive market supply datasets). The modelling framework uses penalised splines to capture temporal nature of the data. These splines utilise correlations between changes in method supply shares that exist within the data. Using splines informed by cross-method correlations allows us to capture the complex shape of the data without over-fitting. Bayesian hierarchical estimation is another key element of this model estimation process. We take advantage of the geographical nature of the data, such that the expected sector, province (subnational level) or country (national level), method-specific supply shares are informed based on the existing geographical structures in the data. This promotes information sharing across locations, and better inform estimates in areas where smaller amounts of data are present. Case studies illustrate how to use the `mcmSupply` for each of the four potential modelling routes, as well as how to estimate contraceptive method supply shares using custom user-supplied data.

The `mcmSupply` package has many benefits to users. Firstly, it is the first of its kind to produce annual estimates with uncertainty for monitoring contraceptive method supply shares over time at the national and subnational levels. On average, most countries carry out DHS surveys every 5-6 years, but in some instances the wait time between surveys may be even longer ([DHS, 2023](#)). The family planning community use contraceptive method supply shares to evaluate the stability and sustainability of a given country’s contraceptive market ([Bradley and Shiras, 2022](#)). Contraceptive method supply share estimates are also pivotal in effectively managing contraceptive commodity supply-chains through a ‘total market approach’ (TMA). A TMA approach to the family planning supply market seeks to engage the public, private commercial, and other sectors in a country to increase family planning users access to vital information, products, and services ([Moazzam, 2015; through the Private Sector , SHOPS](#)). Prior to this package, individuals who required contraceptive method supply shares for a given country relied on estimates from the most recent DHS survey in the country, regardless of its age. The `mcmSupply` package alleviates this issue and provides the family planning community annual estimates with uncertainty for these contraceptive method supply shares. Secondly, contraceptive method supply shares estimates are not only a stand -alone family planning indicator but are also used in the calculation of an another indicator, estimated modern use (EMU) ([Track20, 2020](#)). EMUs aim to measure the proportion of women, aged 15 to 49 years old , who are currently using any modern method of contraception, and are derived from routinely collected family planning service statistics ([Track20, 2020](#)). Currently, EMU calculations depend on an adjustment that relies on the most recent DHS survey to provide estimates of the contraceptive supply share market in a given country. Now this adjustment can instead rely on the annual estimates and projections with uncertainty produced by the `mcmSupply` package. This can serve to improve the overall accuracy of EMUs with respect to their ability to accurately measure modern contraceptive use. In addition, given the probabilistic nature of the `mcmSupply` estimates the associated uncertainty can be propagated into

the EMU calculations.

The last key benefit of the `mcmSupply` package is the speed at which the user is able to access individual countries supply share estimates. Using the single-country models, at either the national or subnational level, provides users with annual estimates of the method supply shares with uncertainty within minutes. This fast estimation approach is computationally efficient while still producing reliable estimates. Using priors informed by multi-country model subcontinental and country-level median parameter estimates, this modelling approach is robust to spurious parameter estimates even when estimating supply shares for countries with fewer surveys available. This computational efficiency without a loss of model accuracy makes the `mcmSupply` package user-friendly and efficient for regular data analysis.

The `mcmSupply` is not without its limitations. We remark that the implementation of these models in JAGS have not been optimised. The multi-country models at the national and subnational levels take hours to run and are intensive on computer memory and CPU. Hence, improvements such as matrix operations rather than for-loops would greatly improve the computation efficiency of this package. Another limitation of this package is that in methods without any survey information, the uncertainty intervals tend to be large. This is especially evident in the single-country models where in the absence of data for a given method, the uncertainty of the associated estimates is even larger than that of the corresponding multi-country model estimates. These limitations inspire our future work, where we seek to improve the computational efficiency of these models. We would also like to investigate the potential for incorporating additional covariates into the models, such as average method pricing for each sector, to improve the uncertainty of model estimates and projections where no DHS survey data is available. Lastly, we wish to design an R-shiny app that promotes the use of these model estimates among statistical non-experts within the family planning community.

## 7 Conclusions

`mcmSupply` is an R package that estimates the modern contraceptive method supply shares at the national and subnational administrative divisions over time for countries participating in the Family Planning 2030 initiative. The package provides the user an easy and accessible way to produce annual estimates with uncertainty using Bayesian hierarchical penalised spline models with cross-method correlations at the national and subnational administration levels. These annual estimates with uncertainty may act as a stand-alone family planning indicator of the stability of the modern contraceptive supply market or be used to produce alternative family planning indicators, such as estimated modern use (EMUs) using service statistics (Track20, 2020). To the best of our knowledge, the package is the first of its kind to estimate these supply shares at both administrative levels. Using an R package to disseminate this work aligns with the findability, accessibility, interoperability, and reusability (FAIR) principles of scientific data (Wilkinson et al., 2016). The data used in the package is cited and explained thoroughly, the code is commented and easy to understand should a user wish to tweak or review any functionalities, and finally it is reusable by the very nature of the R package.

## References

- T. J. Bossert and J. C. Beauvais. Decentralization of health systems in Ghana, Zambia, Uganda and the Philippines: a comparative analysis of decision space. *Health Policy and Planning*, 17(1):14–31, 03 2002. ISSN 0268-1080. doi: 10.1093/heapol/17.1.14. URL <https://doi.org/10.1093/heapol/17.1.14>. [p100]
- S. E. K. Bradley and T. Shiras. Where Women Access Contraception in 36 Low- and Middle-Income Countries and Why It Matters. *Global Health: Science and Practice*, 10(3), 2022. doi: 10.9745/GHSP-D-21-00525. URL <https://www.ghspjournal.org/content/10/3/e2100525>. [p112]
- H. Comiskey, L. Alkema, and N. Cahill. Estimating the proportion of modern contraceptives supplied by the public and private sectors using a Bayesian hierarchical penalized spline model. *Journal of the Royal Statistical Society Series A: Statistics in Society*, 05 2024. ISSN 0964-1998. doi: 10.1093/jrsssa/qnae051. URL <https://doi.org/10.1093/jrsssa/qnae051>. [p99]
- DHS. The DHS Program - Demographic and Health Survey (DHS), 2023. URL <https://dhsprogram.com/methodology/survey-Types/dHs.cfm>. [p99, 112]
- FP2030. Rights and Empowerment Principles for Family Planning, 2021. URL <https://commitments.fp2030.org/principles>. [p99]

- A. Gelman, A. Vehtari, D. Simpson, C. C. Margossian, B. Carpenter, Y. Yao, L. Kennedy, J. Gabry, P.-C. Bürkner, and M. Modrák. Bayesian workflow. 11 2020. URL <http://arxiv.org/abs/2011.01808>. [p100]
- E. Herger Boyle, M. King, and M. Sobek. IPUMS-Demographic and Health Surveys: Version 9 [dataset], 2022. URL <https://doi.org/10.18128/D080.v9>. [p100]
- K. Hornik, F. Leisch, A. Zeileis, and M. Plummer. JAGS: A Program for Analysis of Bayesian Graphical Models Using Gibbs Sampling, 2003. URL <http://www.ci.tuwien.ac.at/Conferences/DSC-2003/>. [p101]
- ICF. Demographic and Health Surveys (various) [Datasets], 2004. [p100]
- V. Knutson, S. Aleshin-Guendel, A. Karlinsky, W. Msemburi, and J. Wakefield. Estimating global and country-specific excess mortality during the Covid-19 pandemic. *The Annals of Applied Statistics*, 17(2):1353–1374, 2023. doi: 10.1214/22-AOAS1673. URL <https://doi.org/10.1214/22-AOAS1673>. [p103]
- Q. Li, T. A. Louis, L. Liu, C. Wang, and A. O. Tsui. Subnational estimation of modern contraceptive prevalence in five sub-Saharan African countries: A Bayesian hierarchical approach. *BMC Public Health*, 19(1), 2019. ISSN 14712458. doi: 10.1186/s12889-019-6545-3. [p100]
- D. Lunn, N. Best, D. Spiegelhalter, G. Graham, and B. Neuenschwander. Combining MCMC with ‘sequential’ PKPD modelling. *Journal of pharmacokinetics and pharmacodynamics*, 36(1):19–38, 2009. ISSN 1573-8744. doi: 10.1007/S10928-008-9109-1. URL <https://pubmed.ncbi.nlm.nih.gov/19132515/>. [p99]
- L. D. Mercer, F. Lu, and J. L. Proctor. Sub-national levels and trends in contraceptive prevalence, unmet need, and demand for family planning in Nigeria with survey uncertainty. *BMC Public Health*, 19 (1), 2019. ISSN 14712458. doi: 10.1186/s12889-019-8043-z. [p100]
- A. Moazzam. Ensuring contraceptive security through effective supply chains: WHO/RHR/17.09. Technical report, World Health Organization, United Nations Population Fund, 2015. Accessed: 2023-03-30. [p112]
- M.-F. Munoz and M. Ali. Commitment to fp2030, 2022. URL [https://cdn.who.int/media/docs/default-source/reproductive-health/who-srh-fp2030.pdf?sfvrsn=780d62cb\\_3](https://cdn.who.int/media/docs/default-source/reproductive-health/who-srh-fp2030.pdf?sfvrsn=780d62cb_3). [p100]
- J. R. New, N. Cahill, J. Stover, Y. P. Gupta, and L. Alkema. Levels and trends in contraceptive prevalence, unmet need, and demand for family planning for 29 states and union territories in India: a modelling study using the Family Planning Estimation Tool. *The Lancet Global Health*, 5(3):e350–e358, 2017. ISSN 2214109X. doi: 10.1016/S2214-109X(17)30033-5. URL [http://dx.doi.org/10.1016/S2214-109X\(17\)30033-5](http://dx.doi.org/10.1016/S2214-109X(17)30033-5). [p100]
- M. Plummer. Cuts in Bayesian graphical models. *Statistics and Computing*, 25(1):37–43, 2015. ISSN 15731375. doi: 10.1007/S11222-014-9503-Z. URL <https://dl.acm.org/doi/10.1007/s11222-014-9503-z>. [p99]
- D. Sharro, L. Hug, P. Danzhen You, P. Leontine Alkema, R. Black, S. Cousens, T. Croft, V. Gaigbe-Togbe, P. Gerland, M. Guillot, K. Hill, B. Masquelier, C. Mathers, J. Pedersen, K. L. Strong, E. Suzuki, J. Wakefield, W. Neff, U. I.-a. G. f. C. M. G. E. Advisory, and its Technical. National, regional, and global levels and trends in neonatal mortality between 1990 and 2017, with scenario-based projections to 2030: a systematic analysis. *The Lancet Global Health*, 7(6):e710–e720, 2019. ISSN 2214109X. doi: 10.1016/S2214-109X(19)30163-9. URL [http://dx.doi.org/10.1016/S2214-109X\(19\)30163-9](http://dx.doi.org/10.1016/S2214-109X(19)30163-9). [p103]
- Y.-S. Su and M. Yajima. *R2jags: Using R to Run ‘JAGS’*, 2021. URL <https://CRAN.R-project.org/package=R2jags>. R package version 0.7-1. [p101]
- S. H. O. through the Private Sector (SHOPS) Plus. Total market approach compendium, 2016. URL <https://shopsplusproject.org/tmacompendium>. [p112]
- Track20. SS to EMU Tool: Process Guide, 2020. URL [http://www.track20.org/pages/track20\\_tools\\_SS\\_to\\_EMU\\_tool.php](http://www.track20.org/pages/track20_tools_SS_to_EMU_tool.php). [p112, 113]
- Track20. Monitoring Progress in Family Planning Estimated Modern Use (EMU): A New Service Statistics-Based Family Planning Indicator, 2020. URL <http://fpet.track20.org/fpet/>. [p112]

- M. Ugarte, T. Goicoa, A. Militino, and M. Durbán. Spline smoothing in small area trend estimation and forecasting. *Computational Statistics & Data Analysis*, 53(10):3616–3629, 2009. URL <https://EconPapers.repec.org/RePEc:eee:csdana:v:53:y:2009:i:10:p:3616-3629>. [p103]
- M. D. Ugarte, T. Goicoa, and A. F. Militino. Spatio-temporal modeling of mortality risks using penalized splines. *Environmetrics*, 21(3-4):270–289, 2010. doi: <https://doi.org/10.1002/env.1011>. URL <https://onlinelibrary.wiley.com/doi/abs/10.1002/env.1011>. [p103]
- A. Vehtari, A. Gelman, D. Simpson, B. Carpenter, and P. C. Burkner. Rank-Normalization, Folding, and Localization: An Improved (Formula presented) for Assessing Convergence of MCMC (with Discussion)\*†. *Bayesian Analysis*, 16(2):667–718, 2021a. ISSN 19316690. doi: 10.1214/20-BA1221. [p100]
- A. Vehtari, A. Gelman, D. Simpson, B. Carpenter, and P.-C. Bürkner. Rank-normalization, folding, and localization: An improved r<sup>+</sup> for assessing convergence of mcmc (with discussion). *Bayesian Analysis*, 16(2), June 2021b. ISSN 1936-0975. doi: 10.1214/20-ba1221. URL <http://dx.doi.org/10.1214/20-BA1221>. [p108]
- H. Wickham. *ggplot2: Elegant Graphics for Data Analysis*. Springer-Verlag New York, 2016. ISBN 978-3-319-24277-4. URL <https://ggplot2.tidyverse.org>. [p101]
- M. D. Wilkinson, M. Dumontier, I. J. Aalbersberg, G. Appleton, M. Axton, A. Baak, N. Blomberg, J. W. Boiten, L. B. da Silva Santos, P. E. Bourne, J. Bouwman, A. J. Brookes, T. Clark, M. Crosas, I. Dillo, O. Dumon, S. Edmunds, C. T. Evelo, R. Finkers, A. Gonzalez-Beltran, A. J. G. Gray, P. Groth, C. Goble, J. S. Grethe, J. Heringa, P. A. C. t Hoen, R. Hooft, T. Kuhn, R. Kok, J. Kok, S. J. Lusher, M. E. Martone, A. Mons, A. L. Packer, B. Persson, P. Rocca-Serra, M. Roos, R. van Schaik, S. A. Sansone, E. Schultes, T. Sengstag, T. Slater, G. Strawn, M. A. Swertz, M. Thompson, J. V. D. Lei, E. V. Mulligen, J. Velterop, A. Waagmeester, P. Wittenburg, K. Wolstencroft, J. Zhao, and B. Mons. Comment: The FAIR Guiding Principles for scientific data management and stewardship. *Scientific Data*, 3, 2016. ISSN 20524463. doi: 10.1038/sdata.2016.18. [p113]
- T. Williamson, S. Duvall, A. A. Goldsmith, K. Hardee, and R. Mbuya-Brown. The Effects of Decentralization on Family Planning: A Framework for Analysis, 2014. URL [https://www.healthpolicyproject.com/pubs/445\\_FPDecentralizationFINAL.pdf](https://www.healthpolicyproject.com/pubs/445_FPDecentralizationFINAL.pdf). [p100]

*Hannah Comiskey*  
Maynooth University  
Hamilton Institute  
Maynooth, Co.Kildare, Ireland  
ORCID: 0000-0003-0034-6725  
[hannah.comiskey.2015@mumail.ie](mailto:hannah.comiskey.2015@mumail.ie)

*Niamh Cahill*  
Maynooth University

# Pomdp: A Computational Infrastructure for Partially Observable Markov Decision Processes

by Michael Hahsler and Anthony R. Cassandra

**Abstract** Many important problems involve decision-making under uncertainty. For example, a medical professional needs to make decisions about the best treatment option based on limited information about the current state of the patient and uncertainty about outcomes. Different approaches have been developed by the applied mathematics, operations research, and artificial intelligence communities to address this difficult class of decision-making problems. This paper presents the pomdp package, which provides a computational infrastructure for an approach called the partially observable Markov decision process (POMDP), which models the problem as a discrete-time stochastic control process. The package lets the user specify POMDPs using familiar R syntax, apply state-of-the-art POMDP solvers, and then take full advantage of R's range of capabilities, including statistical analysis, simulation, and visualization, to work with the resulting models.

## 1 Introduction

Many important problems require decision-making without perfect information, and where decisions made today will affect the future. For example, in diabetes prevention and care, the primary care provider needs to make decisions about screening, early interventions like suggesting lifestyle modification, and eventually medication for disease management based on a patient's available medical history. Especially, screening and lifestyle modifications need to be used early on to be effective in preventing severe and debilitating diseases later on. This is clearly a difficult problem that involves uncertainty and requires a long-term view. We have studied this problem using the partially observable Markov decision process approach in (Kamalzadeh et al., 2021) and, in the absence of solvers for R, we started the development of the pomdp package described in this paper.

A Markov decision process (MDP) is a discrete-time stochastic control process that models how an agent decides on what actions to take when facing an environment whose dynamics can be adequately modeled by a Markov process that can be affected by the agent's behavior (Puterman, 1994). That is, the environment transitions between a set of states where transition probabilities only depend on the current state and are conditioned on the agent's actions. Over time, the agent receives rewards depending on the actions and the environment's state. The agent's objective is to make a plan that maximizes its total reward earned. A plan can be expressed as a mapping of each possible state to the best action in that state. The best possible plan is often called the optimal policy. For MDP problems, the agent is always aware of the state of the environment and can make decisions directly following such a policy.

A partially observable Markov decision process (POMDP) generalizes the concept of the MDP to model more realistic situations where the agent cannot directly observe the environment's state. Here, the agent must infer the current state using observations that are only probabilistically linked to the underlying state. The agent can form a belief about what states it may be in and update its belief when new observations are made. In this setting, the agent has to base its actions on its current belief. A POMDP can be modeled as a *belief MDP* where the underlying Markov model uses belief states instead of the original states of the environment. While the original state space is typically modeled as a finite set of states, making MDPs readily solvable using dynamic programming, the agent's belief is represented by a probability distribution over the states in the form of a continuous probability simplex and are therefore much more challenging to solve. The volume of the believe space that POMDPs are operating in grows exponentially with the number of underlying states. This is called the curse of dimensionality which means that working with problems with a realistic number of states typically requires the use of approximate algorithms.

Karl Johan Åström first described Markov decision processes with a discrete state space and imperfect information in 1965 (Åström, 1965). The model was also studied by the operations research community where the acronym POMDP was introduced (Smallwood and Sondik, 1973). More recently, the POMDP framework was adapted for automated planning problems in artificial intelligence (Kaelbling et al., 1998). The POMDP framework is a popular choice when a known Markov process can adequately approximate system dynamics and the reward function is known. POMDPs have been successfully applied to model various real-world sequential decision processes. Examples include numerous industrial, scientific, business, medical and military applications where an optimal or near-

optimal policy is needed. This includes important applications like machine maintenance scheduling, computer vision, medical diagnosis, and many more. A detailed review of applications can be found in (Cassandra, 1998b).

While the (PO)MDP framework is used to find an optimal or near optimal policy, given a model of system dynamics, the related class of model-free reinforcement learning algorithms, more specifically, temporal difference learning, Q-learning and its deep learning variations (Sutton and Barto, 2018), learn unknown system dynamics and the reward function directly from interactions with the environment. Reinforcement learning methods typically require observable states and perform a large amount of exploration, where the agent performs sub-optimal actions to learn about the environment. Q-learning and some algorithms are already available in R packages like **ReinforcementLearning** (Proellocos and Feuerriegel, 2020). While these model-free approaches are very powerful for many artificial intelligence applications, they may not be appropriate for situations where experts already possess a reasonable amount of knowledge about the system dynamics and where the cost of sub-optimal actions is very high. For example, the cost of administering the wrong medication in a medical setting due to exploration by a pure reinforcement learning approach may not be acceptable and a model-based approach like a POMDP is more appropriate. The R package described in this paper exclusively focuses on planning with POMDP.

While POMDPs are well studied, the complexity of solving all but very small problems limits its application. Recent spectacular advances in artificial intelligence applications have lead to more interest in POMDPs, as shown in the development of new approximate algorithms and by the frameworks available for various programming languages:

- pomdp-solve (Cassandra, 2015) is a C program to solve POMDPs using exact and approximate solvers.
- APPL (APPL Team, 2022) provides the fast point-based POMDP solver SARSOP in C++.
- ZMDP software (Smith, 2009) implements several approximate value iteration algorithms in C++.
- pyPOMDP (Migge and Stollmann, 2013) is a Python 2.x toolbox for solving POMDPs.
- JuliaPOMDP (JuliaPOMDP Team, 2022) is a set of packages for defining and solving MDPs and POMDPs using the Julia programming language.

R activity around POMDPs has also picked up with the packages **sarsop** (Boettiger et al., 2021) and **pomdpSolve** (Hahsler and Cassandra, 2022) which interface the two popular POMDP solver programs APPL and pomdp-solver.

In this paper, we present **pomdp** (Hahsler, 2023) which was co-developed with **pomdpSolve** (Hahsler and Cassandra, 2022) to provide R users with a consistent and flexible infrastructure for solving and working with POMDPs. The package can be used to work with larger POMDP problems but is limited by the capability of the used solvers. Larger problems also typically lead to very complicated policies which can be executed by an automatic agent but are not very helpful for a human user. This paper focuses on features for smaller problems that yield simpler policies. Such models and policies are better suited for human experts who want to understand the problem and are interested in improved decision making. For example, a medical researcher who tries to develop easy-to-follow guidelines for doctors based on experiments with a POMDP model is looking for a relatively simple and robust model with a simple and understandable policy. This typically means to consider a model with few states and a small number of different observations. For example, we have used the package to study diabetes prevention by creating a very small, simplified model to obtain a policy that is actionable in a primary care setting (Kamalzadeh et al., 2021). A second use of smaller models is in a classroom or self-study setting where the **pomdp** package can be used to demonstrate and study how POMDP models, solvers, and resulting policies work.

## 2 Background for partially observable Markov decision processes

A POMDP is a discrete-time stochastic control process that can formally be described by the 7-tuple

$$\mathcal{P} = (S, A, T, R, \Omega, O, \gamma),$$

where

- $S = \{s_1, s_2, \dots, s_n\}$  is the set of partially observable states of the environment,
- $A = \{a_1, a_2, \dots, a_m\}$  is the set of available actions,
- $T$  describes the system dynamics as the set of transition probabilities  $T(s' | s, a)$  the state transition  $s \rightarrow s'$  conditioned on taking action  $a$ .
- $R : S \times A \times S \rightarrow \mathbb{R}$  is the reward function which can depend on the state transition (previous and new state) and the action,

- $\Omega = \{o_1, o_2, \dots, o_k\}$  is the set of possible observations,
- $O$  defines the probabilistic connection of observations with the reached states  $s'$  as the set of observation probabilities  $O(o | a, s')$  conditioned on the action  $a$  taken to reach  $s'$ , and
- $\gamma \in [0, 1]$  is the discount factor modeling how much the agent prefers immediate rewards over later rewards.

The used notation follows largely (Kaelbling et al., 1998). Several variations of this notation can be found in the literature. Sets are often set in calligraphic font and it is also common to see the observation model denoted by  $Z$  instead of  $O$ .

The control process proceeds in discrete time steps called epochs as follows. At each time epoch  $t$ , the environment is in some unknown state  $s \in S$ . The agent chooses an action  $a \in A$ , which causes the environment to transition to state  $s' \in S$  with probability  $T(s' | s, a)$ . Simultaneously, the agent receives an observation  $o \in \Omega$ , which depends on the action and the new state of the environment following the conditional probability distribution  $O(o | a, s')$ . Finally, the agent receives a reward  $R(s, a, s')$  depending on the transition. This process repeats till a specified time horizon is reached. Often, as in the equation below, an infinite horizon is used. The goal of the agent is to plan a policy that prescribes actions that maximize the expected sum of discounted future rewards, i.e., she chooses at each time  $t$  the action that maximizes

$$\mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t r_t \right],$$

where  $r_t = R(s_t, a_t, s_{t+1})$  is the reward at epoch  $t$  which depends on the state transition and the action at that time. Since state transitions are stochastic, the expectation is taken over all trajectories that the process may take. Infinite horizon problems are guaranteed to converge if the discount factor  $\gamma < 1$ . For a finite time horizon, the expectation is calculated over the sum up to the end of the time horizon and a discounted expected final reward (called terminal value) may be added in the final epoch.

In a POMDP, the agent does not know the state the system is in, but it has to use observations to form a belief of what states the system could be in. This belief is called a belief state  $b \in B$  and is represented in the form of a probability distribution over the states.  $B$  is the infinite set of all possible belief states forming a  $|S| - 1$  simplex. The agent starts with an initial belief  $b_0$  (often a uniform distribution) and then updates the belief when new observations are available. In each epoch, after observing  $o$ , the agent can perform a simple Bayesian update where the updated belief for being in state  $s'$  written as  $b'(s')$  is

$$b'(s') = \eta O(o | a, s') \sum_{s \in S} T(s' | s, a) b(s),$$

and

$$\eta = \frac{1}{\sum_{s' \in S} (O(o | a, s') \sum_{s \in S} T(s' | s, a) b(s))}$$

normalizes the new belief state so all probabilities add up to one.

Regular MDPs have (under some assumptions) a deterministic optimal policy that prescribes an optimal action for each state (Puterman, 1994). Even though the actual states are not observable for POMDPs, POMDPs also have a deterministic optimal policy that prescribes an optimal action for each belief state. A policy is a mapping  $\pi : B \rightarrow A$  that prescribes for each belief state an action. The optimal policy is given by

$$\pi^* = \operatorname{argmax}_{\pi} V^{\pi}(b_0)$$

with

$$V^{\pi}(b_0) = \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t r_t \mid \pi, b_0 \right].$$

$V^{\pi}(b_0)$  is called the value function given policy  $\pi$  and the agent's initial belief  $b_0 \in B$ . The value function for any MDP or POMDP is a piecewise linear function that can be described by the highest-reward segments of a set of intersecting hyperplanes. The parameters for these hyperplanes are typically called  $\alpha$ -vectors and are a compact way to specify both, the value function and the policy of the solution of a problem.

For the infinite-horizon case, the policy converges for  $\gamma < 1$  to a policy that is independent of the

time step and the initial belief. In this case, the policy can be visualized as a directed graph called the policy graph. Each node of the graph is related to a hyperplane and represents the part of the belief space where this hyperplane produces the highest reward in the value function. Each node is labeled with the action to be taken given by the policy. The outgoing edges are labeled with observations and specify to what segment of the value function the agent will transition given the previous segment, the action and the observation. The formulation can be easily extended to the finite-horizon case. However, the finite-horizon policy depends on the initial belief and the epoch. The finite-horizon policy forms a policy tree, where each level represents an epoch.

It has to be mentioned that finding optimal policies for POMDPs is known to be a prohibitively difficult problem because the belief space grows exponentially with the number of states. This issue is called the *curse of dimensionality* in dynamic programming. Mundhenk (Mundhenk, 2000) has shown that finding the optimal policy for POMDPs is, in general, an NP<sup>PP</sup>-complete problem which means that it is at least as difficult as the hardest problems in NP. Therefore, exact algorithms can be only used for extremely small problems that are typically of very limited use in practice. More useful algorithms fall into the classes of approximate value iteration and approximate policy iteration (Cassandra, 1998a; Hauskrecht, 2000), which often find good solutions for larger problems. To use POMDPs successfully, the researcher typically needs to experiment with simplifying the problem description and choosing an acceptable level of approximation by the algorithm.

The solution of POMDPs can be used to guide the agent's actions. Automatic agents can follow very complicated policies. Humans often prefer simpler policies, even if they are not optimal but good enough to robustly improve outcomes. Simpler policies also result from problem simplification and allowing for a larger degree of approximation by the solver algorithm.

### 3 Implementation

Package **pomdp** includes a convenient and consistent way for users to define all components of a POMDP model using familiar R syntax, solve the problem using several methods and then analyze and visualize the results. An important design decision is to separate the tasks of defining a problem and analyzing the policy from the actual solver. The separation between the infrastructure in package **pomdp** and the solver code makes sure that additional solvers can be easily added in the future. Solver code is typically interfaced by writing a standard problem definition file, running an external process, and reading the results back. This way of interfacing solvers has several advantages:

- Using an external process, rather than directly interfacing the code in R ensures that memory issues for larger problems do not compromise the running R process itself.
- The separation lets the solver use any available parallelization technique without imposing limitations by R.
- Most existing solver software accepts a standard problem definition file format.
- The problem definition and the results are typically very small and fast to write and read compared to the significant amount of time used by the solver.
- Separating problem definition and result analysis from the actual solver lets the user solve larger problems on a dedicated server.

For communication with the solver, the package supports the widely used POMDP (Cassandra, 2015) file specification and can use POMDPX files (APPL Team, 2022) via package **sarsop**. This means that new algorithms that use these formats can be easily interfaced in the future, and that problems already formulated in these formats can be directly solved using the package. The authors also provide an initial set of solvers with the companion package **pomdpSolve** which provides an easy-to-install distribution of the well-known fast C implementation of a set of solvers originally developed by one of the co-authors (Cassandra, 2015). The package **pomdp** currently provides access to the following algorithms:

- Exact value iteration
  - Enumeration algorithm (Sondik, 1971; Monahan, 1982).
  - Two pass algorithm (Sondik, 1971).
  - Witness algorithm (Littman et al., 1995).
  - Incremental pruning algorithm (Zhang and Liu, 1996; Cassandra et al., 1997).
- Approximate value iteration
  - Finite grid algorithm (Cassandra, 2015), a variation of point-based value iteration to solve larger POMDPs (PBVI; see (Pineau et al., 2003)) without dynamic belief set expansion.

- SARSOP ([Kurniawati et al., 2008](#)), Successive Approximations of the Reachable Space under Optimal Policies, a point-based algorithm that approximates optimally reachable belief spaces for infinite-horizon problems (via the third-party R package `sarsop` ([Boettiger et al., 2021](#))).

While exact methods can only solve very small problems, PBVI and SARSOP can efficiently find approximate solutions for larger problems with thousands of states and hundreds of different observations. `pomdp` uses by default the finite grid algorithm.

The `pomdp` package provides efficient support by using

- sparse matrix representation based on the `Matrix` package ([Bates et al., 2022](#)) for large transition and observation matrices of low density,
- fast matrix operations,
- fast C++ implementations of loops using `Rcpp` ([Eddelbuettel, 2013](#)), and
- parallel execution using `foreach` ([Microsoft and Weston, 2022](#)).

The package implements many auxiliary functions to analyze and visualize POMDPs and their solution. For example, to sample from the belief space, simulate trajectories through a POMDP and estimate beliefs, fast C++ implementations (using `Rcpp` ([Eddelbuettel, 2013](#))) and support for parallel execution using `foreach` ([Microsoft and Weston, 2022](#)) are provided. To represent and visualize policy graphs the widely used and powerful `igraph` package ([Csardi and Nepusz, 2006](#)) with its advanced layout options is used. Interactive policy graphs can be produced based on the `visNetwork` ([Almende B.V. and Contributors and Thieurmel, 2022](#)). While the package does not directly provide functions to create `ggplot2` visualizations ([Wickham, 2016](#)) to avoid installing the large number of packages needed, the manual pages provide examples.

Solving a new POMDP problem with the `pomdp` package consists of the following steps:

1. Define a POMDP problem using the creator function `POMDP()` using R syntax,
2. solve the problem using `solve_POMDP()` which calls an external solver, and
3. analyze and visualize the results with functions like `reward()`, `plot_policy_graph()`, and `plot_value_function()`.

We will now discuss these steps in more detail and then present the complete code for a small toy example.

### 3.1 Defining a POMDP problem

The `POMDP()` creator function has as its arguments the 7-tuple  $(S, A, T, R, \Omega, O, \gamma)$ , the time horizon with terminal values, the initial belief state  $b_0$  and a name for the model. Default values are an infinite time horizon (which has no terminal values), and an initial belief state given by a uniform distribution over all states.

While specifying most parts of the POMDP is straightforward, some arguments can be specified for convenience in several different ways. Transition probabilities, observation probabilities and the reward function can be specified in several ways:

- A named list of dense or sparse matrices or the keywords "identity" and "uniform" representing the probabilities or rewards organized by action.
- As a `data.frame` representing a table with states, actions and the probabilities or reward values created with the helper functions
  - `T_(action, start.state, end.state, probability)`,
  - `O_(action, end.state, observation, probability)` and
  - `R_(action, start.state, end.state, observation, value)`.

`NA` is used to mean that a value applies to all actions, states or observations.

- An R function with the same arguments as `T_()`, `O_()` or `R_()` that returns the probability or reward.

More details can be found in the manual page for the constructor function `POMDP()`.

### 3.2 Accessing Model Data

Several parts of the POMDP description can be defined in different ways. In particular, transition probabilities, observation probabilities, rewards, and the start belief can be defined using dense

matrices, sparse matrices, data frames, functions, keywords or a mixture of all of these. The decision to specify different parts of the description using different formats is typically a result of how it is easier for the user to specify the part of the model. For example, transition and observation matrices can typically be represented efficiently as sparse matrices and the keywords `uniform` and `identity`, while rewards are typically more compactly specified as a data frame with rows describing the reward for a subset of action/state/observation combinations.

To write code that performs computation using this information requires a way to access the data in a unified way. The package provides accessor functions like:

- `start_vector()` translates the initial probability vector description into a numeric vector.
- Transition probabilities, observation probabilities, and rewards can be accessed using functions ending in `_matrix()`. Given an action, a matrix is returned. The user can request a dense or sparse matrix using the logical parameter `sparse`. To reduce the overhead associated with representing dense matrices in sparse format, sparse matrices are only returned if the density of the matrix is below 50%. The user can also specify `sparse = NULL`, which will return the data in the way it was specified by the user (e.g., a data frame). This saves the cost of conversion. Functions ending in `_val()` can be used to access individual values directly.

To allow a user-implemented algorithm direct access to the data in a uniform way, the function `normalize_POMDP()` can be used to create a new POMDP definition where transition probabilities, observation probabilities, rewards, and the start belief are consistently translated to (lists of) matrices and numeric vectors. Similar access facilities for C++ developers are also available in the package source code.

### 3.3 Solving a POMDP

POMDP problems are solved with the function `solve_POMDP()`. This function uses the low-level interface in the companion package `pomdpSolve` to solve a pomdp using the pomdp-solve software and return a solved instance of the POMDP problem. Since the low-level interfaces vary between solvers, pomdp will provide additional functions for other popular solvers. For example, for using the SARSOP solver interfaced in package `sarsop`, a function `solve_SARSOP()` is provided.

Solving POMDPs is often done by trial-and-error while simplifying the problem description to make it tractable. This means that we need to be able to interrupt the solver when it is running too long or when it runs out of memory. To accomplish this, the problem is transferred to the solver by writing a POMDP or POMDPX file, the solver software is then run in a separate process, and the results are read back. This approach results in a more robust interface since the R process is not compromised by a solver that runs out of memory or is interrupted due to too long run time. However, note that writing a large problem description file can be quite slow.

The `solve_POMDP()` and `solve_SARSOP()` functions require a POMDP model and then allow the user to specify or overwrite model parameters that are often used in experimentation like the horizon, the discount rate, and the initial belief state. Additionally, solver-specific parameters like the used algorithm for pomdp-solve can also be specified.

### 3.4 Analyzing the solution

The function `solve_POMDP()` returns a solved instance of the POMDP as a list that contains the original problem definition and an additional element containing the solution including if the solution has converged, the total expected reward given the initial belief, and the  $\alpha$ -vectors representing the value function  $V^\pi$  and the policy  $\pi$ . Keeping the problem definition and the solution together allows the user to resolve an already solved problem multiple times experimenting with different initial beliefs, horizons or discount rates, and also to perform analysis that requires both the problem definition and the solution.

An example of such an analysis is to simulate trajectories for a solved POMDP by following an  $\epsilon$ -greedy policy. An  $\epsilon$ -greedy policy follows the policy given in the solution but with a probability of  $\epsilon$  uses a random action instead, which can lead to exploring parts of the belief space that would not be reached by using only the policy. Such a simulation needs access to the policy in the solution but also to the original problem description (transition and observation probabilities). This simulation is implemented in function `simulate_POMDP()` and includes fast C++ code using `Rcpp` (Eddelbuettel, 2013) and a native R implementation supporting sparse matrix representation and sparse matrix operations. Both implementations support parallelization using `foreach` (Microsoft and Weston, 2022) to speed up the simulation.

Often it is also interesting to test the robustness of a policy on slightly modified problem descriptions or to test the performance of a manually created policy. These experiments are supported using function `add_policy()` which provides a convenient way to combine POMDP problem descriptions with compatible policies.

While the list elements of the solution can be directly accessed, several convenient access and visualization functions are provided. We provide a `plot_value_function()` that visualizes the piecewise linear value function giving the reward over the belief space simplex as a line chart for two-state problems. Function `plot_belief_space()` provides a more flexible visualization of the reward, the policy-based action, or the policy graph node over the whole belief space. A three-state problem has a belief space of the form of a 2-simplex which is a triangle and the visualization uses a ternary plot. The belief space from more than three states cannot be directly visualized, however, projections can be visualized by fixing the probabilities for all but two or three states.

The function `policy()` returns the policy as a table (data frame) consisting of one row for each value function segment with the  $\alpha$ -vector and the prescribed action as the last column. If the policy depends on the epoch, then a list of tables is returned, one for each epoch. If the policy corresponds to a realizable conditional control plan, then the policy can also be converted into an `igraph` object using the function `policy_graph()` and visualized using the function `plot_policy_graph()`. The policy graph shows the prescribed actions and how observations change the agent's belief state. This is often very useful for understanding the policy. For general finite-horizon policies, the policy graph is a policy tree where each level in the tree represents successive epochs. Such trees are often too large to visualize directly, but the `igraph` object can be used in many advanced R packages for network analysis or exported for analysis with external tools.

Further, individual belief updates, the optimal action and the expected reward given a belief can be calculated using `update_belief()`, `optimal_action()`, and `reward()`. Together with the unified accessor functions and the POMDP specifications, the user can use these functions to implement more sophisticated R-based analysis. The source package also contains C++ implementations of these and the accessor functions. These can be used by an advanced R developer to write fast analysis code or implement custom solvers.

### 3.5 Time-dependent POMDPs

For some real-world problems, the transition probabilities, observation probabilities, or rewards may change depending on the epoch. For example, in a medical application, the transition probability modeling the chance of getting an infection may increase with the age of the patient. While the general definition of POMDPs can be easily extended to allow time-dependent transition probabilities, observation probabilities and reward functions to model changes in the modeled system, most existing solvers use fixed matrices.

The package `pomdp` adds a simple mechanism to support time dependence. Time dependence of transition probabilities, observation probabilities and the reward structure can be modeled by considering a set of episodes representing epochs with the same settings and then solving these episodes in reverse order with the accumulated discounted reward of each episode used as the final reward for the preceding episode. Details on how to specify episodes in time-dependent POMDPs can be found in the `pomdp` manual pages.

## 4 Toy Example: The Tiger problem

We will demonstrate how to use the package with a popular toy problem called the Tiger Problem ([Cassandra et al., 1994](#)). This example is often used to introduce students to POMDPs. The problem is defined as:

An agent is facing two closed doors, and a tiger is put with equal probability behind one of the two doors represented by the environment states `tiger-left` and `tiger-right` while treasure is put behind the other door. The available actions are `listen` for tiger noises or opening a door (actions `open-left` and `open-right`). Listening is neither free (the action has a reward of -1) nor is it entirely accurate. There is a 15% probability that the agent hears the tiger behind the left door while it is behind the right door and vice versa. If the agent opens the door with the tiger, it will get hurt (a reward of -100), but if it opens the door with the treasure, it will receive a positive reward of 10. After a door is opened, the problem resets (i.e., the tiger is again randomly assigned to a door), and the agent gets another try. This makes it an infinite horizon problem and we use a discount factor of .75 to guarantee convergence.

#### 4.1 Specifying the Tiger problem

The problem can be specified using the function `POMDP()`.

```
library("pomdp")

Tiger <- POMDP(
  name = "Tiger Problem",
  discount = 0.75,
  states = c("tiger-left", "tiger-right"),
  actions = c("listen", "open-left", "open-right"),
  observations = c("tiger-left", "tiger-right"),
  start = "uniform",

  transition_prob = list(
    "listen" = "identity",
    "open-left" = "uniform",
    "open-right" = "uniform"),

  observation_prob = list(
    "listen" = matrix(c(0.85, 0.15, 0.15, 0.85), nrow = 2, byrow = TRUE),
    "open-left" = "uniform",
    "open-right" = "uniform"),

  reward = rbind(
    R_("listen", NA, NA, NA, -1),
    R_("open-left", "tiger-left", NA, NA, -100),
    R_("open-left", "tiger-right", NA, NA, 10),
    R_("open-right", "tiger-left", NA, NA, 10),
    R_("open-right", "tiger-right", NA, NA, -100)
  )
)
```

Note that we use for each component the most convenient specification method. For observations and transitions, we use a list of distribution keywords and a matrix, while for the rewards, a data frame created with the `R_()` function is used. The `R_()` function accepts the arguments `action`, `start.state`, `end.state`, `observation`, and the reward value. A missing value of `NA` indicates that the reward is valid for any state or observation.

The transition model can be visualized as a graph.

```
g <- transition_graph(Tiger)

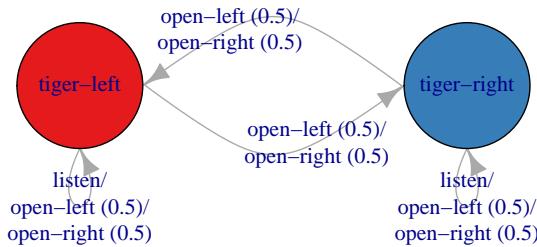
library(igraph)
plot(g,
  layout = rbind(c(-1, 0), c(1, 0)), rescale = FALSE,
  edge.curved = curve_multiple_directed(g, .8),
  edge.loop.angle = pi / 2,
  vertex.size = 65
)
```

The vertices in Figure 1 represent the states and the edges show transitions labeled with actions and the associated transition probabilities in parentheses. Multiple parallel transitions are collapsed into a single arrow with several labels to simplify the visualization. The graph shows that the action `listen` stays with a probability of 1 in the same state (i.e., listening does not move the tiger). The actions `open-left` and `open-right` lead to a reset of the problem which assigns the tiger randomly to a state. This is represented by the transitions with a probability of .5.

For more complicated transition models, individual graphs for each action or interactive graphs using `visNetwork` can also be plotted.

#### 4.2 Solving the Tiger problem for an infinite time horizon

To solve the problem, we use the default method (`pomdp-solve`'s finite grid method interfaced in package `pomdpSolve`) which performs a form of point-based value iteration that can find approximate solutions for larger problems.



**Figure 1:** Transition model of the Tiger problem.

```

sol <- solve_POMDP(Tiger)
sol

#> POMDP, list - Tiger Problem
#> Discount factor: 0.75
#> Horizon: Inf epochs
#> Size: 2 states / 3 actions / 2 obs.
#> Start: uniform
#> Solved:
#>   Method: 'grid'
#>   Solution converged: TRUE
#>   # of alpha vectors: 5
#>   Total expected reward: 1.933439
#>
#> List components: 'name', 'discount', 'horizon', 'states', 'actions',
#>   'observations', 'transition_prob', 'observation_prob', 'reward',
#>   'start', 'info', 'solution'
  
```

The solver returns an object of class POMDP, which contains the solution as an additional list component. The print function displays important information like the used discount factor, the horizon, if the solution has converged and the total expected reward. In this case, the total expected discounted reward for following the policy starting from the initial belief is 1.933. Note that the optimal policy for infinite-horizon does not depend on the initial belief. The reward for other initial beliefs can be calculated using the reward() function. For example, the expected reward for a correct belief that the tiger starts to the left with a probability of 90% is:

```

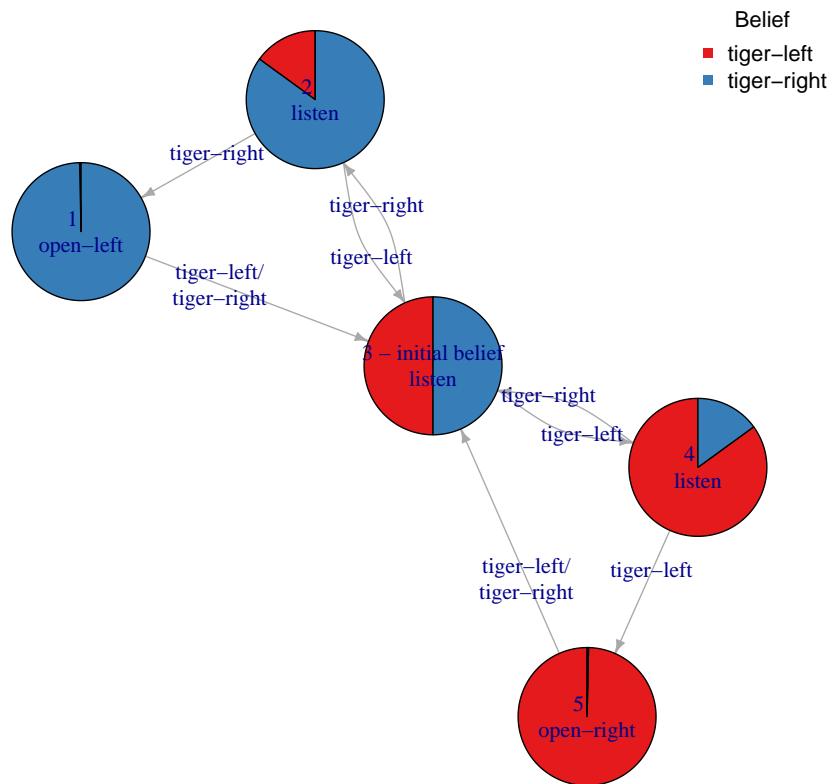
reward(sol, belief = c(0.9, 0.1))

#> [1] 4.779814
  
```

### 4.3 Inspecting the Policy

The policy of a solved POMDP is a set of  $\alpha$ -vectors representing a segment of the value function and the associated best action.

```
policy(sol)
```



**Figure 2:** The policy graph for the converged infinite-horizon solution of the Tiger problem.

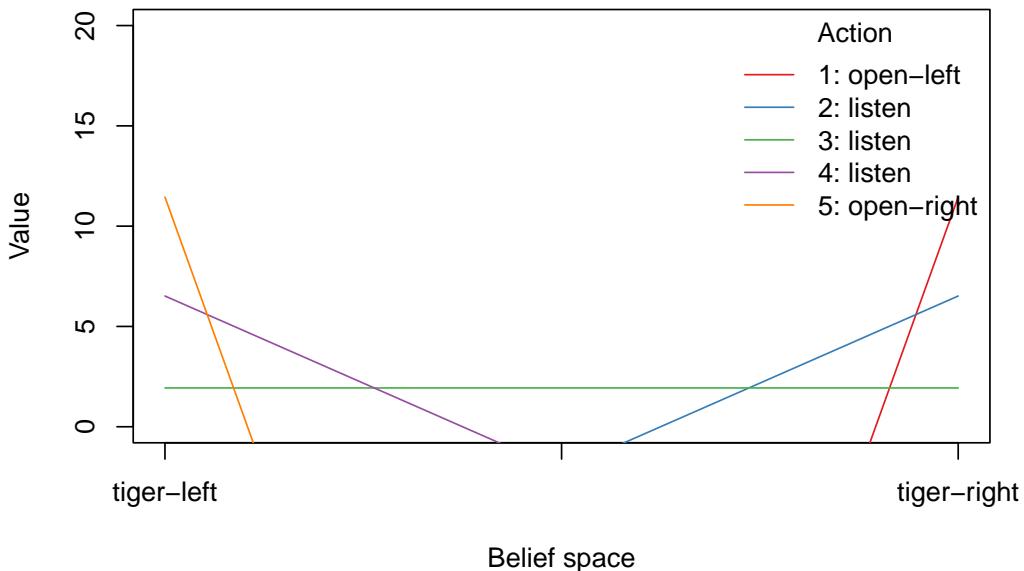
```
#>   tiger-left tiger-right      action
#> 1 -98.549921  11.450079  open-left
#> 2 -10.854299   6.516937   listen
#> 3  1.933439   1.933439   listen
#> 4   6.516937  -10.854299  listen
#> 5  11.450079  -98.549921 open-right
```

The returned policy is a list where each element represents the  $\alpha$ -vectors for an epoch. The policy above has only one list element since the solution converged to a solution that is independent of the epoch.

Smaller policies that correspond to a conditional plan can also be represented as a graph using a custom plot function.

```
plot_policy_graph(sol)
```

The function uses the `igraph` package (Csardi and Nepusz, 2006) to produce the layout. Figure 2 shows the graph for the optimal policy returned by the solver for the Tiger problem. Each node in the policy graph represents an  $\alpha$ -vector and is labeled by the action prescribed by the policy. Each segment covers a part of the belief space which represents how much the agent knows about the location of the tiger based on all previous observations. We use a pie chart inside each node to show a representative belief point that belongs to the segment. This makes it easier to compare the beliefs in different nodes with each other. The representative belief points are found with the function `estimate_belief_for_nodes()` which uses the solver output and searches along policy trajectories.



**Figure 3:** The value function for the solution of the converged Tiger problem.

It is easy to interpret smaller policy graphs. Figure 2 shows that without prior information, the agent starts at the node marked with initial belief. In this case, the agent believes there is a 50/50 chance that the tiger is behind either door. The optimal action is displayed inside the state and, in this case, is to listen. The arcs are labeled with observations. Let us assume that the observation is tiger-left. The agent follows the appropriate arc and ends in a node representing the new range of belief states with a higher probability of the tiger being to the left. However, the optimal action is still to listen. If the agent again hears the tiger on the left then it ends up in a node that has a belief of close to 100% that the tiger is to the left and open-right is the optimal action. The arcs back from the nodes with the open actions to the initial state reset the problem and let the agent start over.

Typically, small and compact policy graphs are preferable in practice because they make the policy easier to understand for the decision maker and also easier to follow. For large, more complicated policy graphs, representation as a graph is difficult leading to issues with node layout and too many crossing vertices. The package can also plot the graph as an interactive HTML widget with movable vertices (see the manual page for `plot_policy_graph()`) to let the user arrange the graph manually. Larger policy graphs can also be exported in common formats like graphML to be displayed and analyzed in large-scale network analysis tools like Gephi (Jacomy et al., 2014).

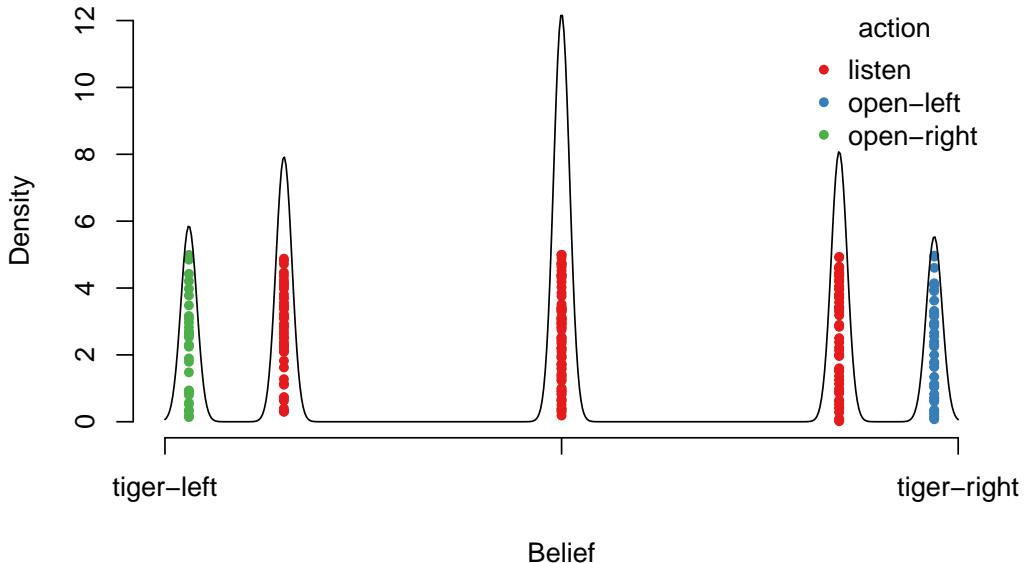
The Tiger problem environment has only two states (tiger-left and tiger-right) with a belief space forming a 1-simplex which is a line going from a probability of 1 that the tiger is left to a probability of 1 that the tiger is right. Therefore, we can visualize the piecewise linear convex value function as a simple line chart with the belief on the x-axis.

```
plot_value_function(sol, ylim = c(0,20))
```

Figure 3 shows the value function. The x-axis represents the belief, the lines represent the nodes in the policy graph (the numbers in the legend match the numbers in the graph in Figure 3), and the piecewise linear value function consists of the line segments with the highest reward. The optimal action for each segment is shown in the legend. This visualization function is mostly provided to study small textbook examples with two states. A more versatile function is `plot_belief_space()` which can produce ternary plots for problems with three or more states by projecting the belief space on three states.

Auxiliary functions provided in the package let the user perform many analyses. For example, we simulate trajectories through the POMDP belief space by following the policy and estimating the distribution of the agent's belief.

```
sim <- simulate_POMDP(sol, n = 50, horizon = 5,
```



**Figure 4:** Belief states reached in 50 simulated trajectories of horizon 5.

```
belief = c(.5, .5), return_beliefs = TRUE)
plot_belief_space(sol, sample = sim$belief_states, ylim = c(0, 12),
                  jitter = 5)
lines(density(sim$belief_states[, 1], bw = .01, from = 0 , to = 1))
axis(2); title(ylab = "Density")
```

Figure 4 shows the five beliefs that are reached in the trajectories as dots and uses jitter and a density estimate to show how much time the agent has spent in the simulation in different parts of the belief space. The color of the dots indicates the actions chosen by the policy.

#### 4.4 Solving the Tiger problem for a finite time horizon

To demonstrate how to solve a POMDP problem with a finite time horizon, we set the horizon to 4 epochs, which means that the agent starts with its initial belief and can perform only four actions. The grid-based method used before finds the optimal policy, but for finite time horizon problems with negative rewards, the value function and the calculated expected reward is only valid when the solution converges. To avoid this issue, we use here the incremental pruning algorithm (Zhang and Liu, 1996; Cassandra et al., 1997).

```
sol <- solve_POMDP(model = Tiger, horizon = 4, method = "incprune")
sol

#> POMDP, list - Tiger Problem
#>   Discount factor: 0.75
#>   Horizon: 4 epochs
#>   Size: 2 states / 3 actions / 2 obs.
#>   Start: uniform
#>   Solved:
#>     Method: 'incprune'
#>     Solution converged: FALSE
#>     # of alpha vectors: 26
#>     Total expected reward: 0.483125
#>
#>   List components: 'name', 'discount', 'horizon', 'states', 'actions',
#>     'observations', 'transition_prob', 'observation_prob', 'reward',
```

```
#>      'start', 'info', 'solution'

policy(sol)

#> [[1]]
#>   tiger-left tiger-right    action
#> 1 -99.321250  10.678750  open-left
#> 2 -11.820719   4.640094  listen
#> 3 -2.734955   2.600990  listen
#> 4 -1.137420   1.595135  listen
#> 5  0.483125   0.483125  listen
#> 6  1.595135  -1.137420  listen
#> 7  2.600990  -2.734955  listen
#> 8  4.640094 -11.820719  listen
#> 9 10.678750 -99.321250 open-right
#>
#> [[2]]
#>   tiger-left tiger-right    action
#> 1 -101.312500  8.687500  open-left
#> 2 -20.550156   5.488906  listen
#> 3 -13.450000   4.700000  listen
#> 4 -3.565469   2.157969  listen
#> 5  0.905000   0.905000  listen
#> 6  2.157969  -3.565469  listen
#> 7  4.700000 -13.450000  listen
#> 8  5.488906 -20.550156  listen
#> 9  8.687500 -101.312500 open-right
#>
#> [[3]]
#>   tiger-left tiger-right    action
#> 1 -100.7500    9.2500  open-left
#> 2 -12.8875    5.2625  listen
#> 3 -1.7500     -1.7500  listen
#> 4  5.2625    -12.8875  listen
#> 5  9.2500    -100.7500 open-right
#>
#> [[4]]
#>   tiger-left tiger-right    action
#> 1      -100        10  open-left
#> 2       -1         -1  listen
#> 3       10        -100 open-right
```

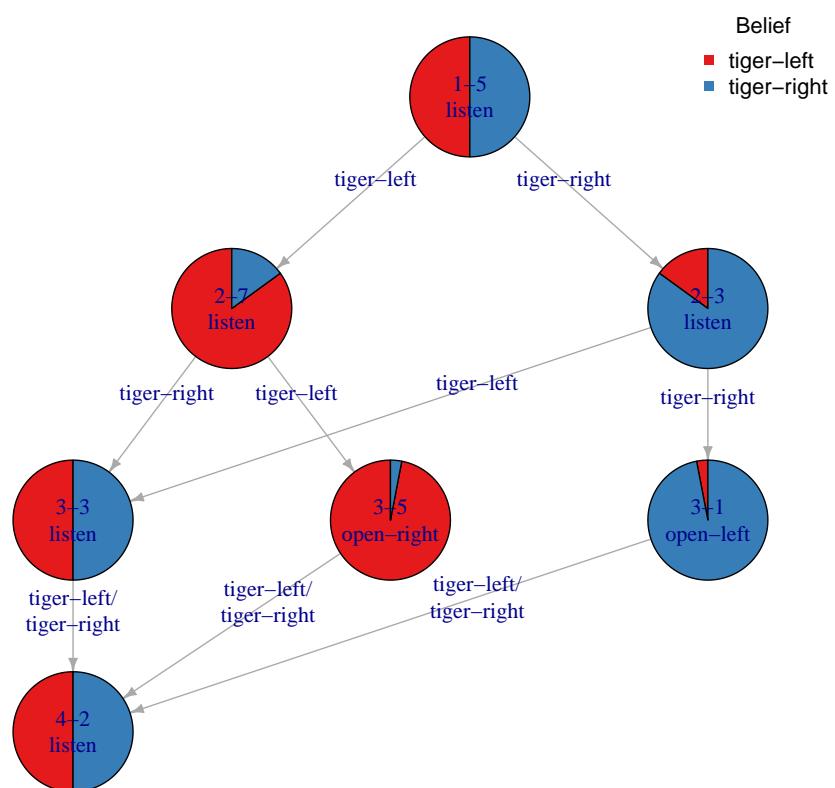
The policy has four elements, one for each epoch. Is easier to understand the policy by visualizing it as a graph.

```
plot_policy_graph(sol)
```

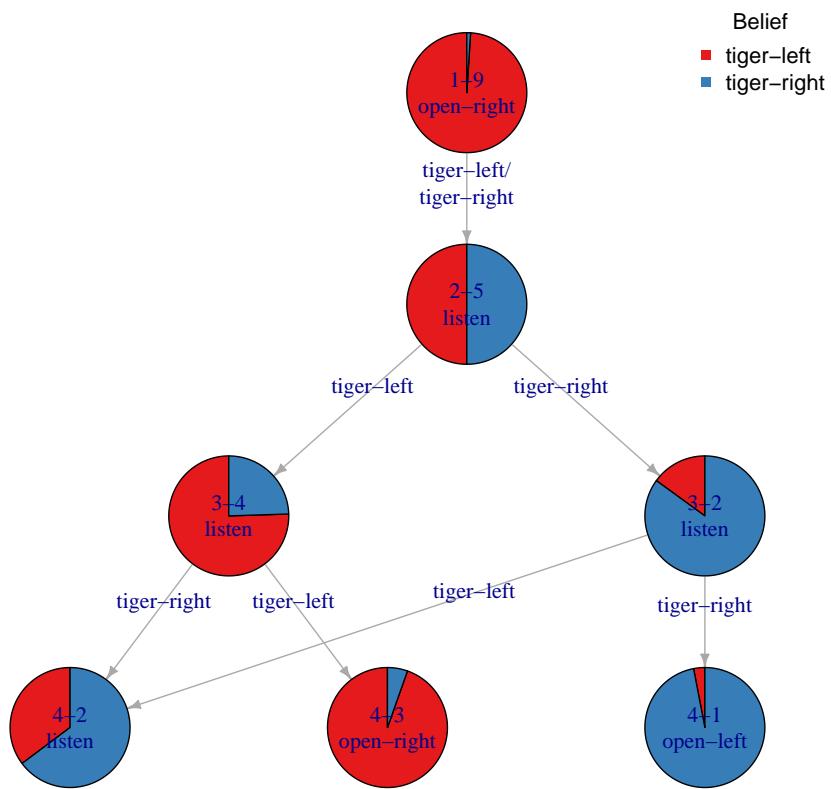
The resulting policy graph is shown in Figure 5 as a tree with four levels, one for each time epoch. The plot function automatically uses a tree layout and adds the epoch as the first number to the node labels. By default, it also simplifies the representation by hiding belief states which cannot be reached from the start belief and, therefore, there are more entries in the policy above than there are nodes in the graph. The root node of the tree represents the initial belief used in the model. The model starts with a uniform initial belief represented by the evenly split pie chart. The policy shows that the optimal strategy is to listen and open a door only if we hear the tiger behind the same door twice. Interestingly, it is optimal never to open a door in the last epoch. The reason is that we cannot reach a sufficiently high belief of the tiger being behind a single door. The expected reward of this policy starting at a uniform initial belief is 0.483.

Policy trees for finite-horizon problems are dependent on the agent's initial belief. To show this, we produce a new policy tree for an initial belief of 99% that the tiger is to the left by overwriting the initial belief in the model definition.

```
sol <- solve_POMDP(model = Tiger, horizon = 4,
                     initial_belief = c(.99, .01), method = "inprune")
reward(sol, belief = c(.99, .01))
```



**Figure 5:** Policy tree for the Tiger problem solved with a horizon of 4 and a uniform initial belief.



**Figure 6:** Policy tree for the Tiger problem solved with a horizon of 4 and an initial belief of 99 percent that the tiger is to the left.

```
#> [1] 9.57875
plot_policy_graph(sol, belief = c(.99, .01))
```

The resulting policy graph with an initial belief indicating that we are very sure that the tiger is to the left is shown in Figure 6. The graph indicates that it is optimal to open the right door right away and then wait if we hear the tiger twice in the same location before we open the other door. Under the strong belief, the agent also expects a much higher reward of 9.579 for the optimal policy.

## 5 Summary

Partially observable Markov decision processes are an important modeling technique useful for many applications. Easily accessible software to solve POMDP problems is crucial to support applied research and instruction in fields including artificial intelligence and operations research. Most existing libraries need advanced technical expertise to install and offer minimal support to analyze the results. The `pomdp` package fills this gap by providing an easily accessible platform to perform experiments and analyze POMDP problems and the resulting policies.

This paper used a minimalist toy example to show the functionality of the package in a concise way. Studying and visualizing complicated policies with hundreds or thousands of belief states is an important topic that has received less attention than improving solver algorithms. R provides a wide range of tools to compare, analyze, and cluster belief states. We plan to investigate the use of these

techniques to support explainability of more complicated policies and will implement corresponding functions in future releases of the package `pomdp`.

## 6 Acknowledgments

Farzad Kamalzadeh participated in the development of an early version of the `pomdp` package and used it for several applications. He was supported by a Graduate Fellowship and a Niemi Center Fellowship, both at SMU. His and Michael Hahsler's work was also supported in part by the National Institute of Standards and Technology (NIST) under grant number 60NANB17D180.

The authors would also like to thank Carl Boettiger for maintaining the package `sarsop` and the anonymous reviewers for their valuable insights.

## References

- Almende B.V. and Contributors and B. Thieurmel. *visNetwork: Network Visualization using 'vis.js' Library*, 2022. URL <https://CRAN.R-project.org/package=visNetwork>. R package version 2.1.2. [p120]
- APPL Team. APPL: Approximate POMDP planning toolkit, 2022. URL <https://bigbird.comp.nus.edu.sg/pmwiki/farm/app1/>. [p117, 119]
- D. Bates, M. Maechler, and M. Jagan. *Matrix: Sparse and Dense Matrix Classes and Methods*, 2022. URL <https://CRAN.R-project.org/package=Matrix>. R package version 1.5-3. [p120]
- C. Boettiger, J. Ooms, and M. Memarzadeh. *sarsop: Approximate POMDP Planning Software*, 2021. URL <https://CRAN.R-project.org/package=sarsop>. R package version 0.6.9. [p117, 120]
- A. R. Cassandra. *Exact and Approximate Algorithms for Partially Observable Markov Decision Processes*. PhD thesis, Providence, RI, USA, 1998a. AAI9830418. [p119]
- A. R. Cassandra. A survey of POMDP applications. Technical Report MCC-INSL-111-98, Microelectronics and Computer Technology Corporation (MCC), 1998b. Presented at the AAAI Fall Symposium. [p117]
- A. R. Cassandra. The POMDP page, 2015. URL <https://www.pomdp.org>. [p117, 119]
- A. R. Cassandra, L. P. Kaelbling, and M. L. Littman. Acting optimally in partially observable stochastic domains. In *Proceedings of the Twelfth National Conference on Artificial Intelligence*, Seattle, WA, 1994. AAAI Classic Paper Award, 2013. [p122]
- A. R. Cassandra, M. L. Littman, and N. L. Zhang. Incremental pruning: A simple, fast, exact method for partially observable Markov decision processes. In *UAI'97: Proceedings of the Thirteenth conference on Uncertainty in artificial intelligence*, pages 54–61, August 1997. [p119, 127]
- G. Csardi and T. Nepusz. The igraph software package for complex network research. *InterJournal, Complex Systems*:1695, 2006. URL <https://igraph.org>. [p120, 125]
- D. Eddelbuettel. *Seamless R and C++ Integration with Rcpp*. Springer, New York, 2013. doi: 10.1007/978-1-4614-6868-4. ISBN 978-1-4614-6867-7. [p120, 121]
- M. Hahsler. *pomdp: Infrastructure for Partially Observable Markov Decision Processes (POMDP)*, 2023. URL <https://github.com/mhahsler/pomdp>. R package version 1.1.3. [p117]
- M. Hahsler and A. R. Cassandra. *pomdpSolve: Interface to 'pomdp-solve' for Partially Observable Markov Decision Processes*, 2022. URL <https://github.com/mhahsler/pomdpSolve>. R package version 1.0.2. [p117]
- M. Hauskrecht. Value-function approximations for POMDPs. *Journal Of Artificial Intelligence Research*, 13:33–94, 2000. doi: <https://doi.org/10.1613/jair.678>. [p119]
- M. Jacomy, T. Venturini, S. Heymann, and M. Bastian. Forceatlas2, a continuous graph layout algorithm for handy network visualization designed for the Gephi software. *PLOS ONE*, 9(6):1–12, 06 2014. doi: 10.1371/journal.pone.0098679. [p126]
- JuliaPOMDP Team. JuliaPOMDP: POMDP packages for Julia, 2022. URL <https://github.com/JuliaPOMDP>. [p117]

- L. P. Kaelbling, M. L. Littman, and A. R. Cassandra. Planning and acting in partially observable stochastic domains. *Artificial Intelligence*, 101(1):99–134, 1998. ISSN 0004-3702. doi: 10.1016/S0004-3702(98)00023-X. [p116, 118]
- F. Kamalzadeh, V. Ahuja, M. Hahsler, and M. E. Bowen. An analytics-driven approach for optimal individualized diabetes screening. *Production and Operations Management*, 30(9):3161–3191, September 2021. ISSN 1937-5956. doi: 10.1111/poms.13422. URL <https://onlinelibrary.wiley.com/doi/abs/10.1111/poms.13422>. [p116, 117]
- H. Kurniawati, D. Hsu, and W. S. Lee. SARSOP: Efficient point-based pomdp planning by approximating optimally reachable belief spaces. In *In Proc. Robotics: Science and Systems*, 2008. [p120]
- M. L. Littman, A. R. Cassandra, and L. P. Kaelbling. Learning policies for partially observable environments: Scaling up. In *Proceedings of the Twelfth International Conference on International Conference on Machine Learning*, ICML'95, page 362–370, San Francisco, CA, USA, 1995. Morgan Kaufmann Publishers Inc. ISBN 1558603778. [p119]
- Microsoft and S. Weston. *foreach: Provides Foreach Looping Construct*, 2022. URL <https://CRAN.R-project.org/package=foreach>. R package version 1.5.2. [p120, 121]
- B. Migge and O. Stollmann. pyPOMDP: POMDP implementation in Python, 2013. URL <https://bitbucket.org/bami/pypomdp/src/master/>. [p117]
- G. E. Monahan. A survey of partially observable Markov decision processes: Theory, models, and algorithms. *Management Science*, 28(1):1–16, 1982. [p119]
- M. Mundhenk. The complexity of optimal small policies. *Math. Oper. Res.*, 25(1):118–129, feb 2000. ISSN 0364-765X. [p119]
- J. Pineau, G. Gordon, and S. Thrun. Point-based value iteration: An anytime algorithm for POMDPs. In *Proceedings of the 18th International Joint Conference on Artificial Intelligence*, IJCAI'03, page 1025–1030, San Francisco, CA, USA, 2003. Morgan Kaufmann Publishers Inc. [p119]
- N. Proellos and S. Feuerriegel. *ReinforcementLearning: Model-Free Reinforcement Learning*, 2020. URL <https://CRAN.R-project.org/package=ReinforcementLearning>. R package version 1.0.5. [p117]
- M. L. Puterman. Markov decision processes: Discrete stochastic dynamic programming. In *Wiley Series in Probability and Statistics*, 1994. [p116, 118]
- R. Smallwood and E. Sondik. The optimal control of partially observable Markov decision processes over a finite horizon. *Operations Research*, 21(5):1071–88, 1973. [p116]
- T. Smith. ZMDP: Software for POMDP and MDP planning, 2009. URL <https://github.com/trey0/zmdp>. [p117]
- E. J. Sondik. *The Optimal Control of Partially Observable Markov Decision Processes*. PhD thesis, Stanford, California, 1971. [p119]
- R. S. Sutton and A. G. Barto. *Reinforcement Learning: An Introduction*. The MIT Press, second edition, 2018. [p117]
- H. Wickham. *ggplot2: Elegant Graphics for Data Analysis*. Springer-Verlag New York, 2016. ISBN 978-3-319-24277-4. URL <https://ggplot2.tidyverse.org>. [p120]
- N. L. Zhang and W. Liu. Planning in stochastic domains: Problem characteristics and approximation. Technical Report HKUST-CS96-31, Hong Kong University, 1996. [p119, 127]
- K. Åström. Optimal control of Markov processes with incomplete state information. *Journal of Mathematical Analysis and Applications*, 10(1):174 – 205, 1965. ISSN 0022-247X. [p116]

Michael Hahsler  
Southern Methodist University  
Department of Computer Science  
Dallas, TX, USA  
<https://michael.hahsler.net>  
ORCID: 0000-0003-2716-1405  
[mhahsler@lyle.smu.edu](mailto:mhahsler@lyle.smu.edu)

*Anthony R. Cassandra  
POMDP, Inc  
Austin, TX, USA  
<https://tonycassandra.com>  
tony.cassandra@gmail.com*

# pencal: an R Package for the Dynamic Prediction of Survival with Many Longitudinal Predictors

by Mirko Signorelli

**Abstract** In survival analysis, longitudinal information on the health status of a patient can be used to dynamically update the predicted probability that a patient will experience an event of interest. Traditional approaches to dynamic prediction such as joint models become computationally unfeasible with more than a handful of longitudinal covariates, warranting the development of methods that can handle a larger number of longitudinal covariates. We introduce the R package `pencal`, which implements a Penalized Regression Calibration (PRC) approach that makes it possible to handle many longitudinal covariates as predictors of survival. `pencal` uses mixed-effects models to summarize the trajectories of the longitudinal covariates up to a prespecified landmark time, and a penalized Cox model to predict survival based on both baseline covariates and summary measures of the longitudinal covariates. This article illustrates the structure of the R package, provides a step by step example showing how to estimate PRC, compute dynamic predictions of survival and validate performance, and shows how parallelization can be used to significantly reduce computing time.

## 1 Introduction

Risk prediction models (Steyerberg, 2009) allow to estimate the probability that an event of interest will occur in the future. Such models are commonly employed in the biomedical field to estimate the probability that an individual will experience an adverse event, and their output can be used to inform patients, monitor their disease progression, and guide treatment decisions.

Traditionally, risk prediction models only used covariate values available at the beginning of the observation period to predict survival. Thus, predictions based on such models could not exploit information gathered at later time points. Because information about the evolution of time-dependent covariates may influence the occurrence of the survival outcome, it is desirable to be able to dynamically update predictions of survival as more longitudinal information becomes available.

Dynamic prediction models employ both baseline and follow-up information to predict survival, and they can be used to update predictions each time new follow-up data are gathered. Three commonly-used statistical methods for dynamic prediction are the time-dependent Cox model, joint modelling of longitudinal and survival data, and landmarking approaches.

The time-dependent Cox model (Therneau and Grambsch, 2000) is an extension of the Cox proportional hazards model that allows for the inclusion of time-dependent covariates. It assumes the value of a time-dependent covariate to be constant between two observation times, and it is only suitable for exogenous time-dependent covariates. The model can be estimated using the R package `survival` (Therneau and Grambsch, 2000).

Joint models for longitudinal and survival outcomes (Henderson et al., 2000) are shared random effects models that combine a submodel for the longitudinal covariates (typically linear mixed models) and one for the survival outcome (usually a Cox model or a parametric survival model). Thanks to the shared random effects formulation, such models are capable to account for a possible interdependence between the longitudinal covariates and the survival outcome. However, the estimation of the shared random effects model is a computationally intensive task that has so far restricted the application of joint models to problems with one or few longitudinal covariates. Over the years, several alternative approaches to the estimation of joint models have been proposed, among which are the R packages `JM` (Rizopoulos, 2010), `JMbayes` (Rizopoulos, 2014), `joineR` (Philipson et al., 2018) and `joineRML` (Hickey et al., 2018).

Lastly, landmarking (Van Houwelingen, 2007) is an approach that dynamically adjusts predictions by refitting the prediction model using all subjects that are still at risk at a given landmark time. Landmarking typically involves two modelling steps. In the first step, repeated measurements of the time-dependent covariates up to the landmark time are summarized using either a summary measure or a suitable statistical model. In the second step, the summaries thus computed are used as predictors of survival alongside with the time-independent covariates. The simplest form of landmarking is the last observation carried forward (LOCF) method, which uses the last available measurement of each longitudinal covariate taken up to the landmark time as summary. The main advantage of this approach is that it is easy to implement, it is computationally straightforward and,

thus, it does not require the development of dedicated software. Limitations of LOCF landmarking include the fact that it discards all previous repeated measurements, failing to make efficient use of the available longitudinal information, and it does not perform any measurement error correction on the longitudinal covariates (this may be particularly desirable for biomarkers and diagnostic tests that are typically subject to measurement error). To overcome these limitations, mixed-effects models can be used to model the trajectories of the longitudinal covariates (Signorelli et al., 2021; Putter and van Houwelingen, 2022; Devaux et al., 2023). While Putter and van Houwelingen (2022) focused on situations with a single longitudinal marker, Signorelli et al. (2021) and Devaux et al. (2023) proposed two methods, respectively called Penalized Regression Calibration (PRC) and DynForest, that can deal with a large number of longitudinal covariates. Notably, the estimation of PRC and DynForest is much more complex than that of LOCF landmarking, warranting dedicated software that can facilitate the implementation of such methods. PRC is implemented in the R package `pencal` (Signorelli, 2023), whereas DynForest in the R package `DynForest` (Devaux et al., 2023).

In this article we introduce the R package `pencal`, which implements the PRC approach. PRC uses mixed-effects models to describe and summarize the longitudinal biomarker trajectories; the summaries thus obtained are used as predictors of survival in a Cox model alongside with any relevant time-independent covariate. To account for the possible availability of a large (potentially high-dimensional) number of time-independent and longitudinal covariates and to reduce the risk of overfitting the training data, PRC uses penalized maximum likelihood to estimate the aforementioned Cox model.

The remainder of the article is organized as follow. In the next section we describe the dynamic prediction problem, the PRC statistical methodology (Signorelli et al., 2021) behind `pencal`, and the problem of evaluating the model's predictive performance. Next we provide a general overview of the R package, discussing the implementation details of its functions for model estimation, prediction and performance validation. Furthermore, we provide a step by step example that shows how to use `pencal` to implement dynamic prediction on a real-world dataset that comprises several longitudinal covariates. We present the results of 4 simulations that assess the relationship between computing time and sample size, number of covariates and number of bootstrap samples used to validate model performance, showing how parallelization may reduce computing time significantly. Lastly, we provide some final remarks, and discuss limitations and possible extensions of the current approach.

## 2 Statistical methods

### 2.1 Input data and notation

We consider a setting where  $n$  subjects are followed from time  $t = 0$  until an event of interest occurs. For each subject  $i \in \{1, \dots, n\}$  we observe the pair  $(t_i, \delta_i)$ , where  $\delta_i$  is a dummy variable that indicates whether the event is observed at time  $t = t_i$  ( $\delta_i = 1$ ), or the observation of the event is right-censored at  $t = t_i$  ( $\delta_i = 0$ ). Thus,  $t_i$  corresponds to the survival time if  $\delta_i = 1$ , and to the censoring time if  $\delta_i = 0$ .

In addition to  $(t_i, \delta_i)$ , we assume that both baseline and follow-up information is collected from the same subjects, and that the number of variables gathered may be large. We consider a flexible unbalanced study design where the number and timing of the follow-up times can differ across subjects. We denote by  $m_i \geq 1$  the number of repeated measurements available for subject  $i$ , and by  $t_{i1}, \dots, t_{im_i}$  ( $t_{ij} \geq 0 \forall j$ ) the corresponding follow-up times. For each subject  $i$  we observe:

1. a vector of  $k$  baseline (time-fixed) predictors  $x_i = (x_{1i}, \dots, x_{ki})$ ;
2.  $m_i$  vectors of  $p$  longitudinal (time-varying) predictors  $y_{ij} = (y_{1ij}, \dots, y_{pj})$ ,  $j = 1, \dots, m_i$  measured at times  $t_{i1}, \dots, t_{im_i}$ . Note that not all longitudinal predictors ought to be measured at every follow-up time  $t_{ij}$ , and the number of available measurements is thus allowed to differ across longitudinal predictors.

### 2.2 Dynamic prediction of survival

Let  $S_i(t) = P(T_i > t)$  denote the survival function, i.e., the probability that subject  $i$  has not experienced the event up to time  $t$ , and let  $S_i(t_B|t_A) = P(T_i > t_B | T_i > t_A)$ ,  $t_B \geq t_A$  denote the conditional probability that subject  $i$  survives up until  $t_B$ , given that they survived up until  $t_A$ . Our goal is to predict the probability of survival of subject  $i$  given all the available information up until a given landmark time  $t_L > 0$ , namely:

$$S_i(t|t_L, x_i, \mathcal{Y}_i(t_L)) = P(T_i > t | T_i > t_L, x_i, \mathcal{Y}_i(t_L)), \quad (1)$$

where  $\mathcal{Y}_i(t_L) = \{y_{i1}, \dots, y_{ir} : t_{ir} \leq t_L\}$  denotes all repeated measurements available up to the landmark time for subject  $i$ .

In practice, often one may be interested in computing the predictions of survival in (1) over a range of  $q$  landmark times  $t_{L1}, t_{L2}, \dots, t_{Lq}$ . By doing so, information from more and more repeated measurements can be incorporated in the prediction model as time passes, and predictions can be updated based on the latest available information. This approach is referred to as *dynamic prediction*, as it involves dynamic updates of model estimates and predictions over time. In this article we illustrate how to use `pencal` to compute predictions of the conditional survival probabilities in (1) for a single landmark time  $t_L$ . Note that implementing a dynamic prediction approach with `pencal` is straightforward, as for each landmark time  $t_{Ls}$  one simply needs to refit PRC over a dataset that comprises all repeated measurements available up until  $t = t_{Ls}$  for the subjects that survived up until the landmark time  $t_{Ls}$  (i.e., including subject  $i$  if and only if  $t_i \geq t_{Ls}$ ).

### 2.3 Penalized regression calibration

PRC (Signorelli et al., 2021) is a statistical method that makes it possible to estimate the conditional survival probabilities in (1) using  $x_i$  and  $\mathcal{Y}_i(t_L)$  as inputs. The estimation of PRC requires a multi-step procedure that comprises the 3 following steps:

1. model the evolution over time of the longitudinal predictors in  $\mathcal{Y}_i(t_L)$  using linear mixed models (LMM, McCulloch and Searle (2004)) or multivariate latent process mixed models (MLPMM, Proust-Lima et al. (2013));
2. use the model(s) fitted at step 1 to compute summaries of the trajectories described by the longitudinal predictors (i.e., the predicted random effects);
3. estimate a penalized Cox model for the survival outcome  $(t_i, \delta_i)$  using as covariates both the baseline predictors  $x_i$  and the predicted random effects computed in step (2).

For simplicity, in this section we describe the version of PRC where in step 1 each longitudinal covariate is modelled using a separate LMM. This version of the model is referred to as PRC LMM in Signorelli et al. (2021). Note that alongside with the PRC LMM approach, Signorelli et al. (2021) also proposed a second approach called PRC MLPMM, where groups of longitudinal covariates are modelled jointly using the MLPMM. This alternative approach can be of interest when multiple longitudinal items are employed to measure the same underlying quantity (for example, multiple antibodies that target the same protein). For the formulation of the PRC MLPMM approach, we refer readers to Sections 2.1-2.3 of Signorelli et al. (2021) as the notation for steps 1 and 2 using the MLPMM is significantly more involved.

Denote by  $\mathcal{I}(t_L) = \{i : t_i > t_L\}$  the set of subjects that survived up until the landmark time  $t_L$ . Let  $y_{si} = (y_{s1}, \dots, y_{sr})$ , where  $t_{ir} \leq t_L$  denotes the last follow-up time before  $t_L$  for subject  $i$ , be the vector that comprises all the measurements of the  $s$ -th longitudinal variable  $Y_s$  available up to the landmark time. In the first step of PRC, we model the evolution over time of each longitudinal covariate  $Y_s$  through a linear regression model

$$y_{si} = W_{si}\beta_s + Z_{si}u_{si} + \varepsilon_{si}, \quad i \in \mathcal{I}(t_L), \quad (2)$$

where  $\beta_s$  is a vector of fixed effect parameters,  $u_{si} \sim N(0, D_s)$  is a vector of random effects,  $\varepsilon_{si} \sim N(0, \sigma_s^2 I_{m_i})$  is the error term vector, and  $W_{si}$  and  $Z_{si}$  are design matrices associated to  $\beta_s$  and  $u_{si}$ . As an example, later in this article we will consider an example where we let  $y_{si}$  depend on the age  $a_{ij}$  of subject  $i$  at each visit, and include a random intercept and random slope in the LMM:

$$y_{sij} = \beta_{s0} + u_{s0} + \beta_{s1}a_{ij} + u_{s1}a_{ij} + \varepsilon_{sij}, \quad (3)$$

where  $(u_{s0}, u_{s1})^T \sim N(0, D_s)$  is a vector of random effects that follows a bivariate normal distribution. We employ maximum likelihood (ML) estimation to estimate  $\beta_s$ ,  $D_s$  and  $\sigma_s^2$  in model (2).

In the second step of PRC, we use the ML estimates from step 1 to derive summaries of the individual longitudinal trajectories for each biomarker. These are the predicted random effects, which can be computed as

$$\hat{u}_{si} = E(u_{si}|Y_{si} = y_{si}) = \hat{D}_s Z_{si}^T \hat{V}_{si}^{-1} (y_{si} - X_{si} \hat{\beta}_s), \quad (4)$$

where  $\hat{V}_{si} = Z_{si} \hat{D}_s Z_{si}^T + \hat{\sigma}_s^2 I$ .

In the third step of PRC, we model the relationship between the survival outcome and the baseline and longitudinal predictors. This is achieved through the specification of a Cox model where we

include the baseline predictors  $x_i$  and the summaries of the longitudinal predictors  $\hat{u}_i = (\hat{u}_{1i}, \dots, \hat{u}_{pi})$  as covariates:

$$h(t_i|x_i, \hat{u}_i) = h_0(t_i) \exp(\gamma x_i + \delta \hat{u}_i), \quad (5)$$

where  $h_0(t_i)$  is the baseline hazard function, and  $\gamma$  and  $\delta$  are vectors of regression coefficients. Since our approach allows for the inclusion of a potentially large number of baseline and longitudinal covariates, we estimate model (5) using penalized maximum likelihood (PML, Verweij and Van Houwelingen (1994)). As penalties we consider the ridge ( $L_2$ ), lasso ( $L_1$ ) and elasticnet penalties. The elasticnet penalty for model (5) is given by

$$\lambda \left[ \alpha \left( \sum_{s=1}^p |\gamma_s| + \sum_{s=1}^p |\delta_s| \right) + (1 - \alpha) \left( \sum_{s=1}^p \gamma_s^2 + \sum_{s=1}^p \delta_s^2 \right) \right], \quad (6)$$

where  $\lambda \geq 0$  and  $\alpha \in [0, 1]$ . The ridge penalty is obtained by setting  $\alpha = 0$ , and the lasso penalty by fixing  $\alpha = 1$ .

## 2.4 Computation of the predicted survival probabilities

Once models (2) and (5) have been estimated, the predicted survival probabilities  $S_i(t|t_L, x_i, \mathcal{Y}_i(t_L))$  are computed using

$$\hat{S}_i(t|t_L, x_i, \mathcal{Y}_i(t_L)) = \exp \left( - \int_0^t \hat{h}_0(s) \exp(\hat{\gamma} x_i + \hat{\delta} \hat{u}_i) \right), \quad (7)$$

where  $\hat{h}_0(s)$  is the estimated baseline hazard function,  $\hat{\gamma}$  and  $\hat{\delta}$  are the PML estimates of  $\gamma$  and  $\delta$  obtained in step 3, and  $\hat{u}_i$  contains the predicted random effects computed in step 2.

For subjects  $i \in \mathcal{I}(t_L)$ , who are included in the training set and survived up until  $t_L$ , computation of (7) is straightforward, since the predicted random effects  $\hat{u}_i$  for such subjects have already been computed in step 2. Predictions of  $S_i(t|t_L, x_i, \mathcal{Y}_i(t_L))$  for a new subject  $i = n + 1$  who survived up until  $t_L$ , but was not part of the training set is a bit more complex: before computing (7), one first needs to compute the predicted random effects for this new subject using (4). Note that such computation is feasible if and only if measurements of both baseline and longitudinal covariates (up to  $t_L$ ) are available for this new subject.

## 2.5 Evaluation of the predictive performance

We consider the time-dependent area under the ROC curve (tdAUC, Heagerty et al. (2000)), the concordance index or C index (Pencina and D'Agostino, 2004) and the Brier score (Graf et al., 1999) as measures of predictive performance. To obtain unbiased estimates of these performance measures, Signorelli et al. (2021) proposed a cluster bootstrap optimism correction procedure (CBOCP) that generalizes the use of the bootstrap as internal validation method to problems involving repeated measurement data. As an alternative to the CBOCP, one may choose to implement a cross-validation approach instead. Should the user opt for such an alternative, we recommend the use of repeated cross-validation over simple cross-validation to achieve a level of accuracy comparable to that of the CBOCP.

## 3 The R package pncal

In this Section we introduce the functions for the estimation of PRC, the computation of the predicted survival probabilities and the validation of predictive performance, providing an overview of the relevant estimation approaches and some important implementation details.

Table 1 provides a side-by-side overview of the functions that can be used to implement the PRC LMM and PRC MLPMM approaches. Note that while two different functions (one for each approach) are needed for the three estimation steps and the computation of the survival probabilities, the evaluation of the predictive performance is implemented in a single function that works with inputs from both approaches.

**Table 1:** Overview of the penca1 functions that implement the different modelling steps for the PRC LMM and PRC MLPMM approaches.

Task	PRC LMM	PRC MLPMM
Step 1: estimate the mixed-effects models	<code>fit_lmms</code>	<code>fit_mlpmms</code>
Step 2: compute the predicted random effects	<code>summarize_lmms</code>	<code>summarize_mlpmms</code>
Step 3: estimate the penalized Cox model	<code>fit_prclmm</code>	<code>fit_prcmlpmm</code>
Computation of predicted survival probabilities	<code>survpred_prclmm</code>	<code>survpred_prcmlpmm</code>
Evaluation of predictive performance	<code>performance_prc</code>	<code>performance_prc</code>

### 3.1 Model estimation and prediction

The first step of PRC involves the estimation of mixed-effects models for the longitudinal outcomes. For the PRC LMM approach, this can be done through the `fit_lmms` function, that proceeds to estimate  $p$  LMMs (one LMM for each of the longitudinal outcomes). Estimation of the LMMs is performed by maximum likelihood through the `lme` function from the R package `nlme` (Pinheiro and Bates, 2000). For the PRC MLPMM approach, the first step involves estimating one MLPMM for each group of longitudinal covariates. Estimation of the MLPMMs is done by maximum likelihood using the modified Marquardt algorithm described in Proust-Lima et al. (2013), as implemented in the `mult1cmm` from the R package `lcmm` (Proust-Lima et al., 2017).

The second step of PRC requires the computation of the predicted random effects. The function `summarize_lmms` implements this for the PRC LMM approach. The function takes the output of `fit_lmms` as input, and proceeds to the computation of the predicted random effects using equation (4). Similarly, the function `summarize_mlpmms` does the same for the PRC MLPMM approach by taking the output of `fit_mlpmms` as input and computing the predicted random effects using the formula given in equation (4) of Signorelli et al. (2021).

The third step of PRC requires the estimation of a Cox model where the baseline covariates and the predicted random effects are used as covariates. Estimation of such model can be performed using the function `fit_prclmm` for the PRC LMM approach and `fit_prcmlpmm` for the PRC MLPMM approach. These functions proceed to the estimation of the aforementioned Cox model by penalized maximum likelihood through the function `cv.glmnet` from the R package `glmnet` (Simon et al., 2011). If the user chooses the ridge or lasso penalty, then the selection of the value of the tuning parameter  $\lambda$  is performed through cross-validation as implemented in `glmnet`. If, instead, the elasticnet penalty is used, `fit_prclmm` and `fit_prcmlpmm` proceed to perform a nested cross-validation procedure to jointly select the optimal values of the tuning parameters  $\alpha$  and  $\lambda$ .

Lastly, the function `survpred_prclmm` can be used to compute the predicted survival probabilities as described in equation (7) for the PRC LMM approach. The corresponding function for the PRC MLPMM approach is `survpred_prcmlpmm`.

### 3.2 Computation of the CBOCP

In penca1, the evaluation of the predictive performance of the fitted model is done by estimating the tDAUC, C index and Brier score through the CBOCP described in Signorelli et al. (2021). For both the PRC LMM and PRC MLPMM approaches, this can be done through the function `performance_prc`. The estimates of the tDAUC, C index and Brier score are computed using functions from the packages `survivalROC` (Heagerty and Saha-Chaudhuri, 2022), `survcomp` (Schröder et al., 2011) and `riskRegression` (Gerdts et al., 2023), respectively.

The computation of the CBOCP requires the choice of the number of bootstrap replicates  $B$  over which the model should be refitted. This can be done by specifying the argument `n.boots` inside the functions that implement the first step of PRC, namely `fit_lmms` for the PRC LMM approach and `fit_mlpmms` for the PRC MLPMM one. The supplied value of `n.boots` is stored in the output of such functions, and all subsequent functions inherit this value, automatically performing the computations necessary for the CBOCP.

If `n.boots = 0` (default), the CBOCP is not computed, and `performance_prc` only returns the naïve estimates of predictive performance. Values of `n.boots`  $\geq 1$  will trigger the computation of the CBOCP, and the output of `performance_prc` will additionally include the estimates of the optimism and the optimism-corrected performance measures. A typical value for  $B$  is 100, but in general we recommend setting  $B$  to a value between 50 and 200 (depending on computing time and the desired level of accuracy, one may also consider larger values of  $B$ ).

### 3.3 User-friendly parallelization

The computation of the CBOCP is by nature repetitive, as it requires to repeat steps 1, 2 and 3 over  $B$  bootstrap samples, and to compute predictions and performance measures both on the bootstrap sample and on the original dataset to estimate the optimism. Due to this repetitiveness, such computations can be easily parallelized to reduce computing time. Moreover, also the estimation of the  $p$  LMMs / MLPMMs in step 1 can be trivially parallelized.

With the goal of making it as easy as possible for users to parallelize such computations, the functions in Table 1 automatically parallelize the aforementioned computations using the `%dopar%` operator from the R package `foreach` (Microsoft and Weston, 2022). The user only needs to specify the number of cores they want to use for the computation using the argument `n.cores`; `pencal` will automatically take care of the parallelization.

### 3.4 Classes, methods and further functionalities

Besides the core functions introduced in Table 1, `pencal` comprises additional functions that are shortly described hereafter.

Three functions are used to simulate the data that are used in the package documentation to illustrate typical usage of `pencal`'s functions. The function `simulate_t_weibull` is used to generate survival times from a Weibull distribution using the inverse transformation method. The functions `simulate_prclmm_data` and `simulate_prclpmm_data` are used to generate data for the estimation of PRC LMM and PRC MLPMM, respectively.

S3 classes and methods are implemented for the outputs of each modelling step:

- step 1: `fit_lmms` and `fit_mlpmms` return objects of class `lmmfit` and `mlpmmfit`, respectively. As the number of longitudinal covariates increases, step 1 of PRC will involve the estimation of many mixed-effects models: to simplify the extraction of the estimates of each mixed model, we provide two summary functions, `summary.lmmfit` and `summary.mlpmmfit`;
- step 2: `summarize_lmms` and `summarize_mlpmms` return objects of class `ranefs`, for which a `summary.ranefs` function is available;
- step 3: `fit_prclmm` outputs an object of class `prclmm`, and `fit_prclpmm` one of class `prclpmm`. For both classes, summary methods (`summary.prclmm` and `summary.mlpmmpf`, respectively) are implemented.

Lastly, it may be sometimes of interest to compare the performance of PRC to that of either a penalized Cox model that only uses baseline values of all covariates, or a penalized Cox model with LOFC landmarking. The `pencox_baseline` function provides an interface to estimate these two models and to compute the associated CBOCP. Its output can be fed to the function `performance_pencox_baseline` to obtain the naïve and optimism-corrected estimates of the tdAUC, C index and Brier score for these two models.

## 4 Dynamic prediction with `pencal`: a step by step example

### 4.1 Loading `pbc2data`

To illustrate how to use `pencal` in practice, we employ data from a study from a clinical trial on primary biliary cholangitis (PBC) conducted by the Mayo Clinic from 1974 to 1984 (Murtaugh et al., 1994). The trial recorded the first of two survival outcomes, namely liver transplantation or death. In this example we focus our attention on the prediction of deaths, treating patients who underwent liver transplantation as right-censored. The data are available in `pencal` in a list called `pbc2data` that can be loaded as follows:

```
library(pencal)
data(pbc2data)
ls(pbc2data)
```

```
#> [1] "baselineInfo"      "longitudinalInfo"
```

`pbc2data` contains two data frames: `baselineInfo` records the survival information  $(t_i, \delta_i)$  and baseline covariates  $x_i$ , whereas `longitudinalInfo` contains repeated measurements of the longitudinal predictors  $y_i$ . For simplicity, we rename the two data frames as `sdata` and `ldata`:

```
sdata = pbc2data$baselineInfo
ldata = pbc2data$longitudinalInfo
```

## 4.2 Input data format

As detailed in the Statistical Methods section, estimation of PRC requires the following input data for each subject  $i = 1, \dots, n$ :

- a pair  $(t_i, \delta_i)$  providing the survival outcome for subject  $i = 1$ ;
- a vector of  $k$  baseline covariates  $x_i$ ;
- an  $m_i \times p$  matrix containing all repeated measurements of the  $p$  longitudinal predictors. Within this matrix each column corresponds to a predictor, and each row to the measurements  $y_{ij}$  collected at time  $t_{ij}$ ,  $j \in \{1, \dots, m_i\}$  for subject  $i$ ;
- any longitudinal covariate needed to construct the design matrices  $W_{si}$  and  $Z_{si}$  that will be used to estimate the LMMs of equation (2).

Such data should be provided to `pencal` using two data frames. The first data frame is a data frame that contains information on the survival outcomes  $(t_i, \delta_i)$  and the baseline covariates  $x_i$ . For the PBC2 example, this is the `sdata` data frame created above:

```
head(sdata)

#>   id      time event baselineAge   sex treatment
#> 1  1 1.095170     1    58.76684 female D-penicil
#> 3  2 14.152338     0    56.44782 female D-penicil
#> 12 3 2.770781     1    70.07447 male  D-penicil
#> 16 4 5.270507     1    54.74209 female D-penicil
#> 23 5 4.120578     0    38.10645 female placebo
#> 29 6 6.853028     1    66.26054 female placebo
```

This data frame should comprise at least 3 variables: a variable named `id` that contains subject identifiers, a variable named `time` containing the values of  $t_i$ , and a dummy variable named `event` that corresponds to  $\delta_i$  (`event = 1` for subjects who experience the event at  $t_i$ , and `event = 0` for right-censored observations). Additionally, it can also contain the baseline covariates  $x_i$  (if any); in the example we have  $k = 3$  baseline covariates: `baselineAge`, `sex` and `treatment`.

The second data frame is a dataset in long format that contains the repeated measurements of the longitudinal predictors  $y_{ij}$  and of any covariate needed to create the design matrices  $W_{si}$  and  $Z_{si}$ . In our example application, such information is stored in `ldata`:

```
head(ldata)

#>   id      age  fuptime serBilir serChol albumin alkaline SGOT platelets
#> 1  1 58.76684 0.0000000  14.5    261    2.60    1718 138.0    190
#> 2  1 59.29252 0.5256817  21.3     NA    2.94    1612  6.2    183
#> 3  2 56.44782 0.0000000  1.1     302    4.14    7395 113.5    221
#> 4  2 56.94612 0.4983025  0.8     NA    3.60    2107 139.5    188
#> 5  2 57.44716 0.9993429  1.0     NA    3.55    1711 144.2    161
#> 6  2 58.55054 2.1027270  1.9     NA    3.92    1365 144.2    122
#>   prothrombin
#> 1          12.2
#> 2          11.2
#> 3          10.6
#> 4          11.0
#> 5          11.6
#> 6          10.6
```

This “longitudinal” data frame should contain the following information:

- a variable named `id` that contains subject identifiers;
- a variable containing the time from baseline  $t_{ij}$  at which the measurement was collected. In `ldata`, we called this variable `fuptime` (short for: follow-up time). Notice that if the longitudinal covariates are measured at  $t_{ij} = 0$ , a row with `fuptime = 0` must be included in `ldata`;
- the  $p$  longitudinal predictors (`serBilir`, `serChol`, `albumin`, `alkaline`, `SGOT`, `platelets` and `prothrombin` in the example);
- the covariates needed to construct  $W_{si}$  and  $Z_{si}$  (`age` in the example).

### 4.3 Choice of the landmark time and data preparation

Before proceeding with the estimation of PRC and the computation of the predicted survival probabilities  $S_i(t|t_L, x_i, \mathcal{Y}_i(t_L))$ , three preliminary steps are needed. The first involves the choice of the landmark time  $t_L$ . Hereafter we choose  $t_L = 2$  years, meaning that we want to predict the survival probability  $S(t|t_L = 2, x_i, \mathcal{Y}_i(2))$  for a subject with baseline covariates  $x_i$  and longitudinal covariates measured up until two years from baseline  $\mathcal{Y}_i(2)$ .

```
# set the landmark time
lmark = 2
```

Once  $t_L$  has been chosen, the second step is to retain for analysis only those subjects that survived up until the landmark time, i.e. all  $i : t_i \geq t_L$ :

```
# remove subjects who had event / were censored before landmark
sdata = subset(sdata, time > lmark)
ldata = subset(ldata, id %in% sdata$id)
```

Lastly, only repeated measurements taken up to the landmark time ( $t_{ij} \leq t_L$ ) should be retained for modelling, whereas measurements taken after the landmark time ( $t_{ij} > t_L$ ) should be discarded:

```
# remove measurements taken after landmark:
ldata = subset(ldata, fuptime <= lmark)
```

### 4.4 Descriptive statistics, data visualization and transformation

After choosing  $t_L = 2$  as landmark time, the number of subjects retained for model estimation is 278, of which 107 experience the event of interest, whereas the remaining 171 are right-censored:

```
# number of subjects retained in the analysis:
nrow(sdata)

#> [1] 278

# number of events (1s) and censored observations (0s):
table(sdata$event)

#>
#>   0   1
#> 171 107
```

The estimated survival probability can be visualized through the Kaplan-Meier estimator in Figure 1 as shown below:

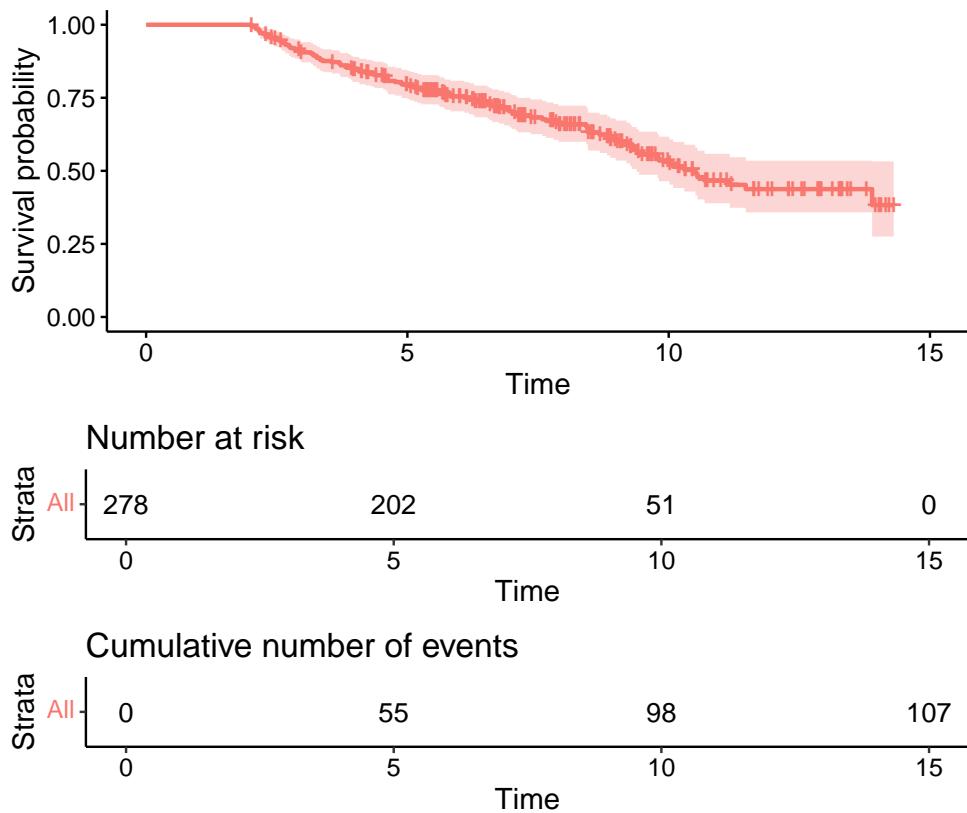
```
library(survival)
library(survminer)
surv.obj = Surv(time = sdata$time, event = sdata$event)
KM = survfit(surv.obj ~ 1, type = "kaplan-meier")
ggsurvplot(KM, data = sdata, risk.table = TRUE, cumevents = TRUE, legend = 'none')
```

Spaghetti plots displaying the trajectory described by a longitudinal covariate, as well as density plots to visualize its marginal distribution, can be created in ggplot2 style using:

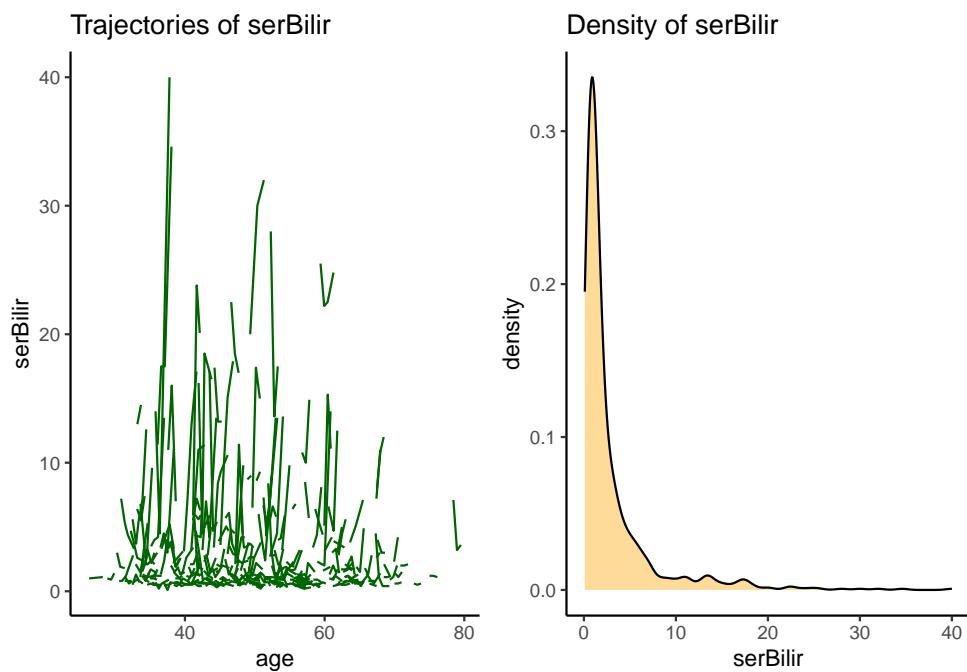
```
library(ggplot2)
library(gridExtra)
traj1 = ggplot(ldata, aes(x = age, y = serBilir, group = id)) +
  geom_line(color = 'darkgreen') + theme_classic() + ggtitle('Trajectories of serBilir')
dens1 = ggplot(ldata, aes(x = serBilir)) +
  geom_density(adjust=1.5, alpha=.4, fill = 'orange') + theme_classic() +
  ggtitle('Density of serBilir')
grid.arrange(traj1, dens1, ncol = 2)
```

The two charts thus created are shown in Figure 2.

We can observe that some longitudinal covariates exhibit strong skewness, as in the case of serBilir. Although in principle the LMM can be used to model variables with skewed distributions, this may sometimes lead to converge problems or poor model fit. It can thus be advisable to transform such covariates to prevent these problems (but note that this is not a required modelling choice, and one may alternatively choose to avoid such transformation). For this reason, we log-transform those longitudinal covariates with skewed distribution before modelling them with LMMs:



**Figure 1:** Chart displaying the Kaplan-Meier estimates of the conditional survival probability  $S(t|2)$ .



**Figure 2:** Spaghetti and density charts for the variable serBilir.

```
lndata$logSerBilir = log(lndata$serBilir)
lndata$logSerChol = log(lndata$serChol)
lndata$logAlkaline = log(lndata$alkaline)
lndata$logSGOT = log(lndata$SGOT)
lndata$logProthrombin = log(lndata$prothrombin)
```

## 4.5 Model estimation

### Step 1: estimation of the mixed-effects models

The first step in the estimation of PRC involves modelling the evolution over time of the longitudinal predictors through mixed-effects models. Hereafter we model each of the longitudinal predictors using the LMM given in equation (3), which comprises a random intercept and a random slope for age. Estimation of this model for each longitudinal predictor can be done using the function `fit_lmms`:

```
lmms = fit_lmms(y.names = c('logSerBilir', 'logSerChol', 'albumin', 'logAlkaline',
                           'logSGOT', 'platelets', 'logProthrombin'),
                  fixefs = ~ age, ranefs = ~ age | id, t.from.base = fuptime,
                  long.data = lndata, surv.data = sdata, n.boots = 0, n.cores = 8,
                  verbose = FALSE)
```

The argument `y.names` is a character vector used to specify the names of the longitudinal predictors in `lndata`. `fixefs` and `ranefs` are formulas used to specify the fixed and random effects part of the LMM using the `nlme` formula notation (Pinheiro et al., 2022; Galecki and Burzykowski, 2013). In the example, `fixefs = ~ age` determines the inclusion of the fixed effects part  $\beta_{s0} + \beta_{s1}a_{ij}$  of model (3), and `ranefs = ~ age | id` the inclusion of the random effects part  $u_{si0} + u_{si1}a_{ij}$  (allowing the random intercept and random slopes to be correlated).

The arguments `long.data` and `surv.data` are used to provide the names of the data frames containing the longitudinal variables (`long.data`) and the survival data and baseline covariates (`surv.data`) in the data formats described previously. The argument `t.from.base` is used to specify the name of the variable in `long.data` that contains the values of time from baseline; this argument is used internally to check that the data have been landmarked properly.

The `n.boots` argument is used to specify the number of bootstrap samples to use for the CBOCP. For the time being we focus on model estimation and on the prediction of the conditional probabilities, setting `n.boots = 0`; later we will show how to compute the CBOCP by setting `n.boots = 50`. Lastly, `n.cores` allows to specify the number of cores to use to parallelize computations (the default is `n.cores = 1`, i.e. no parallelization), and `verbose` is a logical value that indicates whether information messages should be printed in the console (TRUE, default) or not (FALSE). Additional arguments are described in the help page, see `?fit_lmms`.

The parameter estimates from the mixed models can be obtained from the output of step 1 using `summary`. For example, to obtain the parameters of the LMM for `albumin` we can use:

```
summary(lmms, yname = 'albumin', what = 'betas')

#> (Intercept)           age
#> 3.822741450 -0.005710811

summary(lmms, yname = 'albumin', what = 'variances')

#> id = pdLogChol(age)
#>          Variance     StdDev      Corr
#> (Intercept) 8.864945e-02 0.2977405761 (Intr)
#> age         3.447614e-07 0.0005871639 -0.103
#> Residual    1.257161e-01 0.3545646671
```

From the output we can deduce that the ML estimates for the LMM involving `albumin` are  $\hat{\beta} = (3.823, -0.0057)$ ,  $\hat{\sigma}_{u0} = 0.2977$ ,  $\hat{\sigma}_{u1} = 0.00059$ , and  $\hat{\sigma}_{u0,u1} = -0.103 \cdot 0.2977 \cdot 0.00059$ . The usual table with parameter estimates, standard errors and p-values can be obtained with

```
summary(lmms, yname = 'albumin', what = 'tTable')

#>             Value Std.Error DF   t-value      p-value
#> (Intercept) 3.822741450 0.106163343 566 36.008111 1.614057e-148
#> age        -0.005710811 0.002085903 566 -2.737812  6.379534e-03
```

## Step 2: computation of the predicted random effects

After the ML estimates from the LMM have been computed, the second step of PRC involves the computation of the predicted random effects  $\hat{u}_{si}$ . Such computation can be performed using the function `summarize_lmms`:

```
pred_ranefs = summarize_lmms(object = lmms, n.cores = 8, verbose = FALSE)
summary(pred_ranefs)

#> Number of predicted random effect variables: 14
#> Sample size: 278
```

Here, the `object` argument is used to pass the output of `fit_lmms` to `summarize_lmms`; `n.cores` and `verbose` are the same as in `fit_lmms`. The output of `summarize_lmms` is a list that contains, among other elements, a matrix with the predicted random effects called `ranef.orig`, where the row names display the subject identifiers:

```
round(pred_ranefs$ranef.orig[1:4, 1:4], 4)

#>   logSerBilir_b_int logSerBilir_b_age logSerChol_b_int logSerChol_b_age
#> 2      -0.3830        -0.0017       -0.0712        0.0007
#> 3     -0.1171        -0.0006       -0.5985        0.0049
#> 4      0.1686        0.0009       -0.3704        0.0035
#> 5      0.3800        0.0012       -0.2910        0.0029
```

From the output we can deduce, for example, that the predicted random intercept and random slope for `logSerBilir` and subject 4 are  $\hat{u}_{1,0,4} = 0.1686$  and  $\hat{u}_{1,1,4} = 0.0009$ .

## Step 3: estimation of the penalized Cox model

The last step in the estimation of PRC involves the estimation of model (5) through PML. This can be achieved with the `fit_prclmm` function:

```
pencox = fit_prclmm(object = pred_ranefs, surv.data = sdata,
                      baseline.covs = ~ baselineAge + sex + treatment, penalty = 'ridge',
                      standardize = TRUE, n.cores = 8, verbose = FALSE)
```

The `object` argument is used to pass the output of `summarize_lmms` to `fit_prclmm`; `surv.data` is the data frame that contains the information about survival data and baseline covariates; `baseline.covs` is a formula used to define which baseline covariates  $x_i$  should be included in model (5) (with associated regression coefficient  $\gamma$ ). The `penalty` argument is a character that can take one of the following values:

1. `penalty = 'ridge'` to estimate model (5) using the ridge or L2 penalty within the PML estimation;
2. `penalty = 'lasso'` to estimate model (5) using the lasso or L1 penalty;
3. `penalty = 'elnet'` to estimate model (5) using the elasticnet penalty (Zou and Hastie, 2005). If this penalty is chosen, additional arguments such as `n.alpha.elnet` and `n.folds.elnet` can be specified to determine how to select the additional tuning parameter ( $\alpha$ ) used by this penalty through nested cross-validation.

The `standardize` argument is used to determine whether the predicted random effects should be standardized prior to inclusion in the Cox model (default is `TRUE`). By default, `fit_prclmm` does not penalize baseline covariates, but this default behaviour can be changed using the argument `pfac.base.covs` argument (not shown here).

The `n.cores` and `verbose` arguments are the same as in `fit_lmms`. See `?fit_prclmm` for a description of further arguments.

The output of `fit_prclmm` can be summarized through `summary`:

```
summary(pencox)

#> Fitted model: PRC-LMM
#> Penalty function used: ridge
#> Tuning parameters:
#>   lambda alpha
#> 1 0.2126761    0
#> Sample size: 278
```

```
#> Number of events: 107
#> Bootstrap optimism correction: not computed
#> Penalized likelihood estimates (rounded to 4 digits):
#>   baselineAge sexfemale treatmentD-penicil logSerBilir_b_int logSerBilir_b_age
#> 1      0.0476    -0.2872       -0.0157      0.4341     111.3935
#>   logSerChol_b_int logSerChol_b_age albumin_b_int albumin_b_age
#> 1      0.0986     -10.5311      -1.1361     23070.92
#>   logAlkaline_b_int logAlkaline_b_age logSGOT_b_int logSGOT_b_age
#> 1      0.0874     -12.5617       0.238      272.246
#>   platelets_b_int platelets_b_age logProthrombin_b_int logProthrombin_b_age
#> 1      -0.0011     -0.2046      2.8114     -573.3093
```

The PML estimates of  $\gamma$  and  $\delta$  can be obtained through

```
step3summary = summary(pencox)
ls(step3summary)

#> [1] "coefficients" "data_info"      "model_info"    "tuning"

step3summary$coefficients

#>   baselineAge sexfemale treatmentD-penicil logSerBilir_b_int logSerBilir_b_age
#> 1  0.0475985 -0.287243      -0.01567369      0.4340672     111.3935
#>   logSerChol_b_int logSerChol_b_age albumin_b_int albumin_b_age
#> 1  0.0986092 -10.53106     -1.13607      23070.92
#>   logAlkaline_b_int logAlkaline_b_age logSGOT_b_int logSGOT_b_age
#> 1  0.08741668 -12.56169     0.2379553      272.246
#>   platelets_b_int platelets_b_age logProthrombin_b_int logProthrombin_b_age
#> 1  -0.001122804 -0.2046088      2.811446     -573.3093
```

## 4.6 Computing predictions

The function `survpred_prclmm` can be used to compute the conditional survival probabilities  $\hat{S}_i(t|t_L, x_i, \mathcal{Y}_i(t_L))$ ,  $t \geq t_L$  in (7). For the subjects that have been used to estimate PRC, such computation can be performed using `survpred_prclmm` as follows:

```
preds = survpred_prclmm(step1 = lmms, step2 = pred_ranefs, step3 = pencox, times = 3:7)
```

The `step1`, `step2` and `step3` arguments are used to pass the outputs of the 3 estimation steps to the function; the `times` argument is a vector with the prediction times at which one wishes to evaluate the conditional survival probabilities. The predicted survival probabilities are stored in the `predicted_survival` element of the function output:

```
ls(preds)

#> [1] "call"                  "predicted_survival"

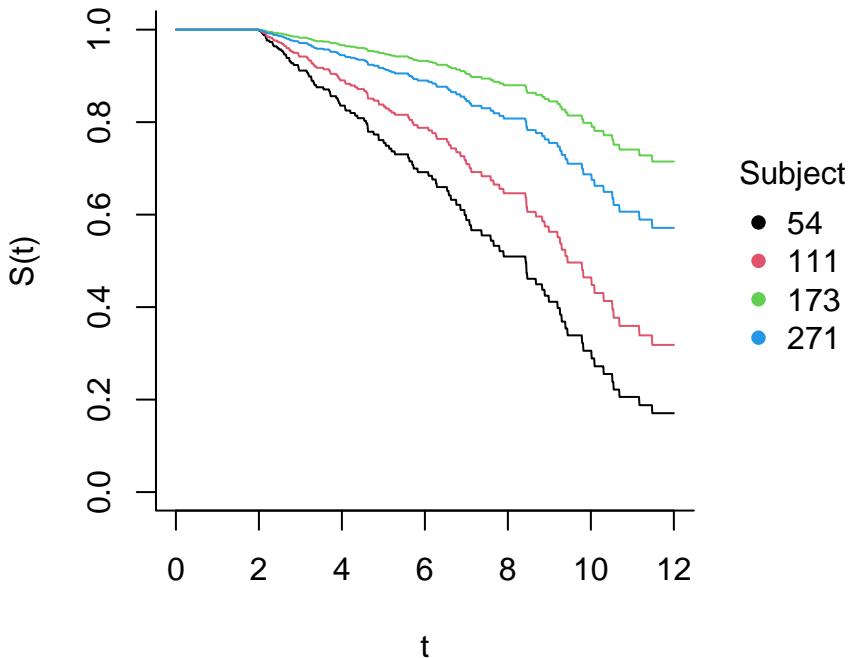
head(preds$predicted_survival)

#>   id      S(3)      S(4)      S(5)      S(6)      S(7)
#> 2  2 0.9398512 0.8867592 0.8329966 0.7813647 0.7008304
#> 3  3 0.8555809 0.7392053 0.6316391 0.5377700 0.4090876
#> 4  4 0.8138498 0.6709525 0.5451302 0.4407852 0.3071611
#> 5  5 0.9400431 0.8981124 0.8492655 0.8020394 0.7277052
#> 6  6 0.9383339 0.8839878 0.8290416 0.7763599 0.6943714
#> 7  7 0.9718724 0.9462254 0.9193945 0.8927316 0.8491681
```

The function `survplot_prc` allows to visualize predictions for a sample of individuals:

```
survplot_prc(step1 = lmms, step2 = pred_ranefs, step3 = pencox,
             ids = c(54, 111, 173, 271), tmax = 12)
```

The `ids` argument is used to indicate the subjects for whom the curve should be displayed, and `tmax` sets the upper limit for the  $x$  axis. The chart, displayed in Figure 3, shows the predicted survival probability  $\hat{S}_i(t|2)$  for subjects 54, 111, 173, and 271. Notice that the survival probability up to the 2 year landmark is 1 because our modelling approach conditions on being still at risk at the landmark time.



**Figure 3:** Predicted survival probabilities  $\hat{S}_i(t|2)$ ,  $t \in [2, 12]$  for subjects  $i \in \{54, 111, 173, 271\}$ .

Prediction for new subjects that have not been used for model estimation is a bit more involved, as it additionally requires to compute the predicted random effects for the new subjects using equation (4) based on the parameter estimates obtained in the first step of PRC. For illustration purposes, suppose that we have data from 3 subjects stored in two data frames called `new_ldata` and `new_sdata`:

```

new_ldata = subset(lpdata, id %in% c(5, 54, 110))
new_sdata = subset(sdata, id %in% c(5, 54, 110), c('id', 'baselineAge', 'sex', 'treatment'))
head(new_ldata)

#>      id     age   fuptime serBilir serChol albumin alkaline SGOT platelets
#> 23    5 38.10645 0.0000000    3.4     279    3.53      671 113.2      136
#> 24    5 38.65130 0.5448472    1.9      NA    3.28      689 103.9      114
#> 25    5 39.17698 1.0705290    2.5      NA    3.34      652 117.8      99
#> 419   54 39.19888 0.0000000   1.3     288    3.40      5487 73.5      254
#> 420   54 39.69171 0.4928266   1.5      NA    3.22      1580 71.3      112
#> 421   54 40.19001 0.9911291   2.6      NA    3.48      2127 86.8      207
#>      prothrombin logSerBilir logSerChol logAlkaline logSGOT logProthrombin
#> 23        10.9  1.2237754  5.631212  6.508769 4.729156  2.388763
#> 24        10.7  0.6418539      NA  6.535241 4.643429  2.370244
#> 25        10.5  0.9162907      NA  6.480045 4.768988  2.351375
#> 419       11.0  0.2623643  5.662960  8.610137 4.297285  2.397895
#> 420       18.0  0.4054651      NA  7.365180 4.266896  2.890372
#> 421       10.9  0.9555114      NA  7.662468 4.463607  2.388763

head(new_sdata)

#>      id baselineAge   sex treatment
#> 23    5     38.10645 female placebo
#> 419   54    39.19888 female D-penicil
#> 859  110    38.91140 female D-penicil

```

Note that the variables and their variable type in `new_ldata` and `new_sdata` should be the same as in the `lpdata` and `sdata`; the only exception to this is that `new_sdata` does not need to contain information about survival, so unlike `sdata` it does not comprise the time and event variables.

To compute predicted probabilities for new subjects, it is once again possible to resort to `survpred_prclmm`; now, it is necessary to specify the arguments `new.longdata` and `new.basecovs` to supply the data about the new subjects to `survpred_prclmm`:

```
pred_new = survpred_prclmm(step1 = lmms, step2 = pred_ranefs, step3 = pencox, times = 3:7,
```

```

new.longdata = new_ldata, new.basecovs = new_sdata)
pred_new$predicted_survival

#>      id      S(3)      S(4)      S(5)      S(6)      S(7)
#> 5    5 0.9460431 0.8981124 0.8492655 0.8020394 0.7277052
#> 54   54 0.9115188 0.8357020 0.7611779 0.6918065 0.5880763
#> 110  110 0.9505706 0.9064581 0.8612933 0.8174139 0.7478898

```

#### 4.7 Evaluation of the predictive performance

As explained in the Statistical Methods Section, estimation of the predictive performance in `pencal` is done through a CBOCP that allows to obtain unbiased estimates of predictive performance as measured by the tdaUC, C index and Brier score.

Before computing the (potentially time-consuming) CBOCP, one may want to first have a look at the naïve (biased) estimates of predictive performance. This can be done using the function `performance_prc`:

```

# naive performance (biased, optimistic estimate)
naive_perf = performance_prc(step2 = pred_ranefs, step3 = pencox, metric = 'tdauc',
                             times = 3:7, n.cores = 8, verbose = FALSE)

#> Warning in performance_prc(step2 = pred_ranefs, step3 = pencox, metric =
#> "tdauc", : The cluster bootstrap optimism correction has not been performed
#> (n.boots = 0). Therefore, only the apparent values of the performance values
#> will be returned.

```

Here `pred_ranefs` is the output of `step2` of PRC and `pencox` the output of step 3. The `metric` argument can be used to specify the performance measures to be computed (possible values are `tdauc`, `c` and `brier`), whereas the `times` argument is used to specify the time points at which the tdaUC and Brier score should be evaluated ( $t = 3, 4, 5, 6, 7$  in this example). Notice that when `fit_lmms` has been run with `n.boots = 0`, `performance_prc` returns a warning to inform users that the CBOCP has not been performed; for now, we can ignore this warning (which we will address soon by refitting PRC with `n.boots = 50`).

```

naive_perf

#> $call
#> performance_prc(step2 = pred_ranefs, step3 = pencox, metric = "tdauc",
#>                   times = 3:7, n.cores = 8, verbose = FALSE)
#>
#> $tdAUC
#>   pred.time tdaUC.naive optimism.correction tdaUC.adjusted
#> 1       3     0.9439             NA             NA
#> 2       4     0.9351             NA             NA
#> 3       5     0.9266             NA             NA
#> 4       6     0.8981             NA             NA
#> 5       7     0.8831             NA             NA

```

From the output we can observe that the naïve estimate of the tdaUC ranges from 0.9439 for predictions of survival at  $t = 3$  up to 0.8831 for predictions at  $t = 7$ . These naïve (in-sample) measurements of predictive performance may be optimistically biased due to overfitting, i.e., the fact that they are evaluated using the same data on which PRC was estimated. To correct for this potential source of bias, below we show how to implement the CBOCP to obtain unbiased estimates of the tdaUC and C index.

Computation of the CBOCP requires to repeat the 3 estimation steps of PRC for each bootstrap samples; this can be done by rerunning the functions `fit_lmms`, `summarize_lmms` and `fit_prclmm` with the same arguments used previously, but setting `n.boots` within `fit_lmms` to an integer value larger than 0. `n.boots` specifies the number of bootstrap samples to use to compute the CBOCP. In the example below we set `n.boots = 50` (note that larger values of `n.boots` can increase the accuracy of the CBOCP estimates, but at the same time they increase computing time).

```

step1 = fit_lmms(y.names = c('logSerBilir', 'logSerChol', 'albumin', 'logAlkaline',
                           'logSGOT', 'platelets', 'logProthrombin'),
                  fixefs = ~ age, ranefs = ~ age | id, t.from.base = fuptime,
                  long.data = ldata, surv.data = sdata, n.boots = 50, n.cores = 8,

```

```

    verbose = FALSE)
step2 = summarize_lmms(object = step1, n.cores = 8, verbose = FALSE)
step3 = fit_prclmm(object = step2, surv.data = sdata,
                     baseline.covs = ~ baselineAge + sex + treatment,
                     penalty = 'ridge', n.cores = 8, verbose = FALSE)

```

Once all computations are finished, it suffices to supply the refitted outputs of step 2 and step 3 to `performance_prc`:

```

# bootstrap-corrected performance (unbiased estimate)
cbocp = performance_prc(step2 = step2, step3 = step3, metric = c('tdauc', 'brier'),
                         times = 3:7, n.cores = 8, verbose = FALSE)
cbocp

#> $call
#> performance_prc(step2 = step2, step3 = step3, metric = c("tdauc",
#>     "brier"), times = 3:7, n.cores = 8, verbose = FALSE)
#>
#> $tdauc
#>   pred.time tdaUC.naive optimism.correction tdaUC.adjusted
#> 1      3     0.9439      -0.0061      0.9378
#> 2      4     0.9351      -0.0150      0.9201
#> 3      5     0.9266      -0.0132      0.9134
#> 4      6     0.8981      -0.0086      0.8895
#> 5      7     0.8831      -0.0118      0.8713
#>
#> $Brier
#>   pred.time Brier.naive optimism.correction Brier.adjusted
#> 1      3     0.0571      0.0147      0.0718
#> 2      4     0.0699      0.0271      0.0970
#> 3      5     0.0844      0.0328      0.1172
#> 4      6     0.0953      0.0348      0.1301
#> 5      7     0.1007      0.0416      0.1423

```

In the outputs above, the columns `tdauc.naive` and `Brier.naive` contain the naïve estimates of the tdAUC and Brier score; `optimism.correction` reports the values of the estimated optimism correction from the CBOCP for the two metrics; finally, `tdauc.adjusted` and `Brier.adjusted` contain the unbiased estimates of the tdAUC and Brier score.

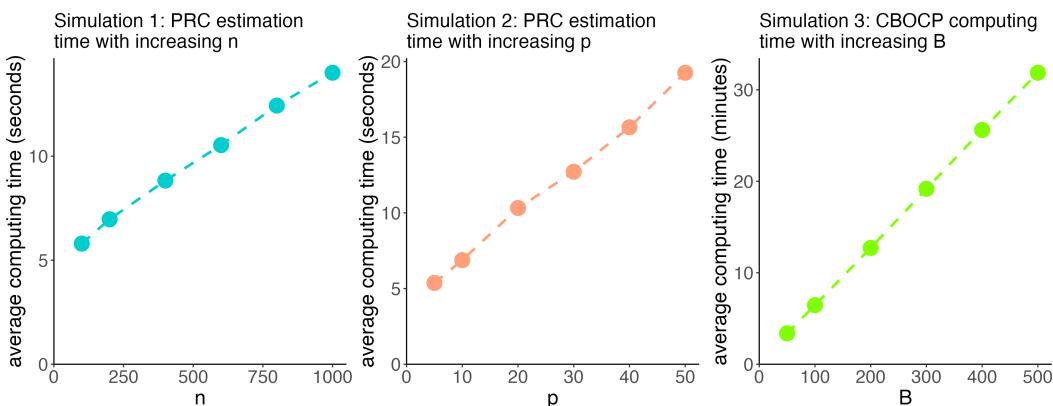
As expected, the unbiased estimates of predictive performance are somewhat worse than the naïve ones. For example, the tdAUC estimate for predictions at  $t = 3$  is 0.9378 instead of the naïve estimate 0.9439. Similarly, the Brier score estimate for predictions at  $t = 3$  is 0.0718 instead of the naïve estimate 0.0571.

## 5 Evaluation of computing time

We now turn our attention to the relationship between the sample size  $n$ , number of longitudinal covariates  $p$  and number of bootstrap replicates  $B$  on computing time. Furthermore, we look into how parallel computing may be used to reduce computing time for the CBOCP. To gain insight into these relationships, we simulate data from the PRC LMM model using the function `simulate_prclmm_data` according to four simulation scenarios:

- in simulation 1 we study the effect of  $n$  on the estimation of PRC. To this aim, we let  $n \in \{100, 200, 400, 600, 800, 1000\}$  and fix  $p = 10$ ;
- in simulation 2 we study the effect of  $p$  on the estimation of PRC by taking  $p \in \{5, 10, 20, 30, 40, 50\}$  and fixing  $n = 200$ ;
- in simulation 3 we shift our attention to the effect of  $B$  on the computing time of the CBOCP. We let  $B \in \{50, 100, 200, 300, 400, 500\}$ , fixing  $n = 200$  and  $p = 10$ ;
- finally, in simulation 4 we compute PRC and the CBOCP on a dataset where  $n = 200$ ,  $p = 50$  and  $B = 50$  using an increasing number of cores, namely  $\{1, 2, 3, 4, 8, 16\}$ .

Computations were performed on an AMD EPYC 7662 processor with 2 GHz CPU, using a single core for simulations 1, 2 and 3, and a number of cores ranging from 1 to 16 in simulation 4. Computing time was measured using the `rbenchmark` package (Kusnirczyk, 2012).



**Figure 4:** Left and center: average computing time (in seconds) of the estimation of the PRC LMM model as a function of the sample size  $n$  (simulation 1, left) and of the number of longitudinal predictors  $p$  (simulation 2, center). Right: average computing time (in minutes) for the computation of the CBOCP as a function of  $B$  (simulation 3).

The charts in Figure 4 display the average computing time over 10 replications of simulations 1, 2 and 3. In simulation 1, the average computing time increases from 5.80 seconds when  $n = 100$  to 14.01 seconds when  $n = 1000$ , so we observe a 2.4-fold increase in computing time as  $n$  increases by a factor of 10. In simulation 2, the average computing time increases from 5.37 seconds when  $p = 5$  to 19.26 seconds when  $p = 50$ , yielding a 3.6-fold increase in computing time as  $p$  increases by a factor of 10. Lastly, in simulation 3 the average computing time increases from 3.38 minutes when  $B = 50$  to 32.32 minutes when  $B = 500$ , with a 9.6-fold time increase corresponding to a 10-fold increase in  $B$ .

Overall, the results of simulations 1 and 2 indicate that computing time for the estimation of PRC increases linearly with, but less than proportionally to,  $n$  and  $p$  (meaning that a  $k$ -fold increase in  $n$  or  $p$  will typically increase computing time by a factor smaller than  $k$ ). Instead, from simulation 3 we can observe that the computing time of the CBOCP increases proportionally to  $B$ .

The results of simulations 1, 2 and 3 show that whereas the estimation of the PRC model itself typically requires only a few seconds, the computation of the CBOCP is more intensive and can require several minutes. This is due to the fact that the CBOCP requires the PRC modelling steps to be repeated on each bootstrap sample, effectively requiring to compute PRC  $B + 1$  times (once on the original dataset +  $B$  times on the  $B$  bootstrap datasets). To reduce the computing time needed to compute the CBOCP, `pencil` enables users to easily parallelize computations through the argument `n.cores` within `fit_lmms`, `summarize_lmms` and `fit_prclmm`. In simulation 4, we show the effect that increasing the number of cores has on the computing time of the CBOCP.

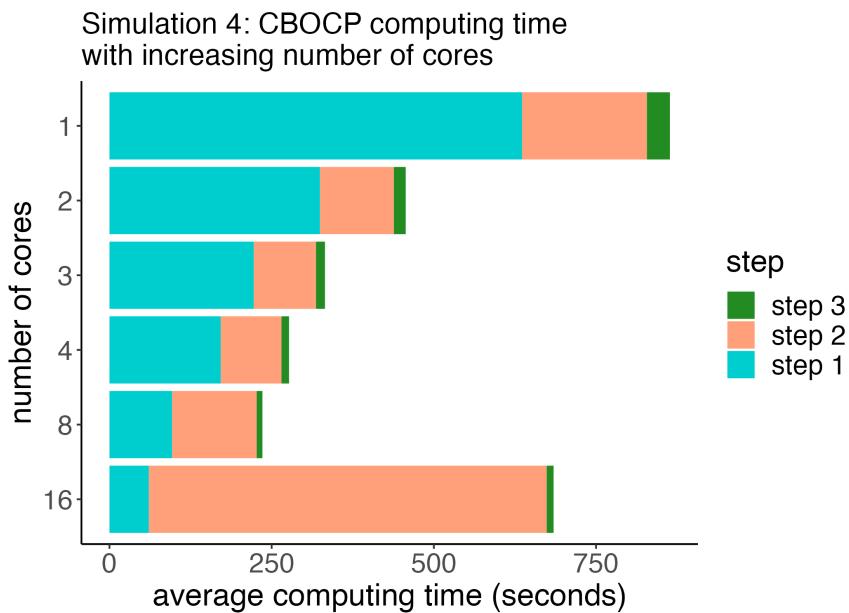
Figure 5 shows the average computing time of the CBOCP over 10 replications of simulation 4. When looking at the total computing time, we can see that increasing the number of cores from 1 to 8 progressively decreases computing time, reducing it from 863.8 seconds without parallelization up to 235.7 seconds when using 8 cores (-68%). The most significant time gains are from 1 to 2 cores (-40.7 seconds) and then from 2 to 3 (-124.4 seconds). Interestingly, further doubling the number of cores from 8 to 16 proves to be detrimental, increasing computing time from 235.7 to 684.5 seconds.

To understand this initially decreasing, but later increasing pattern, it is useful to consider the computing time of each of the modelling steps separately. By looking at Figure 5 we notice that step 1, which involves the estimation of  $p \cdot (B + 1)$  LMMs, is the most time-consuming step; its computing time consistently decreases as the number of cores decreases (from 636.1 to 60.8 seconds).

The same does not apply to step 2, where computing time decreases from 1 to 4 cores (from 192.3 to 93.6 seconds), but then increases considerably when going from 8 (130.4 seconds) to 16 cores (613.2 seconds). This pattern is primarily due to the fact that step 2 mostly involves simple linear algebra: parallelizing this step on a large number of cores may be detrimental, as the (limited) time gain that can be achieved by doing these simple computations in parallel may be more than compensated by the time cost of dispatching the necessary matrices and vectors to many cores and recombining the results at the end of the parallelization.

As concerns step 3, we can see that it is the lightest step in terms of computing time. The pattern is consistently decreasing from 1 (35.4 seconds) to 8 cores (8.8 seconds), with a slight increase when using 16 cores (10.6 seconds).

In conclusion, the results of simulation 4 show how it may be advisable to parallelize computations to compute the CBOCP, but without using an excessive number of cores (specially for step 2). Our



**Figure 5:** Average computing time (in seconds) for the estimation of the PRC LMM model and the computation of the CBOCP as a function of the number of cores.

general advice is to use between 3 and 8 cores for optimal performance, nevertheless we emphasize how the effect of the number of cores on computing time may differ from the patterns in Figure 5 depending on a combination of factors such as  $n$ ,  $p$ , number of repeated measurements per subject, and  $B$ . Furthermore, we notice that within `pencal`, a different number of cores can be chosen for each modelling step; in the example of simulation 4, the optimal performance would be achieved using 16 cores for step 1, 4 cores for step 2, and 8 cores for step 3 (but using 8 cores for all 3 steps isn't much less efficient).

## 6 Summary and discussion

The R package `pencal` provides a user-friendly implementation of Penalized Regression Calibration (PRC, Signorelli et al. (2021)), a statistical method that can be used to implement dynamic prediction of time-to-event outcomes in longitudinal studies where both time-independent and longitudinal (i.e., time-dependent) covariates are available as possible predictors of survival. The package comprises functions for the estimation of PRC and the prediction of survival, as well as functions to compute unbiased estimates of predictive performance through a cluster bootstrap procedure. Because computing such bootstrap procedure may be time-consuming, the package automatically parallelizes repetitive computations using the `%dopar%` operator from the `foreach` package (Microsoft and Weston, 2022).

`pencal` focuses on problems where a single survival outcome is measured with right-censoring. As such, it is not designed to handle interval censoring or competing risks. The modelling of the longitudinal covariates is performed using either the LMM or the MLPMM, which are linear models that are mostly suitable for the analysis of continuous outcomes. Implementing generalized linear mixed models (GLMMs) would make it possible to properly deal with binary and discrete longitudinal covariates, however the estimation of GLMMs and the computation of the predicted random effects are more time-consuming and more prone to convergence problems, two aspects that would particularly complicate the computation of the CBOCP. For this reason we did not pursue GLMMs further, but leave them as a topic of future research. Users dealing with discrete longitudinal covariates may consider log-transforming them before modelling with a LMM within `pencal` (specially if such covariates are right-skewed and/or exhibit overdispersion). Despite this latter limitation, a recent benchmarking study showed that PRC outperformed several alternative modelling approaches when applied to multiple real-world datasets (Signorelli and Retif, 2024).

Two modelling choices deserve particular attention when implementing PRC in specific application contexts. The first refers to the choice of the covariates to include in the fixed and random effects parts of the LMM of Equation (2). In principle, one may want to model the response variable as flexibly as possible, including several fixed effect covariates and multiple random effects in the LMM. However, when doing this one should consider that the purpose of the LMM is to provide subject-specific

summaries of the individual trajectories. Thus, *the primary goal of the LMM in step 1 is that of obtaining predicted random effects that are good summaries of how the trajectory of a given subject differs from the population average.* In practice, such purpose may be more easily achieved using a simple mixed model that allows for a clear interpretation of its random effects, rather than using a complex one where the interpretation of each random effect may be unclear / complicated. The LMM of equation (3) (or, alternatively, the same model with follow-up time included as covariate instead of age) is a good example of simple LMM with clearly interpretable random effects, as the random intercept allows to distinguish subjects with high and low initial levels of the covariate, and the random slope to identify subjects with faster and slower progression rates. Therefore, even though `fit_lmms` makes it possible to consider complex fixed and random effects formulas, we still advise users to consider simpler mixed models in step 1 (and to compare the predictive performance of PRC using either approach, eventually choosing the approach that delivers more accurate predictions if there is a substantial difference).

A second important modelling choice when using `pencal` is which penalty function should be used in step 3. In general, this is a modelling choice that may depend on the specific application and its features (sample size, number of predictors and number of available repeated measurements per subject). Signorelli et al. (2021) performed several simulation studies focused on situations with small sample sizes ( $n = 100$  and  $n = 300$ ) and sparse data generating processes for the survival outcome, whose results showed that the ridge and elasticnet penalty yielded better performance than the lasso penalty. Our experience is that in general the ridge penalty may be preferable both to elasticnet and the lasso in scenarios with small or moderate sample sizes, where little information is available to estimate the  $\alpha$  tuning parameter of elasticnet or to reliably perform variable selection with the lasso. Beyond this, it is always possible to use a data-driven approach to choose which penalty to use by estimating PRC using the 3 different penalties and comparing how this affects predictive performance.

## 6.1 Software and code availability

The R package `pencal` can be downloaded from CRAN at [cran.r-project.org/package=pencal](https://cran.r-project.org/package=pencal). The development version of the package is available on Github at [github.com/mirkosignorelli/pencal-devel](https://github.com/mirkosignorelli/pencal-devel). The code used in the simulations for the evaluation of computing time is available at [github.com/mirkosignorelli/pencal\\_sims/](https://github.com/mirkosignorelli/pencal_sims/).

## 6.2 Acknowledgements

The author kindly acknowledges funding from the Netherlands eScience Center Fellowship Programme.

## References

- A. Devaux, C. Proust-Lima, and R. Genuer. Random Forests for time-fixed and time-dependent predictors: The DynForest R package. *arXiv preprint arXiv:2302.02670*, 2023. [p135]
- A. Gałecki and T. Burzykowski. Linear mixed-effects model. In *Linear Mixed-Effects Models Using R*, pages 245–273. Springer, 2013. [p143]
- T. A. Gerds, J. S. Ohlendorff, and B. Ozenne. riskRegression: Risk Regression Models and Prediction Scores for Survival Analysis with Competing Risks, 2023. URL <https://CRAN.R-project.org/package=riskRegression>. R package version 2023.03.22. [p138]
- E. Graf, C. Schmoor, W. Sauerbrei, and M. Schumacher. Assessment and comparison of prognostic classification schemes for survival data. *Statistics in Medicine*, 18(17–18):2529–2545, 1999. [p137]
- P. J. Heagerty and P. Saha-Chaudhuri. survivalROC: Time-Dependent ROC Curve Estimation from Censored Survival Data, 2022. URL <https://CRAN.R-project.org/package=survivalROC>. R package version 1.0.3.1. [p138]
- P. J. Heagerty, T. Lumley, and M. S. Pepe. Time-dependent ROC curves for censored survival data and a diagnostic marker. *Biometrics*, 56(2):337–344, 2000. [p137]
- R. Henderson, P. Diggle, and A. Dobson. Joint modelling of longitudinal measurements and event time data. *Biostatistics*, 1(4):465–480, 2000. [p134]
- G. L. Hickey, P. Philipson, A. Jorgensen, and R. Kolamunnage-Dona. joineRML: a joint model and software package for time-to-event and multivariate longitudinal outcomes. *BMC Medical Research Methodology*, 18:1–14, 2018. [p134]

- W. Kusnirczyk. *rbenchmark*: Benchmarking routine for R, 2012. URL <https://CRAN.R-project.org/package=rbenchmark>. R package version 1.0.0. [p148]
- C. E. McCulloch and S. R. Searle. *Generalized, Linear, and Mixed Models*. John Wiley & Sons, 2004. [p136]
- Microsoft and S. Weston. *foreach*: Provides Foreach Looping Construct, 2022. URL <https://CRAN.R-project.org/package=foreach>. R package version 1.5.2. [p139, 150]
- P. A. Murtaugh, E. R. Dickson, G. M. Van Dam, M. Malinchoc, P. M. Grambsch, A. L. Langworthy, and C. H. Gips. Primary biliary cirrhosis: prediction of short-term survival based on repeated patient visits. *Hepatology*, 20(1):126–134, 1994. [p139]
- M. J. Pencina and R. B. D’Agostino. Overall C as a measure of discrimination in survival analysis: model specific population value and confidence interval estimation. *Statistics in Medicine*, 23(13): 2109–2123, 2004. [p137]
- P. Philipson, I. Sousa, P. Diggle, P. Williamson, R. Kolamunnage-Dona, R. Henderson, and G. Hickey. *joineR*: Joint Modelling of Repeated Measurements and Time-to-Event Data, 2018. URL <https://cran.r-project.org/package=joineR>. R package version 1.2.8. [p134]
- J. Pinheiro, D. Bates, and R Core Team. *nlme: Linear and Nonlinear Mixed Effects Models*, 2022. URL <https://CRAN.R-project.org/package=nlme>. R package version 3.1-160. [p143]
- J. C. Pinheiro and D. M. Bates. *Mixed-Effects Models in S and S-PLUS*. Springer, 2000. [p138]
- C. Proust-Lima, H. Amieva, and H. Jacqmin-Gadda. Analysis of multivariate mixed longitudinal data: a flexible latent process approach. *British Journal of Mathematical and Statistical Psychology*, 66(3): 470–487, 2013. [p136, 138]
- C. Proust-Lima, V. Philipps, and B. Liquet. Estimation of extended mixed models using latent classes and latent processes: The R package cmm. *Journal of Statistical Software*, 78(2):1–56, 2017. [p138]
- H. Putter and H. C. van Houwelingen. Landmarking 2.0: Bridging the gap between joint models and landmarking. *Statistics in Medicine*, 41(11):1901–1917, 2022. [p135]
- D. Rizopoulos. JM: An R package for the joint modelling of longitudinal and time-to-event data. *Journal of Statistical Software*, 35:1–33, 2010. [p134]
- D. Rizopoulos. The R package JMbayes for fitting joint models for longitudinal and time-to-event data using MCMC. *arXiv preprint arXiv:1404.7625*, 2014. [p134]
- M. S. Schröder, A. C. Culhane, J. Quackenbush, and B. Haibe-Kains. survcomp: an R/Bioconductor package for performance assessment and comparison of survival models. *Bioinformatics*, 27(22): 3206–3208, 2011. [p138]
- M. Signorelli. *pencal*: Penalized Regression Calibration (PRC) for the Dynamic Prediction of Survival, 2023. URL <https://cran.r-project.org/package=pencal>. R package version 2.1.0. [p135]
- M. Signorelli and S. Retif. An empirical appraisal of methods for the dynamic prediction of survival with numerous longitudinal predictors. *arXiv preprint arXiv:2403.14336*, 2024. [p150]
- M. Signorelli, P. Spitali, C. Al-Khalili Sgyziarto, The Mark-MD Consortium, and R. Tsonaka. Penalized regression calibration: a method for the prediction of survival outcomes using complex longitudinal and high-dimensional data. *Statistics in Medicine*, 40(27):6178–6196, 2021. [p135, 136, 137, 138, 150, 151]
- N. Simon, J. Friedman, T. Hastie, and R. Tibshirani. Regularization paths for Cox’s proportional hazards model via coordinate descent. *Journal of statistical software*, 39(5):1, 2011. [p138]
- E. W. Steyerberg. *Clinical Prediction Models*. Springer, 2009. [p134]
- T. M. Therneau and P. M. Grambsch. Modeling survival data: extending the Cox model. 2000. [p134]
- H. C. Van Houwelingen. Dynamic prediction by landmarking in event history analysis. *Scandinavian Journal of Statistics*, 34(1):70–85, 2007. [p134]
- P. J. Verweij and H. C. Van Houwelingen. Penalized likelihood in Cox regression. *Statistics in Medicine*, 13(23-24):2427–2436, 1994. [p137]
- H. Zou and T. Hastie. Regularization and variable selection via the elastic net. *Journal of the Royal Statistical Society: Series B*, 67(2):301–320, 2005. [p144]

*Mirko Signorelli*  
Mathematical Institute, Leiden University  
Niels Bohrweg 1, 2333 CA Leiden (NL)  
<https://mirkosignorelli.github.io>  
ORCiD: 0000-0002-8102-3356  
[msignorelli.papers@gmail.com](mailto:msignorelli.papers@gmail.com)

# clarify: Simulation-Based Inference for Regression Models

by Noah Greifer, Steven Worthington, Stefano Iacus, and Gary King

**Abstract** Simulation-based inference is an alternative to the delta method for computing the uncertainty around regression post-estimation (i.e., derived) quantities such as average marginal effects, average adjusted predictions, and other functions of model parameters. It works by drawing model parameters from their joint distribution and estimating quantities of interest from each set of simulated values, which form a simulated “posterior” distribution of the quantity from which confidence intervals can be computed. `clarify` provides a simple, unified interface for performing simulation-based inference for any user-specified derived quantities as well as wrappers for common quantities of interest. `clarify` supports a large and growing number of models through its interface with the `marginaleffects` package and provides native support for multiply imputed data.

## 1 Introduction

Although regression models are frequently used in empirical research to study relationships among variables, often the quantity of substantive interest is not one of the coefficients of the model, but rather a quantity derived from the coefficients, such as predicted values or average marginal effects. Quantifying the uncertainty of these derived quantities (i.e., computing standard errors, confidence intervals, and p-values) requires additional processing. Several methods of doing so exist, including the delta method, the bootstrap, and simulation-based inference. `clarify` implements simulation-based inference, which we describe below along with these other methods.

The delta method involves computing a first-order Taylor series approximation to the variance of the derived quantity, and standard Wald-based inference relies on computing quantiles based on the Normal distribution and using them to compute p-values and confidence intervals. `clarify` implements an alternative to the delta method — simulation-based inference — which involves simulating a “posterior” distribution of the derived quantities. Simulation-based inference does not require understanding Taylor series or the calculus that underlies them, which can make it more palatable to non-technical audiences and easier to learn for students without necessarily sacrificing statistical performance (King et al., 2000; Zelner, 2009). Some studies have found that simulation-based inference performs as well or better than the delta method for computing derived quantities (i.e., with respect to achieving close to nominal coverage for confidence intervals), especially for complicated derived quantities and in smaller samples (MacKinnon et al., 2004; Hole, 2007; Herron, 1999). Its empirical performance has been particularly well-studied in the context of mediation analysis, in which the quantities of interest are products and ratios of regression coefficients, where it has been shown to perform well relative to the delta method due to the non-Normality of these quantities (Tofghi and MacKinnon, 2016; Preacher and Selig, 2012).

The methodology `clarify` relies on was developed by Krinsky and Robb (1986) and is described in King et al. (2000) and Herron (1999). Simulation-based inference involves taking draws from a specified joint distribution of model parameters, computing derived quantities from these draws, and collecting the derived quantities in a “posterior” distribution, from which uncertainty measures (standard errors and confidence intervals) can be computed. This method assumes the model parameters are drawn from a multivariate Normal (or T) distribution with means at the estimated values and covariance equal to the asymptotic covariance matrix of the estimated values, a standard assumption motivated by the central limit theorem that underlies usual inference on the original model parameters. Arriving at the posterior distribution does not require taking any derivatives or making any approximations beyond those usually used for inference on model parameter estimates, except for the approximation due to Monte Carlo error induced by sampling from a finite number of simulations (which can always be reduced by increasing the number of draws at the cost of increased computing time).

The nonparametric bootstrap is another alternative to the delta method for inference that does not require its analytic approximations (Efron and Tibshirani, 1986); bootstrapping typically involves re-sampling individuals from the sample, fitting the model in each bootstrap sample, and computing the quantity of interest from each model. Although bootstrapping tends to work well in practice, especially for complex and non-Normal estimators, refitting the model repeatedly can be prohibitively time-consuming and computationally expensive, especially for complicated models or large datasets. Simulation-based inference only requires the model to be fit once, and the simulations involve taking draws from a distribution produced from the single set of estimated parameters, making it much quicker in practice and allowing the user to capitalize on the already valid estimation of the model pa-

rameters. Methods for computing valid confidence intervals in cases where the quantity of interest has a complicated distribution are better developed when bootstrapping; however (Efron and Tibshirani, 1986).

More formally, we fit a regression model  $y_i = f(x_i; \beta)$ , such as a linear or other generalized linear model with model coefficients  $\beta$ . We assume

$$\hat{\beta} \sim \text{MVN}(\beta, \Sigma_{\hat{\beta}})$$

where  $\hat{\beta}$  is the vector of estimates of  $\beta$  and  $\Sigma_{\hat{\beta}}$  is their asymptotic covariance matrix. We define a function  $\tau(\beta)$  that represents a quantity of interest derived from the model parameters, and compute its estimate  $\tau(\hat{\beta})$  as  $\tilde{\tau}$ . To perform simulation-based inference, we take  $M$  draws  $\tilde{\beta}^{(j)}$  for  $j \in (1, \dots, M)$  from a multivariate Normal distribution with mean vector  $\mu = \hat{\beta}$  and covariance  $\Sigma = \hat{\Sigma}_{\hat{\beta}}$ , where  $\hat{\Sigma}_{\hat{\beta}}$  is an estimate of the asymptotic covariance matrix of the parameter estimates. We use the distribution of  $\tilde{\tau} = \tau(\tilde{\beta})$  as the “posterior” distribution of  $\widehat{\tau(\beta)}$ , and compute its variance as

$$\widehat{\sigma}_{\widehat{\tau(\beta)}}^2 = \frac{1}{M-1} \sum_{j=1}^M (\tilde{\tau}^{(j)} - \bar{\tilde{\tau}})^2$$

and quantile  $100(1-\alpha)\%$  confidence interval limits as  $\left[ \tilde{\tau}_{\left(\frac{\alpha}{2}\right)}, \tilde{\tau}_{\left(1-\frac{\alpha}{2}\right)} \right]$ , where  $\tilde{\tau}_{(q)}$  is the  $q$ th value of  $\tilde{\tau}$  when arranged in ascending order (i.e., the  $q$ th quantile of the empirical cumulative distribution function). Simulation-based Wald-type confidence intervals can be computed as  $\left[ \tau(\hat{\beta}) + \widehat{\sigma}_{\widehat{\tau(\beta)}} Z_{\frac{\alpha}{2}}, \tau(\hat{\beta}) + \widehat{\sigma}_{\widehat{\tau(\beta)}} Z_{1-\frac{\alpha}{2}} \right]$ , where  $Z_q$  is the  $q$ th quantile of a standard Normal distribution. The delta method-based Wald-type confidence intervals use this formula but with the first-order Taylor approximation to the asymptotic variance:  $\widehat{\sigma}_{\widehat{\tau(\beta)}}^2 = \nabla \tau(\hat{\beta}) \Sigma_{\hat{\beta}} \nabla \tau'(\hat{\beta})$ , where  $\nabla \tau(\hat{\beta})$  is the gradient of  $\tau(\beta)$  with respect to  $\beta$  evaluated at  $\hat{\beta}$ .

To compute a p-value for a hypothesis test involving the quantity of interest, i.e.,  $H_0 : \tau(\beta) = \tau_0$  with a given null value  $\tau_0$ , we can invert the confidence interval (Thulin, 2021); that is, we find the largest value of  $\alpha$  such that  $\tau_0$  is within the confidence interval and use that  $\alpha$  as the p-value for the test. For Wald-based inference (either using the simulation-based variance or delta method-based variance), this is equivalent to performing a standard two-sided Z-test using the test statistic  $Z = (\tau(\hat{\beta}) - \tau_0) / \widehat{\sigma}_{\widehat{\tau(\beta)}}$ . One benefit of using the quantile p-values for inference is that equivalent tests of the same hypothesis will always yield identical p-values; for example, testing the equality of two derived quantities will yield the same p-value when comparing the difference between the quantities against a null hypothesis of 0 and the ratio of the quantities against a null hypothesis of 1, as each of these hypotheses is true if and only if the other is true.

One would expect simulation-based quantile inference, simulation-based Wald inference, and delta method-based Wald inference to align when the posterior is Normally distributed around the estimate, in which case any discrepancies would be due to Monte Carlo error in the simulated values (and therefore would shrink with increasing draws). However, for low values of  $\alpha$ , it may require many draws for the simulation-based intervals to stabilize; delta method-based intervals are not subject to this error. There are a few cases in which the results might diverge: in some cases, the first-order Taylor series approximation to the variance may be poor, though in practice the approximation error is small and shrinks quickly with increasing sample size. When the posterior distribution is non-Normal but symmetric around the estimate, the quantile intervals may be more accurate (i.e., in the sense of achieving closer to nominal coverage) because they do not rely on quantiles from the Normal distribution (Tofighi and MacKinnon, 2016).

Another potential advantage quantile intervals can have over Wald intervals is that when *some* monotonic transformation of the estimate has a symmetric distribution centered around the transformed estimate, the quantile intervals can achieve correct coverage without requiring knowledge of *which* transformation is required (Efron and Tibshirani, 1986); this is true of the quantile-based p-values as well. When the distribution is not centered around the estimate and no monotonic transformation will make it so, though, neither quantile-based nor Wald-based intervals would be expected to perform well, and quantile intervals could yield even worse coverage than Wald-based intervals, a phenomenon that occurs in the context of bootstrapping (Efron and Tibshirani, 1986)<sup>1</sup>. An informal falsification test for whether such a monotonic transformation exists is whether the median of the

<sup>1</sup>We thank an anonymous reviewer for pointing out a scenario in which this could occur: for a quantity of interest with a right-skewed sampling distribution, one would prefer an estimate to the right of the quantity's true value to have a confidence interval skewed to the left to capture the bulk of the sampling distribution, but in practice, a quantile confidence interval would also be skewed to the right. While a symmetric Wald-based interval may not have adequate coverage, the quantile-based interval could perform even worse.

simulated estimates is aligned with the point estimate; if it is not, there is no monotonic transformation that will yield a symmetric quantile interval with the desired coverage.

## 2 Related Software

Similar functionality exists in the **CLARIFY** package in Stata<sup>2</sup> (Tomz et al., 2003) and used to be available in the **Zelig** R package (Imai et al., 2008), though there are differences in these implementations. **clarify** provides additional flexibility by allowing the user to request any derived quantity, in addition to providing shortcuts for common quantities, including predictions at representative values, average marginal effects, and average dose-response functions (described below). **clarify** relies on and can be seen as a companion to the **marginaleffects** package (Arel-Bundock et al., Forthcoming), which offers similar functionality but primarily uses the delta method for calculating uncertainty (though simulation-based inference is supported in a more limited capacity as well).

## 3 Using **clarify**

There are four steps to using **clarify**:

1. Fit the model to the data using modeling functions in supported packages.
2. Use `sim()` to take draws from the multivariate distribution of the estimated model coefficients.
3. Use `sim_apply()` or its wrappers `sim_setx()`, `sim_ame()`, and `sim_adrf()` to compute derived quantities using each simulated set of coefficients.
4. Use `summary()` and `plot()` to summarize and visualize the distribution of the derived quantities and perform inference on them.

In the sections below, we will describe how to implement these steps in detail. First, we will load **clarify** using `library()`.

```
library(clarify)
```

For a running example, we will use the `lalonde` dataset in the **MatchIt** package (Ho et al., 2011), which contains data on 614 participants enrolled in a job training program or sampled from a survey (Dehejia and Wahba, 1999). The treatment variable is `treat`, and the outcome is `re78`, and all other variables are confounders. Although the original use of this dataset was to estimate the effect of `treat` on `re78`, we will use it more generally to demonstrate all of **clarify**'s capabilities. In addition, we will use a transformation of the outcome variable to demonstrate applications to nonlinear models, for which the benefits of simulation-based inference are more apparent.

```
data("lalonde", package = "MatchIt")

# Create a binary outcome variable
lalonde$re78_0 <- ifelse(lalonde$re78 > 0, 1, 0)

head(lalonde)

#>      treat age educ race married nodegree re74 re75        re78 re78_0
#> NSW1     1  37   11 black      1       1   0    0  9930.0460      1
#> NSW2     1  22    9 hispan     0       1   0    0  3595.8940      1
#> NSW3     1  30   12 black      0       0   0    0 24909.4500      1
#> NSW4     1  27   11 black      0       1   0    0  7506.1460      1
#> NSW5     1  33    8 black      0       1   0    0  289.7899      1
#> NSW6     1  22    9 black      0       1   0    0 4056.4940      1
```

<sup>2</sup>Despite the similar name, the R package **clarify** and the Stata package **CLARIFY** differ in several ways, one of which is that the estimates reported by **clarify** in R are those computed using the original model coefficients, whereas those reported by **CLARIFY** in Stata are those computed as the average of the simulated distribution. The R implementation avoids the “simulation-induced bias” described by Rainey (2023).

### 3.1 1. Fitting the model

The first step is to fit the model. `clarify` can operate on a large set of models (those supported by `marginalEffects`), including generalized linear models, multinomial models, multivariate models, and instrumental variable models, many of which are available in other R packages. Even if `clarify` does not offer direct support for a given model, there are ways to use its functionality regardless (explained in more detail below).

Because we are computing derived quantities, it is not critical to parameterize the model in such a way that the coefficients are interpretable, e.g., by using a model with interpretable coefficients or centering predictors. Below, we will fit a probit regression model for the outcome given the treatment and confounders. Coefficients in probit regression do not have a straightforward interpretation, but that does not matter; our quantities of interest can be expressed as derived quantities—functions of the model parameters, such as predictions, counterfactual predictions, and averages and contrasts of them.

```
fit <- glm(re78_0 ~ treat * married + age + educ + race +
            nodegree + re74 + re75, data = lalonde,
            family = binomial("probit"))
```

### 3.2 2. Drawing from the coefficient distribution

After fitting the model, we will use `sim()` to draw coefficients from their sampling distribution. The sampling distribution is assumed to be multivariate Normal or multivariate T with appropriate degrees of freedom, with a mean vector equal to the estimated coefficients and a covariance matrix equal to the asymptotic covariance matrix extracted from the model. The arguments to `sim()` are listed below:

```
sim(fit = , n = , vcov = , coefs = , dist = )
```

- `fit` – the fitted model object, the output of the call to the fitting function (e.g., `glm()`)
- `n` – the number of simulated values to draw; by default, 1000. More values will yield more replicable and precise results at the cost of speed.
- `vcov` – either the covariance matrix of the estimated coefficients, a function used to extract it from the model (e.g., `sandwich::vcovHC()` for the robust covariance matrix), or a string or formula giving a code for extracting the covariance matrix, which is passed to `marginalEffects::get_vcov()`. If left unspecified, the default covariance matrix will be extracted from the model.
- `coefs` – either a vector of coefficients to be sampled or a function to extract them from the fitted model. If left unspecified, the default coefficients will be extracted from the model. Typically this does not need to be specified.
- `dist` – the name of the distribution from which to draw the sampled coefficients. Can be "normal" for a Normal distribution or `t(#)` for a T-distribution, where # represents the degrees of freedom. If left unspecified, `sim()` will decide on which distribution makes sense given the characteristics of the model (the decision is made by `insight::get_df()` with `type = "wald"`). Typically this does not need to be specified.

If one's model is not supported by `clarify`, one can omit the `fit` argument and just specify the `vcov` and `coefs` arguments, which will draw the coefficients from the distribution named in `dist` ("normal" by default).

`sim()` uses a random number generator to draw the sampled coefficients from the sampling distribution, so a seed should be set using `set.seed()` to ensure results are replicable across sessions. Using more iterations (i.e., increasing `n`) yields results that will be more stable across runs even when a seed is not set.

The output of the call to `sim()` is a `clarify_sim` object, which contains the sampled coefficients, the original model fit object if supplied, and the coefficients and covariance matrix used to sample.

```
set.seed(1234)

# Drawing 1000 simulated coefficients using an HC2 robust
# covariance matrix
s <- sim(fit, n = 1000,
          vcov = "HC2")

s
```

```
#> A `clarify_sim` object
#> - 11 coefficients, 1000 simulated values
#> - sampled distribution: multivariate normal
#> - original fitting function call:
#>
#> glm(formula = re78_0 ~ treat * married + age + educ + race +
#>     nodegree + re74 + re75, family = binomial("probit"), data = lalonde)
```

### 3.3 3. Computing derived quantities

After sampling the coefficients, one can compute derived quantities on each set of sampled coefficients and store the result, which represents the “posterior” distribution of the derived quantity, as well as on the original coefficients, which are used as the final estimates. The core functionality is provided by `sim_apply()`, which accepts a `clarify_sim` object from `sim()` and a function to compute and return one or more derived quantities, then applies that function to each set of simulated coefficients. The arguments to `sim_apply()` are below:

```
sim_apply(sim = , FUN = , verbose = , cl = , ...)
```

- `sim` – a `clarify_sim` object; the output of a call to `sim()`.
- `FUN` – a function that takes in either a model fit object or a vector of coefficients and returns one or more derived quantities. The first argument should be named `fit` to take in a model fit object or `coefs` to take in coefficients.
- `verbose` – whether to display a progress bar.
- `cl` – an argument that controls parallel processing, which can be the number of cores to use or a cluster object resulting from `parallel::makeCluster()`.
- `...` – further arguments to `FUN`.

The `FUN` argument can be specified in one of two ways: either as a function that takes in a model fit object (e.g., a `glm` or `lm` object, the output of a call to `glm()` or `lm()`) or a function that takes in a vector of coefficients. The latter will always work but the former only works for supported models. When the function takes in a model fit object, `sim_apply()` will first insert each set of sampled coefficients into the model fit object and then supply the modified model to `FUN`.

For example, we will let our derived quantity of interest be the predicted probability of the outcome for participant PSID1. We specify our `FUN` function as follows:

```
sim_fun1 <- function(fit) {
  predict(fit, newdata = lalonde[["PSID1", ]], type = "response")
}
```

The `fit` object supplied to this function will be one in which the coefficients have been set to their values in a draw from their sampling distribution as generated by `sim()`. We then supply the function to `sim_apply()` to simulate the sampling distribution of the predicted value of interest:

```
est1 <- sim_apply(s, FUN = sim_fun1, verbose = FALSE)

est1

#> A `clarify_est` object (from `sim_apply()`)
#> - 1000 simulated values
#> - 1 quantity estimated:
#> PSID1 0.9757211
```

The resulting `clarify_est` object contains the simulated estimates in matrix form as well as the estimate computed on the original coefficients. We will examine the posterior distribution shortly, but first we will demonstrate computing a derived quantity from the coefficients directly.

The `race` variable is a factor, and the `black` category is used as the reference level, so it is not immediately clear whether there is a difference between the coefficients `racehispan` and `racewhite`, which represent the non-reference categories `hispan` and `white`. To compare these two directly, we can use `sim_apply()` to compute a derived quantity that corresponds to the difference between them.

```

sim_fun2 <- function(coefs) {
  hispan <- unname(coefs["racehispan"])
  white <- unname(coefs["racewhite"])

  c("w - h" = white - hispan)
}

est2 <- sim_apply(s, FUN = sim_fun2, verbose = FALSE)

est2

#> A `clarify_est` object (from `sim_apply()`)
#> - 1000 simulated values
#> - 1 quantity estimated:
#> w - h -0.09955915

```

The function supplied to FUN can be arbitrarily complicated and return as many derived quantities as one wants, though the slower each run of FUN is, the longer it will take to simulate the derived quantities. Using parallel processing by supplying an argument to cl can sometimes dramatically speed up evaluation.

There are several functions in `clarify` that serve as convenience wrappers for `sim_apply()` to automate some common derived quantities of interest. These include:

- `sim_setx()` – computing predicted values and first differences at representative or user-specified values of the predictors
- `sim_ame()` – computing average adjusted predictions, contrasts of average adjusted predictions, and average marginal effects
- `sim_adrf()` – computing average dose-response functions and average marginal effects functions

These are described in their own sections below. In addition, there are functions that have methods for `clarify_est` objects, including `cbind()` for combining two `clarify_est` objects together and `transform()` for computing quantities that are derived from the already-computed derived quantities. These are also described in their own sections below.

### 3.4 4. Summarize and visualize the simulated distribution

To examine the uncertainty around and perform inference on our estimated quantities, we can use `plot()` and `summary()` on the `clarify_est` object.

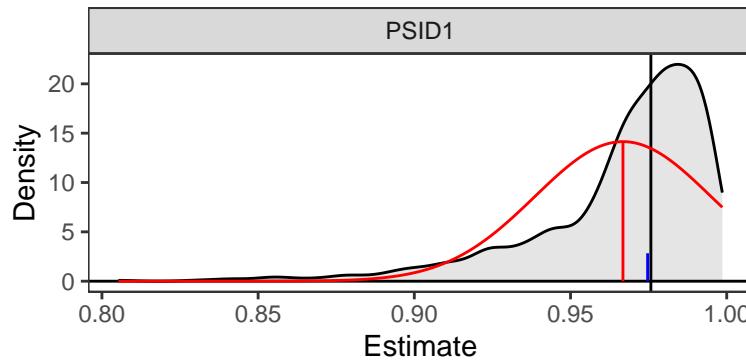
`plot()` displays a density plot of the resulting estimates across the simulations, with markers identifying the point estimate (computed using the original model coefficients as recommended by Rainey (2023)) and, optionally, uncertainty bounds (which function like confidence or credible interval bounds). The arguments to `plot()` are below:

```
plot(x = , parm = , ci = , level = , method = , reference =)
```

- `x` – the `clarify_est` object (the output of a call to `sim_apply()`).
- `parm` – the names or indices of the quantities to be plotted if more than one was estimated in `sim_apply()`; if unspecified, all will be plotted.
- `ci` – whether to display lines at the uncertainty bounds. The default is `TRUE` to display them.
- `level` – if `ci` is `TRUE`, the desired two-sided confidence level. The default is `.95` so that the bounds are at the `.025` and `.975` quantiles when `method` (see below) is `"quantile"`.
- `method` – if `ci` is `TRUE`, the method used to compute the bounds. Allowable methods include the Normal approximation (`"wald"`) or using the quantiles of the resulting distribution (`"quantile"`). The Normal approximation involves multiplying the standard deviation of the estimates (i.e., which functions like the standard error of the sampling distribution) by the critical Z-statistic computed using `(1-level)/2` to create a symmetric margin of error around the point estimate. The default is `"quantile"` to instead use quantile-based bounds.
- `reference` – whether to display a normal density over the plot for each estimate and an indicator line for the median of the estimate. The default is `FALSE` to omit them.

Below, we plot the first estimate we computed above, the predicted probability for participant PSID1:

```
plot(est1, reference = TRUE, ci = FALSE)
```



Overlaid on the plot in red is a Normal distribution with the same mean and standard deviation as the simulated values; this is requested by setting `reference = TRUE`. From the plot, one can see that the distribution of simulated values is non-Normal, asymmetrical, and not centered around the estimate, with no values falling above 1 because the outcome is a predicted probability. Given its non-Normality, the quantile-based bounds may be more appropriate than those resulting from the Normal approximation, as the bounds computed from the Normal approximation would be outside the bounds of the estimate. The blue reference line for the median of the estimates is close to the point estimate, suggesting it is possible for a monotonic transformation to have a symmetric distribution around the estimate<sup>3</sup>. The plot itself is a ggplot object that can be modified using ggplot2 syntax.

We can use `summary()` to display the value of the point estimate, the uncertainty bounds, and other statistics that describe the distribution of estimates. The arguments to `summary()` are below:

```
summary(object = , parm = , level = , method = , null = )
```

- `object` – the `clarify_est` object (the output of a call to `sim_apply()`).
- `parm` – the names or indices of the quantities to be displayed if more than one was estimated in `sim_apply()`; if unspecified, all will be displayed.
- `level` – the desired two-sided confidence level. The default is .95 so that the bounds are at the .025 and .975 quantiles when `method` (see below) is "quantile".
- `method` – the method used to compute the uncertainty bounds. Allowable methods include a Normal approximation ("wald") or using the quantiles of the resulting distribution ("quantile"). See `plot()` above.
- `null` – an optional argument specifying the desired null value in a hypothesis test for the estimates. If specified, a p-value will be computed using either a standard Z-test (if `method` is "wald") or an inversion of the uncertainty interval (described below). The default is not to display any p-values.

We can use `summary()` with the default arguments on our first `clarify_est` object to view the point estimate and quantile-based uncertainty bounds.

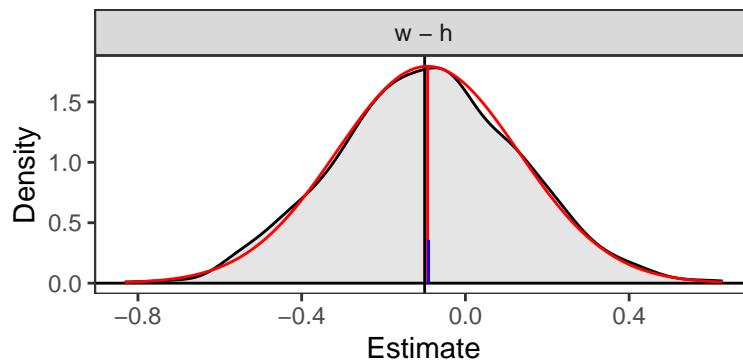
```
summary(est1)
```

```
#>      Estimate 2.5 % 97.5 %
#> PSID1    0.976  0.890  0.996
```

Our second estimated quantity, the difference between two regression coefficients, is closer to Normally distributed, as the plot below demonstrates (and would be expected theoretically), so we will use the Normal approximation to test the hypothesis that the difference differs from 0.

```
plot(est2, reference = TRUE, ci = FALSE)
```

<sup>3</sup>In fact, we know the inverse link function for the model (i.e., the Normal distribution function `qnorm()`) is such a transformation in this case; `marginaleffects` and other packages that implement the delta method for confidence intervals around model predictions typically automatically compute confidence intervals on the transformed predictions before transforming the intervals back using the model's link function. As long as such a transformation exists, the quantile intervals will be as valid as those that rely on transforming and back-transforming.



```
summary(est2, method = "wald", null = 0)

#>      Estimate 2.5 % 97.5 % Std. Error Z value P-value
#> w - h -0.0996 -0.5352 0.3361     0.2223 -0.45    0.65
```

The uncertainty intervals and p-values in the `summary()` output are computed using the Normal approximation because we set `method = "wald"`, and the p-value for the test that our estimate is equal to 0 is returned because we set `null = 0`. Note that the Normal approximation should be used only when the simulated posterior distribution is both close to Normal and centered around the estimate (i.e., when the mean of the simulated values [red vertical line] coincides with the estimate computed on the original coefficients [black vertical line]). In such cases, however, the delta method will likely perform as well, if not better, and all of its other benefits apply (i.e., it is computationally quicker and not subject to Monte Carlo error).

## 4 `sim_apply()` Wrappers: `sim_setx()`, `sim_ame()`, `sim_adrf()`

`sim_apply()` can be used to compute the simulated posterior distribution for an arbitrary derived quantity of interest, but there are some quantities that are common in applied research and may otherwise be somewhat challenging to program by hand, so `clarify` provides shortcut functions to make computing these quantities simple. These functions include `sim_setx()`, `sim_ame()`, and `sim_adrf()`. Each of these can be used only when regression models compatible with `clarify` are supplied to the original call to `sim()`.

Like `sim_apply()`, each of these functions is named `sim_*`(`), which signifies that they are to be used on an object produced by sim() (i.e., a clarify_sim object). (Multiple calls to these functions can be applied to the same clarify_sim object and combined; see the cbind() section below.) These functions are described below.`

### 4.1 `sim_setx()`: predictions at representative values

`sim_setx()` provides an interface to compute predictions at representative and user-supplied values of the predictors. For example, we might want to know what the effect of treatment is for a “typical” individual, which corresponds to the contrast between two model-based predictions (i.e., one under treatment and one under control for a unit with “typical” covariate values). This functionality mirrors the `setx()` and `setx1()` functionality of `Zelig` (which is where its name originates) and provides similar functionality to functions in `modelbased`, `emmeans`, `effects`, and `ggeffects`.

For each predictor, the user can specify whether they want predictions at specific values or at “typical” values, which are defined in `clarify` as the mode for unordered categorical and binary variables, the median for ordered categorical variables, and the mean for continuous variables. Predictions for multiple predictor combinations can be requested by specifying values that will be used to create a grid of predictor values, or the grid itself can be supplied as a data frame of desired predictor profiles. In addition, the “first difference,” defined here as the difference between predictions for two predictor combinations, can be computed.

The arguments to `sim_setx()` are as follows:

```
sim_setx(sim = , x = , x1 = , outcome = , type = , verbose = , cl = )
```

- `sim` – a `clarify_sim` object; the output of a call to `sim()`.

- $x$  – a named list containing the requested values of the predictors, e.g., `list(v1 = 1:4, v2 = "A")`, or a data frame containing the desired profiles. Any predictors not included will be set at their “typical” value as defined above.
- $x1$  – an optional named list or data frame similar to  $x$  except with the value of one predictor changed. When specified, the first difference is computed between the covariate combination defined in  $x$  (and only one combination is allowed when  $x1$  is specified) and the covariate combination defined in  $x1$ .
- `outcome` – a string containing the name of the outcome of interest when a multivariate (multiple outcome) model is supplied to `sim()` or the outcome category of interest when a multinomial model is supplied to `sim()`. For univariate (single outcome) and binary outcomes, this is ignored.
- `type` – a string containing the type of predicted value to return. In most cases, this can be left unspecified to request predictions on the scale of the outcome.
- `verbose` – whether to display a progress bar.
- `cl` – an argument that controls parallel processing, which can be the number of cores to use or a cluster object resulting from `parallel::makeCluster()`.

Here, we will use `sim_setx()` to examine predicted values of the outcome for control and treated units, at `re75` set to 0 and 20000, and `race` set to “black”.

```
est3 <- sim_setx(s,
  x = list(treat = 0:1,
            re75 = c(0, 20000),
            race = "black"),
  verbose = FALSE)
```

When we use `summary()` on the resulting output, we can see the estimates and their uncertainty intervals (calculated using quantiles by default).

```
summary(est3)

#>                               Estimate 2.5 % 97.5 %
#> treat = 0, re75 = 0          0.667 0.558 0.772
#> treat = 1, re75 = 0          0.712 0.617 0.790
#> treat = 0, re75 = 20000     0.938 0.700 0.994
#> treat = 1, re75 = 20000     0.953 0.747 0.996
```

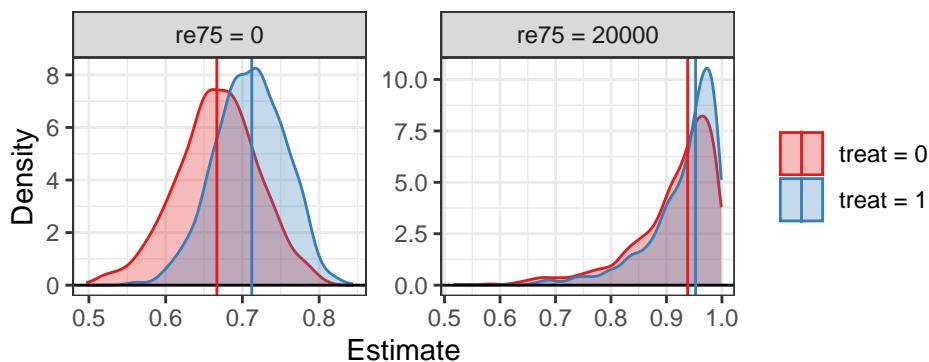
To see the complete grid of the predictor values used in the predictions, which helps to identify the “typical” values of the other predictors, we can access the “`setx`” attribute of the object:

```
attr(est3, "setx")

#>                               treat married      age      educ race nodegree      re74
#> treat = 0, re75 = 0          0        0 27.36319 10.26873 black      1 4557.547
#> treat = 1, re75 = 0          1        0 27.36319 10.26873 black      1 4557.547
#> treat = 0, re75 = 20000      0        0 27.36319 10.26873 black      1 4557.547
#> treat = 1, re75 = 20000      1        0 27.36319 10.26873 black      1 4557.547
#>                               re75
#> treat = 0, re75 = 0          0
#> treat = 1, re75 = 0          0
#> treat = 0, re75 = 20000     20000
#> treat = 1, re75 = 20000     20000
```

We can plot the distributions of the simulated values using `plot()`, which also separates the predictions by the predictor values (it is often clearer without the uncertainty bounds). The `var` argument controls which variable is used for facetting the plots.

```
plot(est3, var = "re75", ci = FALSE)
```



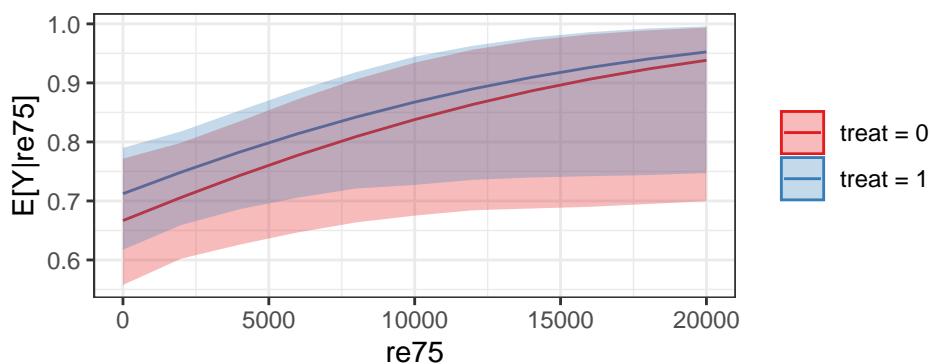
One can see again how a delta method or Normal approximation may have yielded uncertainty intervals outside the plausible range of the estimate without applying knowledge of the correct transformation to avoid doing so.

If a continuous variable with many levels is included in the grid of the predictors, something like a dose-response function for a typical unit can be generated. Below, we set `re75` to vary from 0 to 20000 in steps of 2000.

```
est4 <- sim_setx(s,
  x = list(treat = 0:1,
            re75 = seq(0, 20000, by = 2000),
            race = "black"),
  verbose = FALSE)
```

When we plot the output, we can see how the predictions vary across the levels of `re75`:

```
plot(est4)
```



We will return to display average dose-response functions using `sim_adrf()` later.

Finally, we can use `sim_setx()` to compute first differences, the contrast between two covariate combinations. We supply one covariate profile to `x` and another to `x1`, and `sim_setx()` simulates the two predicted values and their difference. Below, we simulate the first difference for a treated and control unit who has `re75` of 0 and typical values of all other covariates:

```
est5 <- sim_setx(s,
  x = list(treat = 0, re75 = 0),
  x1 = list(treat = 1, re75 = 0),
  verbose = FALSE)
```

When we use `summary()`, we see the estimates for the predicted values and their first difference ("FD"):

```
summary(est5)
```

```
#>           Estimate   2.5 % 97.5 %
#> treat = 0    0.7856  0.7039  0.8558
#> treat = 1    0.8213  0.7111  0.8995
#> FD          0.0357 -0.0598  0.1188
```

It is possible to compute first differences without using `x1` by using `transform()`, which we describe later.

## 4.2 `sim_ame()`: average adjusted predictions and average marginal effects

Using predicted values and effects at representative values is one way to summarize regression models, but another way is to compute average adjusted predictions (AAPs), contrasts of AAPs, and average marginal effects (AMEs). The definitions for these terms may vary and the names for these concepts differ across sources, but here we define AAPs as the average of the predicted values for all units after setting one predictor to a chosen value, and we define AMEs for binary predictors as the contrast of two AAPs and for continuous predictors as the average of the instantaneous rate of change in the AAP corresponding to a small change in the predictor from its observed values across all units<sup>4</sup> (Long and Freese, 2014).

The arguments to `sim_ame()` are as follows:

```
sim_ame(sim = , var = , subset = , by = , contrast = , outcome = ,
        type = , eps = , verbose = , cl = )
```

- `sim` – a `clarify_sim` object; the output of a call to `sim()`.
- `var` – the name of the focal variable over which to compute the AAPs or AMEs, or a list containing the values for which AAPs should be computed.
- `subset` – a logical vector, evaluated in the original dataset used to fit the model, defining a subset of units for which the AAPs or AMEs are to be computed.
- `by` – the name of one or more variables for which AAPs should be computed within subgroups. Can be supplied as a character vector of variable names or a one-sided formula.
- `contrast` – the name of an effect measure used to contrast AAPs. For continuous outcomes, "diff" requests the difference in means, but others are available for binary outcomes, including "rr" for the risk ratio, "or" for the odds ratio, and "nnt" for the number needed to treat, among others. If not specified, only AAPs will be computed if the variable named in `var` is categorical or specific values of the focal variable are specified in `var`. Ignored when the variable named in `var` is continuous and no specific values are specified because the AME is the only quantity computed. When `var` names a multi-category categorical variable, `contrast` cannot be used; see the section describing `transform()` for computing contrasts with them.
- `outcome` – a string containing the name of the outcome of interest when a multivariate (multiple outcome) model is supplied to `sim()` or the outcome category of interest when a multinomial model is supplied to `sim()`. For univariate (single outcome) and binary outcomes, this is ignored.
- `type` – a string containing the type of predicted value to return. In most cases, this can be left unspecified to request predictions on the scale of the outcome (e.g., probabilities for binary outcomes).
- `eps` – the value by which the observed values of the variable named in `var` are changed when it is continuous to compute the AME. This usually does not need to be specified.
- `verbose` – whether to display a progress bar.
- `cl` – an argument that controls parallel processing, which can be the number of cores to use or a cluster object resulting from `parallel::makeCluster()`.

Here, we will use `sim_ame()` to compute the AME of `treat` just among those who were treated (in causal inference, this is known as the average treatment effect in the treated, or ATT (Greifer and Stuart, 2023)). We will request our estimate to be on the risk ratio scale.

---

<sup>4</sup>In `marginaleffects`, AAPs are computed using `avg_predictions()`, AMEs for binary variables are computed using `avg_comparisons()`, and AMEs for continuous variables are computed using `avg_slopes()`. AAPs are sometimes known as average “counterfactual” predictions.

```
est6 <- sim_ame(s,
  var = "treat",
  subset = treat == 1,
  contrast = "rr",
  verbose = FALSE)
```

We can use `summary()` to display the estimates and their uncertainty intervals. Here, we will also use `null` to include a test for the null hypothesis that the risk ratio is equal to 1.

```
summary(est6, null = c(`RR` = 1))

#>           Estimate 2.5 % 97.5 % P-value
#> E[Y(0)]    0.687  0.608  0.760      .
#> E[Y(1)]    0.755  0.685  0.809      .
#> RR          1.100  0.949  1.255  0.21
```

Here we see the estimates for the AAPs,  $E[Y(0)]$  for the expected value of the outcome setting `treat` to 0 and  $E[Y(1)]$  for the expected value of the outcome setting `treat` to 1, and the risk ratio `RR`. The p-value on the test for the risk ratio aligns with the uncertainty interval containing 1.

If we instead wanted the risk difference or odds ratio, we would not have to re-compute the AAPs. Instead, we can use `transform()` to compute a new derived quantity from the computed AAPs. The section on `transform()` demonstrates this.

We can compute the AME for a continuous predictor. Here, we will consider `age` (just for demonstration; this analysis does not have a valid interpretation).

```
est7 <- sim_ame(s,
  var = "age",
  verbose = FALSE)
```

We can use `summary()` to display the AME estimate and its uncertainty interval.

```
summary(est7)

#>           Estimate 2.5 % 97.5 %
#> E[dY/d(age)] -0.00605 -0.00940 -0.00259
```

The AME is named  $E[dY/d(\text{age})]$ , which signifies that a derivative has been computed (more precisely, the average of the unit-specific derivatives). This estimate can be interpreted like a slope in a linear regression model, but as a single summary of the effect of a predictor it is often too coarse to capture nonlinear relationships. The section below explains how to compute average dose-response functions for continuous predictors, which provide a more complete picture of their effects on an outcome.

Below, we will examine effect modification of the ATT by the predictor `married` using the `by` argument to estimate AAPs and their ratio within levels of `married`:

```
est6b <- sim_ame(s,
  var = "treat",
  subset = treat == 1,
  by = ~married,
  contrast = "rr",
  verbose = FALSE)

summary(est6b)

#>           Estimate 2.5 % 97.5 %
#> E[Y(0)|0]    0.691  0.612  0.768
#> E[Y(1)|0]    0.733  0.655  0.796
#> RR[0]         1.061  0.909  1.234
#> E[Y(0)|1]    0.668  0.556  0.769
#> E[Y(1)|1]    0.848  0.676  0.940
#> RR[1]         1.270  0.948  1.583
```

The presence of effect modification can be tested by testing the contrast between the effects computed within each level of the `by` variable; this is demonstrated in the section on `transform()` below.

### 4.3 sim\_adrf(): average dose-response functions

A dose-response function for an individual is the relationship between the set value of a continuous focal predictor and the expected outcome. The average dose-response function (ADRF) is the average of the dose-response functions across all units. Essentially, it is a function that relates the value of the predictor to the corresponding AAP of the outcome, the average value of the outcome were all units to be set to that level of the predictor. ADRFs can be used to provide additional detail about the effect of a continuous predictor beyond a single AME.

A related quantity is the average marginal effect function (AMEF), which describes the relationship between a continuous focal predictor and the AME at that level of the predictor. That is, rather than describing how the outcome changes as a function of the predictor, it describes how the *effect* of the predictor on the outcome changes as a function of the predictor. It is essentially the derivative of the ADRF and can be used to identify at which points along the ADRF the predictor has an effect.

The ADRF and AMEF can be computed using `sim_adrf()`. The arguments are below:

```
sim_adrf(sim = , var = , subset = , by = , contrast = , at = ,
n = , outcome = , type = , eps = , verbose = , cl = )
```

- `sim` – a `clarify_sim` object; the output of a call to `sim()`.
- `var` – the name of the focal variable over which to compute the ADRF or AMEF.
- `subset` – a logical vector, evaluated in the original dataset used to fit the model, defining a subset of units for which the ADRF or AMEF is to be computed.
- `by` – the name of one or more variables for which the ADRF or AMEF should be computed within subgroups. Can be supplied as a character vector of variable names or a one-sided formula.
- `contrast` – either "adrf" or "amef" to request the ADRF or AMEF, respectively. The default is to compute the ADRF.
- `at` – the values of the focal predictor at which to compute the ADRF or AMEF. This should be a vector of values that the focal predictor can take on. If unspecified, a vector of `n` (see below) equally spaced values from the minimum to the maximum value of the predictor will be used. This should typically be used only if quantities are desired over a subset of the values of the focal predictor.
- `n` – if `at` is unspecified, the number of points along the range of the focal predictor at which to compute the ADRF or AMEF. More yields smoother functions, but will take longer and require more memory. The default is 21.
- `outcome` – a string containing the name of the outcome of interest when a multivariate (multiple outcome) model is supplied to `sim()` or the outcome category of interest when a multinomial model is supplied to `sim()`. For univariate (single outcome) and binary outcomes, this is ignored.
- `type` – a string containing the type of predicted value to return. In most cases, this can be left unspecified to request predictions on the scale of the outcome.
- `eps` – the value by which the observed values of the variable named in `var` are changed when it is continuous to compute the AMEF. This usually does not need to be specified.
- `verbose` – whether to display a progress bar.
- `cl` – an argument that controls parallel processing, which can be the number of cores to use or a cluster object resulting from `parallel::makeCluster()`.

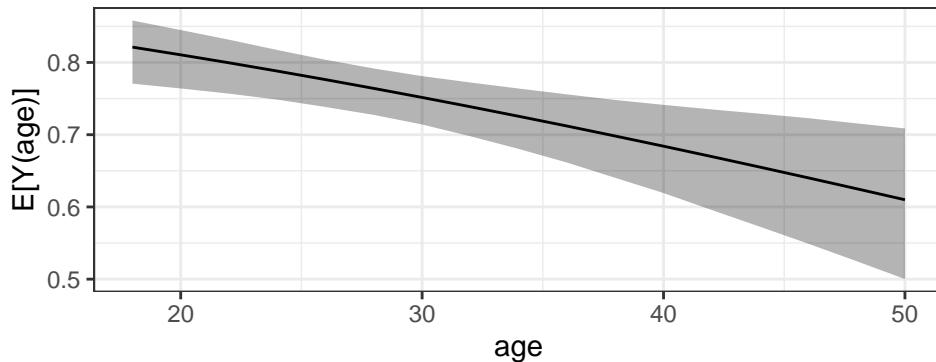
Here, we will consider `age` (just for demonstration; this analysis does not have a valid interpretation) and compute the ADRF and AMEF of `age` on the outcome. We will only examine ages between 18 and 50, even though the range of `age` goes slightly beyond these values. First, we will compute the ADRF of `age`, which examines how the outcome would vary on average if one set all units' values of `age` to each value between 18 and 50 (here we only use even ages to speed up computation).

```
age_seq <- seq(18, 50, by = 2)

est8 <- sim_adrf(s,
                  var = "age",
                  contrast = "adrf",
                  at = age_seq,
                  verbose = FALSE)
```

We can plot the ADRF using `plot()`.

```
plot(est8)
```



From the plot, we can see that as age increases, the expected outcome decreases.

We can also examine the AAPs at the requested ages using `summary()`, which will display all the estimated AAPs by default, so we will request just the first 4 (ages 18 to 24):

```
summary(est8, parm = 1:4)

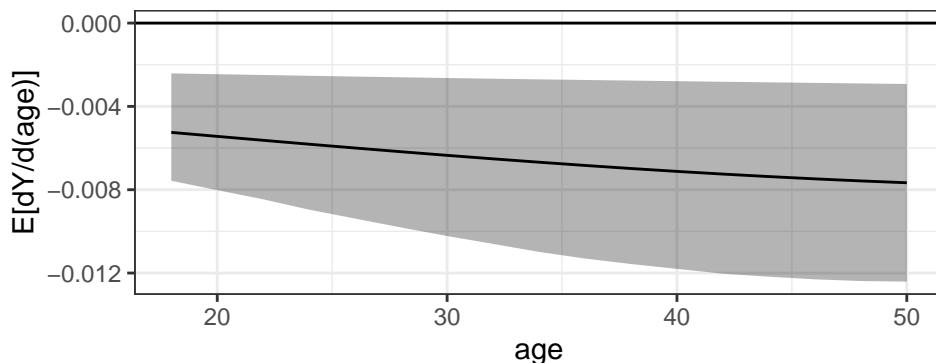
#>           Estimate 2.5 % 97.5 %
#> E[Y(18)]    0.821 0.771  0.858
#> E[Y(20)]    0.811 0.764  0.845
#> E[Y(22)]    0.800 0.757  0.832
#> E[Y(24)]    0.788 0.749  0.817
```

Next, we will compute the AMEF, the effect of age at each level of age.

```
est9 <- sim_adrf(s,
  var = "age",
  contrast = "amef",
  at = age_seq,
  verbose = FALSE)
```

We can plot the AMEF using `plot()`:

```
plot(est9)
```



From the plot, we can see the AME of age decreases slightly but is mostly constant across values of age, and the uncertainty intervals for the AMEs consistently exclude 0.

## 5 Transforming and Combining Estimates

Often, our quantities of interest are not just the outputs of the functions above, but comparisons between them. For example, to test for moderation of a treatment effect, we may want to compare AMEs in multiple groups defined by the moderator. Or, it might be that we are interested in an effect described using a different effect measure than the one originally produced; for example, we

may decide we want the risk difference AME after computing the risk ratio AME. The functions `transform()` and `cbind()` allow users to transform quantities in a single `clarify_est` object and combine two `clarify_est` objects. These are essential for computing quantities that themselves are derived from the derived quantities computed by the `sim_*`() functions.

### 5.1 `transform()`

`transform()` is a generic function in R that is typically used to create a new variable in a data frame that is a function of other columns. For example, to compute the binary outcome we used in our model, we could have run the following<sup>5</sup>:

```
lalonde <- transform(lalonde,
                      re78_0 = ifelse(re78 == 0, 1, 0))
```

Similarly, to compute a derived or transformed quantity from a `clarify_est` object, we can use `transform()`. Here, we will compute the risk difference AME of `treat`; previously, we used `sim_ame()` to compute the AAPs and the risk ratio.

```
est6 <- transform(est6,
                   RD = `E[Y(1)]` - `E[Y(0)]`)
```

Note that we used tics (`) around the names of the AAPs; this is necessary when they contain special characters like parentheses or brackets. An alternative is to use the shortcut names `.b#`, where # is replaced with a number (e.g., as `.b1`, `.b2`, etc.) corresponding to the index of the quantity referenced. For example, because `E[Y(1)]` and `E[Y(0)]` are the second and first computed quantities, respectively, the above code could be replaced with

```
est6 <- transform(est6,
                   RD = .b2 - .b1)
```

which will yield identical results<sup>6</sup>.

When we run `summary()` on the output, the new quantity, which we named “RD”, will be displayed along with the other estimates. We will also set a null value for this quantity.

```
summary(est6, null = c(`RR` = 1, `RD` = 0))

#>      Estimate   2.5 %  97.5 % P-value
#> E[Y(0)]  0.6866  0.6081  0.7596   .
#> E[Y(1)]  0.7551  0.6850  0.8088   .
#> RR        1.0998  0.9485  1.2554   0.21
#> RD        0.0685 -0.0382  0.1580   0.21
```

As mentioned previously, one benefit of using simulation-based inference with p-values computed from inverting the confidence intervals is that the p-values for testing the same hypothesis with the risk difference and risk ratio (and any other effect measure for comparing a pair of values) will always exactly align, thereby ensuring inference does not depend on the effect measure used. In contrast, Wald-type inference (based on either the simulation-derived or delta method standard error) is not invariant to transformations of the quantity of interest.

The same value would be computed if we were to have called `sim_ame()` on the same `clarify_sim` object and requested the risk difference using `contrast = "diff"`; using `transform()` saves time because the AAPs are already computed and stored in the `clarify_est` object.

We can use `transform()` along with the `by` variable in `sim_ame()` to compute the contrast between quantities computed within each subgroup of `married`. Previously we used `by` to compute the risk ratio ATT within levels of `married`; here we will compute the ratio of these risk ratios to assess the presence of effect modification.

<sup>5</sup>Users familiar with the `tidyverse` will note the similarities between `transform()` and `dplyr::mutate()`; only `transform()` can be used with `clarify_est` objects.

<sup>6</sup>Note that if a quantity is named `.b#`, e.g., `.b1`, it can only be referred to using the positional shortcut and not its name. That is, the positional shortcut takes precedence over the names of the quantities.

```

est6b |>
  transform(`RR[1]/RR[0]` = `RR[1]` / `RR[0]`) |>
  summary(parm = c("RR[0]", "RR[1]", "RR[1]/RR[0]"),
          null = 1)

#>           Estimate 2.5 % 97.5 % P-value
#> RR[0]      1.061  0.909  1.234   0.434
#> RR[1]      1.270  0.948  1.583   0.094 .
#> RR[1]/RR[0] 1.196  0.908  1.516   0.174
#> ---
#> Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

`RR[1]/RR[0]` contains the ratio of the risk ratios for `married = 1` and `married = 0`. Here we also include a test for whether each of the risk ratios and their ratio differs from 1, which is equivalent to testing whether the risk ratios differ across levels of `married`.

## 5.2 `cbind()`

`cbind()` is another generic R function that is typically used to combine two or more datasets columnwise (i.e., to widen a dataset). In `clarify`, `cbind()` can be used to combine two `clarify_est` objects so that the estimates can be examined jointly and so that it is possible to compare them directly. For example, if we were to compute AMEs in two subgroups using `subset` and wanted to compare them, we would call `sim_ame()` twice, one for each subset (though in practice it is more effective to use `by`; this is just for illustration), as demonstrated below:

```

# AME of treat with race = "black"
est10b <- sim_ame(s, var = "treat", subset = race == "black",
                    contrast = "diff", verbose = FALSE)
summary(est10b)

#>           Estimate 2.5 % 97.5 %
#> E[Y(0)]  0.6677  0.5813  0.7529
#> E[Y(1)]  0.7439  0.6661  0.8016
#> Diff      0.0762 -0.0359  0.1700

# AME of treat with race = "hispan"
est10h <- sim_ame(s, var = "treat", subset = race == "hispan",
                    contrast = "diff", verbose = FALSE)
summary(est10h)

#>           Estimate 2.5 % 97.5 %
#> E[Y(0)]  0.8266  0.7146  0.8990
#> E[Y(1)]  0.8971  0.7888  0.9527
#> Diff      0.0704 -0.0223  0.1387

```

Here, we computed the risk difference for the subgroups `race = "black"` and `race = "hispan"`. If we wanted to compare the risk differences, we could combine them and compute a new quantity equal to their difference. We will do that below.

First, we need to rename the quantities in each object so they do not overlap; we can do so using `names()`, which has a special method for `clarify_est` objects.

```

names(est10b) <- paste(names(est10b), "b", sep = "_")
names(est10h) <- paste(names(est10h), "h", sep = "_")

```

Next, we use `cbind()` to bind the objects together.

```

est10 <- cbind(est10b, est10h)
summary(est10)

#>           Estimate 2.5 % 97.5 %
#> E[Y(0)]_b  0.6677  0.5813  0.7529
#> E[Y(1)]_b  0.7439  0.6661  0.8016

```

```
#> Diff_b      0.0762 -0.0359  0.1700
#> E[Y(0)]_h   0.8266  0.7146  0.8990
#> E[Y(1)]_h   0.8971  0.7888  0.9527
#> Diff_h      0.0704 -0.0223  0.1387
```

Finally, we can use `transform()` to compute the difference between the risk differences:

```
est10 <- transform(est10,
  `Dh - Db` = Diff_h - Diff_b)
summary(est10, parm = "Dh - Db")

#>           Estimate    2.5 %   97.5 %
#> Dh - Db -0.00575 -0.06833  0.04103
```

Importantly, `cbind()` can only be used to join together `clarify_est` objects computed using the same simulated coefficients (i.e., resulting from the same call to `sim()`). This preserves the covariance among the estimated quantities, which is critical for maintaining valid inference. That is, `sim()` should only be called once per model, and all derived quantities should be computed using its output.

## 6 Using `clarify` with Multiply Imputed Data

Multiple imputation is a popular method of estimating quantities of interest in the presence of missing data and involves creating multiple versions of the original dataset, each with the missing values imputed with estimates from an imputation model. Simulation-based inference in multiply imputed data is relatively straightforward. Simulated coefficients are drawn from the model estimated in each imputed dataset separately, and then the simulated coefficients are pooled into a single set of simulated coefficients. In Bayesian terms, this would be considered “mixing draws” and is the recommended approach for Bayesian analysis with multiply imputed data (Zhou and Reiter, 2010).

Using `clarify` with multiply imputed data is simple. Rather than using `sim()`, we use the function `misim()`. `misim()` functions just like `sim()` except that it takes in a list of model fits (i.e., containing a model fit to each imputed dataset) or an object containing such a list (e.g., a `mira` object from `mice::with()` or a `mimira` object from `MatchThem::with()`). `misim()` simulates coefficient distributions within each imputed dataset and then appends them together to form a single combined set of coefficient draws.

`sim_apply()` and its wrappers accept the output of `misim()` and compute the desired quantity using each set of coefficients. When these functions rely on using a dataset (e.g., `sim_ame()`, which averages predicted outcomes across all units in the dataset used to fit the model), they automatically know to associate a given coefficient draw with the imputed dataset that was used to fit the model that produced that draw. In user-written functions supplied to the `FUN` argument of `sim_apply()`, it is important to correctly extract the dataset from the model fit. This is demonstrated below.

The final estimates of the quantity of interest are computed as the mean of the estimates computed in each imputed dataset (i.e., using the original coefficients, not the simulated ones), which is the same quantity that would be computed using standard pooling rules. This is not always valid for non-collapsible estimates, like ratios, and so care should be taken to ensure the mean of the resulting estimates has a valid interpretation (this is related to the transformation-induced bias described by Rainey (2017)).

The arguments to `misim()` are as follows:

```
misim(fitlist = , n = , vcov = , coefs = , dist = )
```

- `fitlist` – a list of model fits or an accepted object containing them (e.g., a `mira` object from `mice::with()`)
- `n` – the number of simulations to run for each imputed dataset. The default is 1000, but fewer can be used because the total number of simulated quantities will be `m * n`, where `m` is the number of imputed datasets.
- `vcov, coefs, dist` – the same as with `sim()`, except that a list of such arguments can be supplied to be applied to each imputed dataset.

Below we illustrate using `misim()` and `sim_apply()` with multiply imputed data. We will use the `africa` dataset from the `Amelia` package.

```

library(Amelia)
data("africa", package = "Amelia")

# Multiple imputation
a.out <- amelia(x = africa, m = 10, cs = "country",
                  ts = "year", logs = "gdp_pc", p2s = 0)

# Fit model to each dataset
model.list <- with(a.out, lm(gdp_pc ~ infl * trade))

# Simulate coefficients, 100 draws per imputation
si <- misim(model.list, n = 100)

si

#> A `clarify_misim` object
#> - 4 coefficients, 10 imputations with 100 simulated values each
#> - sampled distributions: multivariate t(116)

```

The function we will be applying to each imputed dataset will be one that computes the AME of `infl`. (We will run the same analysis afterward using `sim_ame()`.)

```

sim_fun <- function(fit) {
  # Extract the original dataset using get_predictors()
  X <- insight::get_predictors(fit)

  p0 <- predict(fit)

  # Predictions after perturbing infl slightly
  p1 <- predict(fit, newdata = transform(X, infl = infl + 1e-5))

  c(AME = mean((p1 - p0) / 1e-5))
}

est_mi <- sim_apply(si, FUN = sim_fun, verbose = FALSE)

summary(est_mi)

#>     Estimate 2.5 % 97.5 %
#> AME      -5.75 -8.82  -2.26

```

Note that `sim_apply()` “knows” which imputation produced each set of simulated coefficients, so using `insight::get_predictors()` on the `fit` supplied to `sim_fun()` will use the right dataset. Care should be taken when analyses restrict each imputed dataset in a different way (e.g., when matching with a caliper in each one), as the resulting imputations may not refer to a specific target population and mixing the draws may be invalid.

Below, we can use `sim_ame()`:

```

est_mi2 <- sim_ame(si, var = "infl", verbose = FALSE)

summary(est_mi2)

#>             Estimate 2.5 % 97.5 %
#> E[dY/d(infl)]    -5.75 -8.82  -2.26

```

We get the same results, as expected.

Note that `misim()` is compatible with model fit objects from `mice`, `Amelia`, `MatchThem`, and any other package that produces a list of model fit objects with each corresponding to the output of a model fit to an imputed dataset.

## 7 Comparison to Other Packages

Several packages offer methods for computing interpretable quantities from regression models, including `emmeans` (Lenth, 2024), `margins` (Leeper, 2021), `modelbased` (Makowski et al., 2020), and

`marginaleffects` (Arel-Bundock et al., Forthcoming). Many of the quantities computed by these packages can also be computed by `clarify`, the primary difference being that `clarify` uses simulation-based inference rather than delta method-based inference.

`marginaleffects` offers the most similar functionality to `clarify`, and `clarify` depends on functionality provided by `marginaleffects` to accommodate a wide variety of regression models. `marginaleffects` also offers simulation-based inference using `marginaleffects::inferences()` and support for arbitrary user-specified post-estimation functions using `marginaleffects::hypotheses()`. However, `clarify` and `marginaleffects` differ in several ways. The largest difference is that `clarify` supports iterative building of more and more complex hypotheses through the `transform()` method, which quickly computes new quantities and transformations from the existing computed quantities, whereas `marginaleffects` only supports a single transformation and, as of version 0.20.0, cannot use simulation-based inference for these quantities.

Because of `clarify`'s focus on simulation, it provides functionality directly aimed at improving simulation-based inference, including plots to view the distributions of simulated values and support for parallel processing. `clarify` also provides support for simulation-based inference of multiply imputed data, which does not require any special pooling rules.

There are areas and cases where `marginaleffects` may be the better choice than `clarify` or where the differences between the packages are of little consequence. `marginaleffects` focuses on providing a complete framework for post-estimation using model predictions, whereas `clarify` is primarily focused on supporting user-defined functions, with commonly used estimators offered as a convenience. In cases where the delta method is an acceptable approximation (e.g., for quantities computed from linear models or other quantities known to be approximately Normally distributed in finite samples), using the delta method through `marginaleffects` will be much faster, more accurate, and more replicable than the simulation-based inference `clarify` provides. For the quantities easily computed by `marginaleffects` that support simulation-based inference through `marginaleffects::inferences()`, using `marginaleffects` can provide a more familiar and flexible syntax than `clarify` might offer. Ultimately, the user should use the package that supports their desired syntax and mode of inference.

## 8 Conclusion

`clarify` provides functionality to facilitate simulation-based inference for deriving quantities from regression models. This framework provides an alternative to the delta method that can yield confidence intervals with coverage closer to nominal for some quantities of interest. While we do not claim that simulation-based inference should be universally preferred over delta method-based inference, there are cases in which it can retain some advantageous properties, and we hope the availability of these methods in `clarify` encourages additional research on when those properties can be realized and facilitates empirical work that takes advantage of these properties.

## References

- V. Arel-Bundock, N. Greifer, and A. Heiss. How to interpret statistical models using `marginaleffects` in R and Python. *Journal of Statistical Software*, Forthcoming. [p156, 172]
- R. H. Dehejia and S. Wahba. Causal effects in nonexperimental studies: Reevaluating the evaluation of training programs. *Journal of the American Statistical Association*, 94(448):1053–1062, 12 1999. doi: 10.1080/01621459.1999.10473858. [p156]
- B. Efron and R. Tibshirani. Bootstrap methods for standard errors, confidence intervals, and other measures of statistical accuracy. *Statistical Science*, 1(1):54–75, 1986. ISSN 0883-4237. URL <https://www.jstor.org/stable/2245500>. [p154, 155]
- N. Greifer and E. A. Stuart. Choosing the causal estimand for propensity score analysis of observational studies. 2023. doi: 10.48550/arXiv.2106.10577. [p164]
- M. C. Herron. Postestimation uncertainty in limited dependent variable models. *Political Analysis*, 8(1):83–98, 1999. doi: 10.1093/oxfordjournals.pan.a029806. URL [https://www.cambridge.org/core/product/identifier/S104719870008524/type/journal\\_article](https://www.cambridge.org/core/product/identifier/S104719870008524/type/journal_article). [p154]
- D. E. Ho, K. Imai, G. King, and E. A. Stuart. Matchit: Nonparametric preprocessing for parametric causal inference. *Journal of Statistical Software, Articles*, 42(8):1–28, 2011. doi: 10.1863/jss.v042.i08. Citation Key: JSSv042i08 tex.ids= hoMatchItNonparametricPreprocessing2011. [p156]

- A. R. Hole. A comparison of approaches to estimating confidence intervals for willingness to pay measures. *Health Economics*, 16(8):827–840, 08 2007. doi: 10.1002/hec.1197. URL <https://onlinelibrary.wiley.com/doi/10.1002/hec.1197>. [p154]
- K. Imai, G. King, and O. Lau. Toward a common framework for statistical analysis and development. *Journal of Computational and Graphical Statistics*, 17(4):892–913, 12 2008. doi: 10.1198/106186008X384898. [p156]
- G. King, M. Tomz, and J. Wittenberg. Making the most of statistical analyses: Improving interpretation and presentation. *American Journal of Political Science*, 44(2):347–361, 2000. doi: 10.2307/2669316. tex.ids= kingMakingMostStatistical2000a publisher: [Midwest Political Science Association, Wiley]. [p154]
- I. Krinsky and A. L. Robb. On approximating the statistical properties of elasticities. *The Review of Economics and Statistics*, 68(4):715, 11 1986. doi: 10.2307/1924536. URL <http://dx.doi.org/10.2307/1924536>. [p154]
- T. J. Leeper. *margins: Marginal Effects for Model Objects*, 2021. R package version 0.3.26. [p171]
- R. V. Lenth. *emmeans: Estimated Marginal Means, aka Least-Squares Means*, 2024. URL <https://CRAN.R-project.org/package=emmeans>. R package version 1.10.2. [p171]
- J. S. Long and J. Freese. *Regression models for categorical dependent variables using Stata*. Stata Press Publication, StataCorp LP, College Station, Texas, third edition edition, 2014. OCLC: ocn890178695. [p164]
- D. P. MacKinnon, C. M. Lockwood, and J. Williams. Confidence limits for the indirect effect: Distribution of the product and resampling methods. *Multivariate Behavioral Research*, 39(1):99–128, 01 2004. doi: 10.1207/s15327906mbr3901\_4. URL [http://www.tandfonline.com/doi/abs/10.1207/s15327906mbr3901\\_4](http://www.tandfonline.com/doi/abs/10.1207/s15327906mbr3901_4). [p154]
- D. Makowski, M. S. Ben-Shachar, I. Patil, and D. Lüdecke. Estimation of model-based predictions, contrasts and means. CRAN, 2020. URL <https://github.com/easystats/modelbased>. [p171]
- K. J. Preacher and J. P. Selig. Advantages of monte carlo confidence intervals for indirect effects. *Communication Methods and Measures*, 6(2):77–98, Apr 2012. ISSN 1931-2458, 1931-2466. doi: 10.1080/19312458.2012.679848. URL <http://www.tandfonline.com/doi/abs/10.1080/19312458.2012.679848>. [p154]
- C. Rainey. Transformation-induced bias: Unbiased coefficients do not imply unbiased quantities of interest. *Political Analysis*, 25(3):402–409, 07 2017. doi: 10.1017/pan.2017.11. [p170]
- C. Rainey. A careful consideration of CLARIFY: Simulation-induced bias in point estimates of quantities of interest. *Political Science Research and Methods*, pages 1–10, Apr. 2023. ISSN 2049-8470, 2049-8489. doi: 10.1017/psrm.2023.8. [p156, 159]
- M. Thulin. *Modern statistics with R: from wrangling and exploring data to inference and predictive modelling*. Eos Chasma Press, Uppsala, Sweden, 2021. [p155]
- D. Tofghi and D. P. MacKinnon. Monte carlo confidence intervals for complex functions of indirect effects. *Structural Equation Modeling: A Multidisciplinary Journal*, 23(2):194–205, 03 2016. doi: 10.1080/10705511.2015.1057284. URL <https://doi.org/10.1080/10705511.2015.1057284>. Publisher: Routledge \_eprint: <https://doi.org/10.1080/10705511.2015.1057284>. [p154, 155]
- M. Tomz, J. Wittenberg, and G. King. Clarify: Software for interpreting and presenting statistical results. *Journal of Statistical Software*, 8:1–30, 01 2003. doi: 10.18637/jss.v008.i01. [p156]
- B. A. Zelner. Using simulation to interpret results from logit, probit, and other nonlinear models. *Strategic Management Journal*, 30(12):1335–1348, 2009. doi: 10.1002/smj.783. URL <https://onlinelibrary.wiley.com/doi/abs/10.1002/smj.783>. tex.ids= zelnerUsingSimulationInterpret2009a \_eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1002/smj.783>. [p154]
- X. Zhou and J. P. Reiter. A note on bayesian inference after multiple imputation. *The American Statistician*, 64(2):159–163, 05 2010. doi: 10.1198/tast.2010.09109. [p170]

Noah Greifer  
Harvard University  
Institute for Quantitative Social Science

*Cambridge, MA*  
*ngreifer.github.io*  
*ORCiD: 0000-0003-3067-7154*  
*ngreifer@iq.harvard.edu*

*Steven Worthington*  
*Harvard University*  
*Institute for Quantitative Social Science*  
*Cambridge, MA*  
*ORCiD: 0000-0001-9550-5797*  
*sworthington@iq.harvard.edu*

*Stefano Iacus*  
*Harvard University*  
*Institute for Quantitative Social Science*  
*Cambridge, MA*  
*ORCiD: 0000-0002-4884-0047*  
*siacus@iq.harvard.edu*

*Gary King*  
*Harvard University*  
*Institute for Quantitative Social Science*  
*Cambridge, MA*  
*ORCiD: 0000-0002-5327-7631*  
*king@harvard.edu*

# FuzzySimRes: Epistemic Bootstrap – an Efficient Tool for Statistical Inference Based on Imprecise Data

by Maciej Romanik, Przemysław Grzegorzewski, Abbas Parchami

**Abstract** The classical Efron's bootstrap is widely used in many areas of statistical inference, including imprecise data. In our new package FuzzySimRes we adapted the bootstrap methodology to the epistemic fuzzy data, i.e. fuzzy perceptions of the usual real-valued random variables. The epistemic bootstrap algorithms deliver real-valued samples generated randomly from the initial fuzzy sample. Then, these samples can be utilized directly in various statistical procedures. Moreover, we implemented a practically oriented simulation procedure to generate synthetic fuzzy samples and provided a real-life epistemic dataset ready to use for various techniques of statistical analysis. Some examples of their applications, together with the comparisons of the epistemic bootstrap algorithms and the respective benchmarks, are also discussed.

## 1 Introduction

Efron's bootstrap (Efron and Tibshirani, 1993) is a simple but very powerful tool. This useful resampling method is successfully applied in statistical inference, including estimation, hypotheses testing, and other data analysis techniques, e.g., Davison and Hinkley (1997); James et al. (2021); Romanik (2019).

In our package **FuzzySimRes** we adapted the classical bootstrap algorithm to a special kind of imprecise data, i.e. the epistemic random fuzzy numbers (see Couso and Dubois (2014)), which might be treated as fuzzy perceptions of the usual real-valued random variables. This way, a special resampling methodology, known as the epistemic bootstrap, can be introduced (Grzegorzewski and Romanik, 2021, 2022a,c). Following the suggested methods we can generate random real-valued samples based on the initial fuzzy sample. Such a “change of a viewpoint” from the “fuzzy world” to its “clear” (i.e. real-value) counterpart can be a very useful and important tool. This allows all commonly used classical statistical methods (developed for real-valued samples), including statistical tests, estimation procedures, etc., to be directly and easily adapted to fuzzy epistemic samples.

Please note that statistical inference for the fuzzy data is usually underdeveloped, poses some problems, and leads to discussions about the intuitions, solutions, etc. (e.g., concerning the different approaches to the p-value).

We provide some useful functions in our package **FuzzySimRes**. They are related to a practically oriented simulation of various types of fuzzy numbers (FNs), the epistemic bootstrap itself, and its applications related to the estimation of important statistical measures of the initial sample, and the one- and two-sample statistical tests. Additionally, we provide the real-life dataset of the epistemic FNs, which can be useful in comparing various approaches to fuzzy statistical inference. Based on the two general epistemic bootstrap functions, users of **FuzzySimRes** can build their own “epistemic bootstrap statistical tools” to fit their purposes (e.g., the necessity of using tests other than the Kolmogorov-Smirnov one).

In the following, we briefly compare **FuzzySimRes** package with other existing ones and introduce a necessary notation. Then, the functions implemented in the package are illustrated with the respective examples. Finally, the outcomes of these functions are compared taking into account some benchmarks for different statistical problems using both the synthetic and real-life data.

### 1.1 A brief review of related packages

There are some packages related to fuzzy numbers and their statistical analysis. Firstly, we should mention **FuzzyNumbers** (Gagolewski and Caha, 2021). This library aims to provide S4 classes and methods for FNs. They can be used to construct different types of FNs (e.g. triangular or trapezoidal ones), compute arithmetic operators for fuzzy values, calculate their approximations, and find different characteristics of FNs (like the possibility and necessity values, expected interval, ambiguity, membership functions among many others) for arbitrary FNs or some of their special types, etc. Notably, our package **FuzzySimRes** uses S4 objects describing FNs derived from this package. However, there are no special functions devoted to simulations or resampling methods in **FuzzyNumbers** package. Therefore, it can be seen as a kind of “foundation” to deal with FNs.

The next package, **FuzzySTs** (Berkachy and Donzé, 2020) is a collection of various statistical tools, like fuzzification methods, numerical estimations of fuzzy statistical measures and bootstrap distribution of the likelihood ratio, testing hypotheses by fuzzy confidence intervals and estimation of the fuzzy p-values for epistemic fuzzy data. These approaches are related to fuzzy notions, like fuzzy p-values, resulting in the strictly fuzzy output (Berkachy and Donzé, 2019).

And **SAFD** (Trutschig et al., 2013) package joins two kinds of functions. Similarly to **FuzzyNumbers**, it provides basic operations on FNs (like the sum, mean, etc.), but it also contains some strictly statistical functions. They allow us to simulate FNs and perform bootstrap tests for the equality of the means. As for the simulation function, this is an implementation of the second procedure described by González-Rodríguez et al. (2009), where a respective basis is perturbated stochastically to generate a new polygonal fuzzy number. There are important theoretical and practical differences (Parchami et al., 2024) between this approach and the one applied in **FuzzySimRes** package. The statistical tests in **SAFD** library are exclusively based on the classical bootstrap as described by Colubi (2009); Montenegro et al. (2004).

On the other hand, **Sim.PLFN** (Parchami, 2017) can be seen as a kind of “ancestor” of our package. It allows only to simulation of some kinds of FNs some kinds of FNs, especially so-called piecewise linear FNs (Coroianu et al., 2013), and calculates a few basic operators (like the sum, mean, and variance) for them.

Another package, **FuzzyStatTra** (Lubiano and de la Rosa de Saa, 2017), also provides basic statistical functions for FNs like calculation of the mean and medians, indexes, and various distance measures. Some simulation procedures are also included there, but they are intended for special cases of dependent and independent components described by Sinova et al. (2016). Therefore, they can not be considered as “multi-purpose” generation functions when the probability distributions are selected by the user.

We should also mention our previous package, **FuzzyResampling** (Romaniuk et al., 2022) that provides various resampling algorithms other than the classical bootstrap (Romaniuk and Grzegorzewski, 2023). The main aim of these approaches is to overcome a problem with repetition of a few distinct values (which is commonly seen in the case of the Efron’s bootstrap) and to create FNs, which are “similar” (in the sense of some characteristics of FNs) but not “the same” as values from the initial sample (Grzegorzewski et al., 2020; Romaniuk and Hryniwicz, 2019; Grzegorzewski and Romaniuk, 2022b). Additionally, the tests for the means related to the approach presented by Lubiano et al. (2016) but based on these new resampling methods are also provided.

Nevertheless, **FuzzySimRes** has some unique features. Firstly, it adds very useful simulation functions (as its acronym – Fuzzy Simulations and Resampling – suggests) to complete **FuzzyNumbers** in this field. These procedures are very intuitive and practically oriented as noted by Parchami et al. (2024) (contrary to, e.g., **SAFD** that adds some random noise without keeping track of important characteristics of the input FNs). Secondly, the so-called epistemic bootstrap is implemented there. Apart from ready-to-use general epistemic bootstrap functions, special procedures are provided for the estimation of parameters of the underlying statistical model, together with an interface that can be used with various classical statistical tests. The epistemic bootstrap is a relatively new idea and the respective algorithms were not implemented in other publicly available software packages (including R itself). It should be noted, that this approach is completely different when compared with the ontic-oriented resampling procedures from **FuzzyResampling** that can be seen as a “generalization” of the classical bootstrap procedure (Grzegorzewski et al., 2020; Grzegorzewski and Romaniuk, 2022b).

## 1.2 Epistemic fuzzy numbers

In the following, we recall some basic concepts and notations concerning fuzzy numbers. For a more detailed introduction, we refer the reader to, e.g., Ban et al. (2015).

A **fuzzy number** (abbreviated further as FN) is an imprecise value characterized by a mapping  $\tilde{A} : \mathbb{R} \rightarrow [0, 1]$  (a **membership function**), such that its  $\alpha$ -cut defined by

$$\tilde{A}_\alpha = \begin{cases} \{x \in \mathbb{R} : \tilde{A}(x) \geq \alpha\} & \text{if } \alpha \in (0, 1], \\ cl\{x \in \mathbb{R} : \tilde{A}(x) > 0\} & \text{if } \alpha = 0, \end{cases} \quad (1)$$

is a nonempty compact interval for each  $\alpha \in [0, 1]$ . Operator  $cl$  in (1) denotes the closure. Every FN is completely characterized both by its membership function  $\tilde{A}(x)$  and a family of  $\alpha$ -cuts  $\{\tilde{A}_\alpha\}_{\alpha \in [0,1]}$ . There are two special  $\alpha$ -cuts: the **core**  $\tilde{A}_1 = \text{core}(\tilde{A})$ , which contains all values fully compatible with the concept described by  $\tilde{A}$ , and the **support**  $\tilde{A}_0 = \text{supp}(\tilde{A})$  on real line, for which values are compatible to some extent with the concept modeled by  $\tilde{A}$ . A family of all FNs will be denoted by  $\mathbb{F}(\mathbb{R})$ .

There are many possible shapes of the membership functions. A special family of the **LR-fuzzy numbers** is defined by

$$\tilde{A}(x) = \begin{cases} 0 & \text{if } x < a_1, \\ L\left(\frac{x-a_1}{a_2-a_1}\right) & \text{if } a_1 \leq x < a_2, \\ 1 & \text{if } a_2 \leq x < a_3, \\ R\left(\frac{a_4-x}{a_4-a_3}\right) & \text{if } a_3 \leq x < a_4, \\ 0 & \text{if } x \geq a_4, \end{cases} \quad (2)$$

where  $L, R : [0, 1] \rightarrow [0, 1]$  are continuous and strictly increasing function such that  $L(0) = R(0) = 0$  and  $L(1) = R(1) = 1$ , and  $a_1, a_2, a_3, a_4 \in \mathbb{R}$ , where  $a_1 \leq a_2 \leq a_3 \leq a_4$ . If  $L$  and  $R$  are linear functions, i.e.  $L\left(\frac{x-a_1}{a_2-a_1}\right) = \frac{x-a_1}{a_2-a_1}$  and  $R\left(\frac{a_4-x}{a_4-a_3}\right) = \frac{a_4-x}{a_4-a_3}$ , we get a **trapezoidal fuzzy number** (denoted further on as TPFN). Moreover, if  $a_2 = a_3$  then we have a **triangular fuzzy number** (abbreviated as TRFN). In these two cases, we can simply write  $A = (a_1, a_2, a_3, a_4)$  (for TPFN) or  $A = (a_1, a_2, a_4)$  (for TRFN) to fully describe such FNs.

Another type of the LR-fuzzy number is known as the  **$k$ -knot piecewise linear fuzzy number** (Coroianu et al., 2019) (or polygonal fuzzy number, see Bález-Sánchez et al. (2012), abbreviated further as PLFN), which is suitable especially in an approximation of more complex FNs. In this case,  $L$  and  $R$  functions are polygons consisting of  $k \in \mathbb{N}$  segments.

Fuzzy numbers are used to model the results of various experiments that cannot be precisely described, qualified, or measured. But in many cases, we have to draw conclusions and make decisions based on data whose uncertainty comes both from randomness (which classical statistics copes with) and lack of precision (for which the fuzzy set theory is perfect for modeling). To model such data one can use **fuzzy random variables**, known also as **random fuzzy numbers** (Parchami et al., 2024).

It should be noted here that we can look at fuzzy random variables from two different perspectives: **ontic** or **epistemic** (e.g. Couso and Dubois (2014)). The first concerns data that appear to be essentially fuzzy in value, while the second refers to situations where, although precise (accurate) data values exist, they are imprecisely observed (e.g. due to imperfections in measuring devices, inaccuracies caused by people performing the measurements, or how results are reported), so their true actual values remain unknown. This second kind of imprecision is widespread in real-life problems met in engineering, science, and other applications, so further on we limit our attention only to epistemic data.

Following the definition by Kwakernaak (1978) and Kruse (1982) a fuzzy random variable  $\tilde{X}$  can be considered as a *fuzzy perception* of the unknown random variable  $X$ , called the *original* of  $\tilde{X}$ . More precisely, given a probability space  $(\Omega, \mathcal{F}, P)$ , a mapping  $X : \Omega \rightarrow \mathbb{F}(\mathbb{R})$  is said to be a fuzzy random variable (f.r.v.) if for each  $\alpha \in [0, 1]$   $(\inf X_\alpha) : \Omega \rightarrow \mathbb{R}$  and  $(\sup X_\alpha) : \Omega \rightarrow \mathbb{R}$  are real-valued random variables on  $(\Omega, \mathcal{F}, P)$ . Similarly, a **fuzzy random sample**  $\tilde{X}_1, \dots, \tilde{X}_n$  is a fuzzy perception of a random sample  $X_1, \dots, X_n$  of the usual real-valued random variables. For more details, we refer the reader to Kwakernaak (1978); Kruse (1982).

### 1.3 Epistemic vs classical bootstrap

There are important differences between the classical Efron's bootstrap (Efron and Tibshirani, 1993) and its epistemic counterpart (Grzegorzewski and Romaniuk, 2021, 2022c, 2024). In the classical bootstrap approach, the initial sample is then directly resampled. Therefore, in the case of fuzzy input, the output also consists of the same FNs as in the primary sample (with possible repetitions or omitting some of them). This procedure can be very useful in statistical inference (see, e.g., (Gil et al., 2006; Lubiano et al., 2016; Montenegro et al., 2004)) but the respective statistical tests (or other statistical procedures like the estimation) have to be specially developed for fuzzy-valued samples. Therefore, there is a need to construct "almost completely new" statistical solutions taking into account various distance measures for fuzzy sets existing in the literature, more complex definitions of the expected value, possible problems with difference operator, etc. (Ban et al., 2015; Heilpern, 1992). Resampling procedures existing in **FuzzyResampling** package can be seen as a kind of generalization of this classical bootstrap (in the same manner as the smoothed bootstrap in the case of real-valued samples). They aim to preserve some important characteristics of FNs (like the value, ambiguity, etc.) but with an alternation of FNs from the initial sample into "new" values occurring in the generated samples (Grzegorzewski et al., 2020; Grzegorzewski and Romaniuk, 2022b). However, we are still obtaining fuzzy-valued outputs for these methods.

On the other hand, in the epistemic bootstrap, a completely real-valued (i.e. "crisp") sample is generated from a fuzzy-valued initial sample. It allows to use of directly highly developed statistical tools for real-valued data (various statistical tests, point or interval estimators, etc.) without the

need for transforming them into a “new fuzzy world”. Consequently, knowing statistical tools with suitable good properties, the areas of possible applications of epistemic fuzzy data may substantially expand. To explain it better, consider the following goodness-of-fit testing problem. In Lubiano et al. (2016) and Lubiano et al. (2017), the outcomes of the well-known questionnaire TIMSS-PIRLS 2011 performed by Spanish primary school pupils were considered, while in Ramos-Guajardo et al. (2019) experts’ perceptions about different characteristics of the Gamoneda blue cheeses were discussed. In both cases, researchers dealt with subjective valuations expressed in natural language, which are inherently imprecise, and therefore modeled using ontic fuzzy sets. Thus, the problems mentioned above required the construction of appropriate statistical tools that would enable inferences to be made based on this type of data. Meanwhile, the epistemic variants of the classical Kolmogorov-Smirnov and Cramer-von Mises tests were directly used for fuzzy data concerning the lifetimes of street light equipment (Hesamian and Taheri, 2013) and electronic circuit thickness (Faraz and Shapiro, 2010) in Grzegorzewski and Romaniuk (2024). The obtained results were consistent with predictions concerning these real-life samples, like the behavior of the probability distributions of their originals (Gibbons and Chakraborti, 2010). The example related to the electronic circuit thickness is also considered further in this paper. Some other applications can be also found in (Grzegorzewski and Romaniuk, 2022a,c,b, 2024).

Moreover, using brute computational force, we can easily improve the quality of the outputs. However, the results are quite satisfactory also for the limited number of  $\alpha$ -cuts. For instance, using even 10  $\alpha$ -cuts leads to the p-values for the epistemic versions of the goodness-of-fit tests (like the Kolmogorov-Smirnov or Cramer-von Mises tests) very close to their respective benchmarks (Grzegorzewski and Romaniuk, 2024).

## 2 Overview of FuzzySimRes package

Firstly, we briefly discuss the functions implemented in **FuzzySimRes** package. They can be roughly divided into four groups:

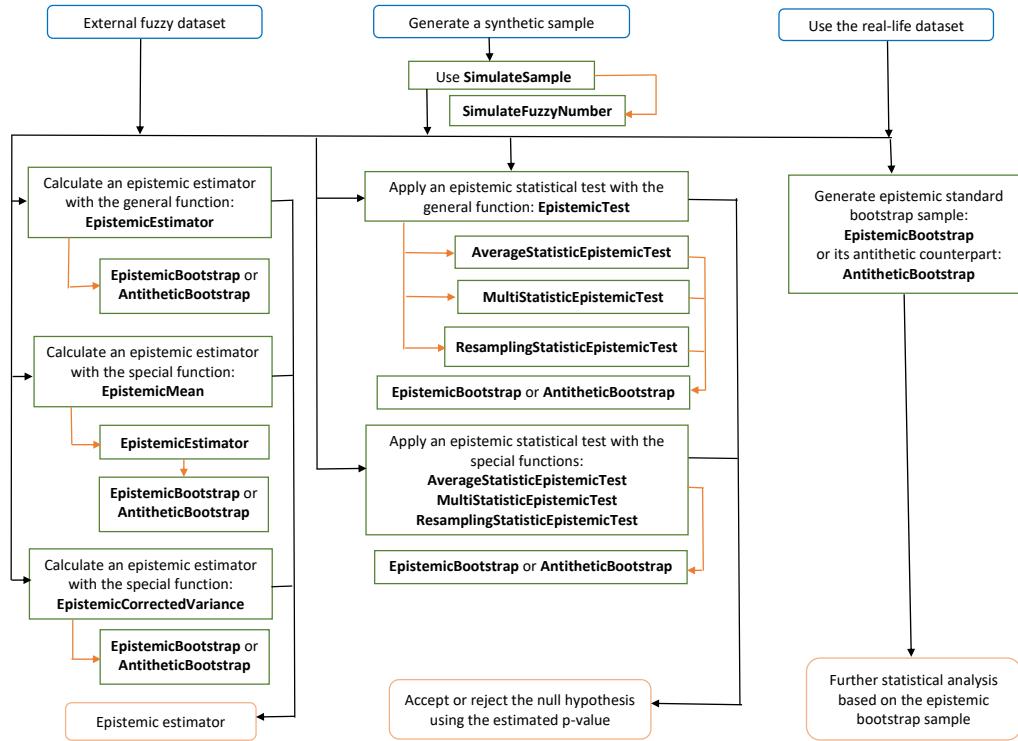
1. random generation of FNs of various types,
2. general epistemic bootstrap procedures,
3. epistemic estimation of basic population characteristics from fuzzy samples,
4. interface to statistical tests based on the epistemic bootstrap.

Moreover, a set of real-life epistemic fuzzy data is also included in the package. All examples in R can be reproduced using the supplementary file.

Taking into account the above-mentioned types of functions, there are many possible applications of **FuzzySimRes** package:

1. Generation of synthetic fuzzy samples according to the specified probability distributions. Such samples can be then used to check the validity and quality of new statistical tools for FNs in a strictly controlled “environment”, e.g., to plot power curves for a statistical test (Grzegorzewski and Romaniuk, 2022c) or to check the influence of different model parameters on the estimated p-values (Grzegorzewski and Romaniuk, 2024).
2. Estimation is one of the key problems in statistical inference. The same applies to fuzzy-valued data, especially in the epistemic case. Then, our considerations about the mean or the standard deviation related to the respective originals of the fuzzy random sample can give us the necessary insight into the parameters of the underlying statistical model (Grzegorzewski and Romaniuk, 2021, 2022c).
3. Statistical tests are the next important subject in statistical inference. To accept or reject the null hypothesis, the respective statistical test has to be developed. As it was previously mentioned, the epistemic bootstrap allows for direct application of the widely known real-valued tests instead of their “fuzzy-oriented” counterparts. Therefore, e.g., the classical goodness-of-fit Kolmogorov-Smirnov or Cramer-von Mises tests can be directly used with the interface provided by **FuzzySimRes** package (Grzegorzewski and Romaniuk, 2022a, 2024).
4. Real-life fuzzy data are also important to develop statistical procedures. Synthetic samples are very useful, but some problems are only visible when the data are provided by a “true source”. In **FuzzySimRes** package, there is a special set of such data used to construct a fuzzy statistical control chart (Faraz and Shapiro, 2010) and check the quality of statistical tests based on the epistemic bootstrap (Grzegorzewski and Romaniuk, 2024).

The general workflow for some possible applications (black lines) and the internal order of invoking functions (orange lines) from **FuzzySimRes** package can be found in Fig. 1.



**Figure 1:** General workflow for possible applications and invoking the functions from **FuzzySimRes** package.

## 2.1 Generation of the initial sample

In many cases, synthetic samples of predefined properties are necessary to analyze statistical methods numerically. Two functions in **FuzzySimRes** allow the generation of random fuzzy variables. The first one

```
SimulateFuzzyNumber(originalPD, parOriginalPD, incrCorePD,
  parIncrCorePD, suppLeftPD, parSuppLeftPD,
  suppRightPD, parSuppRightPD, knotNumbers = 0,
  type = "trapezoidal", ...)
```

is used to generate randomly a single TPFN (for type = "trapezoidal"), TRFN (type = "triangular"), or PLFN (type = "PLFN", respectively). All these types of FNs utilize the respective S4 objects from **FuzzyNumbers**.

To simulate a TPFN  $\tilde{X}$ , five independent real-valued random variables are necessary:  $X$  for its "true value" (i.e., its *original*),  $C^l, C^r$  – the left and right increment of the core,  $S^l, S^r$  – the left and right increment of the support, respectively. To generate these random variables the functions derived from **stats** (R Core Team, 2023) with the respective parameters are used (see Table 1), e.g., to draw randomly the original  $X$ , the function *originalPD* with the parameters *parOriginalPD* should be applied.

Random variable	Function	Parameters
$X$	<i>originalPD</i>	<i>parOriginalPD</i>
$C^l, C^r$	<i>incrCorePD</i>	<i>parIncrCorePD</i>
$S^l$	<i>suppLeftPD</i>	<i>parSuppLeftPD</i>
$S^r$	<i>suppRightPD</i>	<i>parSuppRightPD</i>

**Table 1:** Random variables used to simulate a TPFN.

As a result we obtain a random TPFN given by  $(X - C^l - S^l, X - C^l, X + C^r, X + C^r + S^r)$  (see also Grzegorzewski and Romaniuk (2022a) for the similar procedure). Obviously, for a TRFN we have  $C^l = C^r = 0$  without using the respective parameters in *SimulateFuzzyNumber*. In the case of a PLFN, the number of knots *knotNumbers* should be greater than zero, and then the specially truncated

probability distributions for both arms of the support are applied. The function `SimulateFuzzyNumber` returns both the generated FN (as value in the output list) and its random original  $X$  (as original).

Let us initialize a random seed and generate a TPFN with the “true origin” described by the normal distribution with the expected value  $\mu = 0$  and standard deviation  $\sigma = 1$  (denoted by  $N(\mu, \sigma)$ ), the increments of the core given by the uniform distribution on the interval  $(0, 0.6)$  (denoted by  $U(0, 0.6)$ ) and the increments of the support from  $U(0, 1)$ :

```
# seed PRNG
> set.seed(123456)

> SimulateFuzzyNumber(originalPD="rnorm",parOriginalPD=list(mean=0,sd=1),
+ incrCorePD="runif",parIncrCorePD=list(min=0,max=0.6),
+ suppLeftPD="runif",parSuppLeftPD=list(min=0,max=1),
+ suppRightPD="runif",parSuppRightPD=list(min=0,max=1),
+ type="trapezoidal")

$value
Trapezoidal fuzzy number with:
  support=[-0.316967,1.10087],
  core=[0.480817,0.902528].
```

The second function generates a sample of  $n$  independent FNs similarly to `SimulateFuzzyNumber`:

```
SimulateSample(n = 1,originalPD,parOriginalPD,incrCorePD,
  parIncrCorePD,suppLeftPD,parSuppLeftPD,
  suppRightPD,parSuppRightPD,knotNumbers = 0,
  type = "trapezoidal")
```

This function returns a list of simulated FNs together with a vector of their respective originals. Let us generate 10 TPFNs given by the same distributions as in the previous example and print the second simulated value and its “true origin”:

```
# seed PRNG
> set.seed(123456)

> sample1 <- SimulateSample(n=10,originalPD="rnorm",
+ parOriginalPD=list(mean=0,sd=1),
+ incrCorePD="runif",parIncrCorePD=list(min=0,max=0.6),
+ suppLeftPD="runif",parSuppLeftPD=list(min=0,max=1),
+ suppRightPD="runif",parSuppRightPD=list(min=0,max=1),
+ type="trapezoidal")

> sample1$original[2]

[1] -1.301602

> sample1$value[2]

$X2
Trapezoidal fuzzy number with:
  support=[-1.937,-0.229014],
  core=[-1.40214,-0.822808].
```

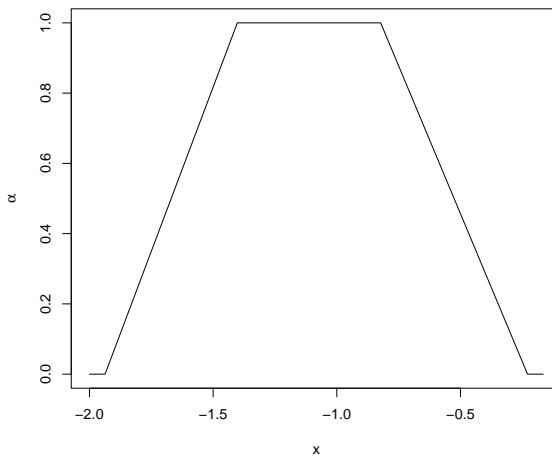
```
> plot(sample1$value[[2]])
```

The obtained graph of this exemplary FN can be found in Fig. 2.

## 2.2 Epistemic bootstrap

All of the functions described in further sections use two main procedures related to the epistemic bootstrap (Grzegorzewski and Romanuk, 2021, 2022a,c, 2024).

The first one



**Figure 2:** Example of the generated TPFN.

```
EpistemicBootstrap(fuzzySample, cutsNumber = 1,...)
```

applies the **standard epistemic bootstrap** (abbreviated as *std*) to a single value or a whole list of FNs given by *fuzzySample*. This procedure firstly generates uniformly a list of  $\alpha$ -cuts (their number is specified by *cutsNumber*). Then, it generates a sample from each of the input FNs, corresponding to the aforementioned list of the  $\alpha$ -cuts. A final output is given as a real-valued matrix, with the number of rows that is equal to *cutsNumber*, and the number of columns designated by the initial sample size. This way, we obtain  $b$  real-valued bootstrap samples  $\mathbb{X}^{*j} = (X_1^{*j}, \dots, X_n^{*j})$ , based on the initial fuzzy sample  $\tilde{\mathbb{X}} = (\tilde{X}_1, \dots, \tilde{X}_n)$ , where  $j = 1, \dots, b$  and  $b$  is equal to *cutsNumber*.

Let us apply the epistemic bootstrap with 3  $\alpha$ -cuts for the previously generated *sample1*, and then show the output rounded to 4 decimal places:

```
> set.seed(123456)

> epistemicOutput <- EpistemicBootstrap(sample1$value, cutsNumber = 3)

> round(epistemicOutput,digits = 4)

      X1      X2      X3      X4      X5      X6      X7      X8      X9      X10
0.7978 0.5323 -1.0784 0.6744 0.8553 0.9501 1.1755 0.6460 -0.5347 -0.0049 1.5937
0.7536 0.5253 -1.4512 1.4851 0.8546 0.9337 1.4773 0.6576 -0.3149 -0.0022 1.8909
0.3913 0.1991 -0.4767 1.1215 0.9443 -0.0089 0.8814 0.4538 0.0633 -0.9345 1.1242
```

The first column shows  $\alpha$ -cuts drawn randomly, while the rest columns contain values generated from each  $\alpha$ -cut.

The second function

```
AntitheticBootstrap(fuzzySample, cutsNumber = 1,...)
```

applies the so-called **antithetic epistemic bootstrap** (denoted further on by *anti*). Instead of drawing a single value from the given  $\alpha$ -cut of each FN, we generate two values: one from this  $\alpha$ -cut and the other from  $(1 - \alpha)$ -cut, and then we determine their average. As indicated in Grzegorzewski and Romaniuk (2022a,c), the antithetic approach improves the quality of some statistical inference methods. An example of how to use this function can be found in the supplementary file.

The epistemic bootstrap produces a real-valued sample based on the initial fuzzy values. Therefore, it can be easily applied to estimate various statistical measures of the input values (like the mean) or to conduct many “classical” (i.e. real-valued) statistical tests.

### 2.3 Estimation of parameters

Estimation of basic population parameters (like the mean) is a fundamental task of most statistical inference problems. Given fuzzy data, we can easily adapt the epistemic bootstrap to estimate the quantities of interest (see (Grzegorzewski and Romaniuk, 2022c)).

A general function

```
EpistemicEstimator(fuzzySample, estimator = "sd", cutsNumber = 1, bootstrapMethod = "std",
  trueValue = NA, ...)
```

can be used to determine the desired estimate from the `fuzzySample` of the specified function in `estimator`. Both the classical epistemic approach (`bootstrapMethod = "std"`) and its antithetic counterpart (`bootstrapMethod = "anti"`) are available. Since the mean is the most used statistical parameter, it can be obtained using a special function

```
EpistemicMean(fuzzySample, cutsNumber = 1, bootstrapMethod = "std", trueValue = NA, ...)
```

instead of the general command.

Besides estimates, the standard error (SE) and the mean squared error (MSE) of the considered estimators are also calculated. The SE is estimated (for  $b > 1$ ) using the formula

$$\widehat{SE} = \sqrt{\frac{1}{b-1} \sum_{k=1}^b (\hat{\theta}(\mathbb{X}^{*k}) - \bar{\hat{\theta}})^2}, \quad (3)$$

where  $\hat{\theta}(\mathbb{X}^{*k})$  is the estimator of  $\theta$  based on the epistemic bootstrap sample for the  $k$ -th  $\alpha$ -cut, and  $\bar{\hat{\theta}}$  is the overall mean for  $\hat{\theta}(\mathbb{X}^{*1}), \dots, \hat{\theta}(\mathbb{X}^{*b})$ . If the true (but usually unknown) value of  $\theta$  is set with `trueValue`, then the MSE is estimated by

$$\widehat{MSE} = \frac{1}{b} \sum_{k=1}^b (\hat{\theta}(\mathbb{X}^{*k}) - \theta)^2. \quad (4)$$

Let us estimate the median and its SE for `sample1` using 100  $\alpha$ -cuts and the classical epistemic bootstrap:

```
> set.seed(56789)

> EpistemicEstimator(sample1$value, estimator = "median", cutsNumber = 100)

$value
[1] 0.6287525

$SE
[1] 0.1705336

$MSE
[1] NA
```

To estimate the variance using bootstrap, instead of the classical well-known formula, its more sophisticated and specially corrected variant ([Grzegorzewski and Romaniuk, 2022c](#)) can be used with the function

```
EpistemicCorrectedVariance(fuzzySample, cutsNumber = 1, bootstrapMethod = "std", ...)
```

As noted in [Grzegorzewski and Romaniuk \(2022c\)](#), this estimator can more closely approximate the desired value, e.g., we have

```
> set.seed(56789)

> EpistemicCorrectedVariance(sample1$value, cutsNumber = 100$)

[1] 0.8729738
```

## 2.4 Statistical tests

The real-valued samples generated by the epistemic bootstrap can be also used for hypothesis testing. However, given several bootstrap samples, one has to clarify how to merge the obtained test statistics or the p-values ([Grzegorzewski and Romaniuk, 2022a](#)). `FuzzySimRes` contains a general function

```
EpistemicTest(sample1, sample2, algorithm = "avs", ...)
```

that can be used to activate one of the specially tailored procedures.

By setting `algorithm = "avs"` the **averaging statistic** (abbreviated as *avs*) is activated and the function

```
AverageStatisticEpistemicTest(sample1,sample2,bootstrapMethod = "std",
  test = "ks.test",cutsNumber = 1,criticalValueFunction = "KSTestCriticalValue",...)
```

is used. Similarly, by setting `algorithm = "ms"` the function

```
MultiStatisticEpistemicTest(sample1,sample2,bootstrapMethod = "std",
  test = "ks.test",cutsNumber = 1,combineMethod = "simes",...)
```

and **multi-statistic** method (denoted by *ms*) are applied. Finally, for `algorithm = "res"` the **resampling** algorithm (abbreviated as *res*) together with the function

```
ResamplingStatisticEpistemicTest(sample1,sample2,bootstrapMethod = "std",
  test = "ks.test",cutsNumber = 1,K = 1,combineMethod = "simes",...)
```

run. The above functions can be applied to both one-sample and two-sample statistical tests, where the relevant samples are entered as lists of fuzzy values. For the one-sample case, `sample2=NULL` should be set.

To use a statistical test, one has to specify the name of the respective function in `test` (e.g., `test="ks.test"` for `ks.test` from **stats** activates the Kolmogorov-Smirnov goodness-of-fit test, abbreviated further on as the KS test). User-defined functions can be also used if they have at least one or two parameters (`x` for one- or `x,y` for two-sample case, namely) and return a list of at least two values (statistic for the output test statistic, and `p.value` for the calculated p-value). In the case of the *avs* approach, the additional parameter `criticalValueFunction` is required with the name of the function calculating the p-value for a specified critical level of the considered test statistic. For the KS test, such a procedure is given by `KSTestCriticalValue` available in **FuzzySimRes**.

To conduct the test, the classical epistemic approach (`bootstrapMethod = "std"`) or its antithetic version (`bootstrapMethod = "anti"`) can be applied. The p-values (in the case of *ms* or *res* methods) are aggregated with the algorithm specified in `combineMethod`. Besides `combineMethod="mean"`, i.e. the simple averaging of p-values, all other methods are as in the package **palasso** (Rauschenberger et al., 2020).

Let us generate the second sample with the small shift in location and compare it with the previously generated `sample1` using the two-sample KS test with the *anti* and *ms* approaches for 100  $\alpha$ -cuts:

```
> set.seed(56789)

> sample2 <- SimulateSample(n=10,originalPD="rnorm",
+   parOriginalPD=list(mean=0.5,sd=1),
+   incrCorePD="runif",parIncrCorePD=list(min=0,max=0.6),
+   suppLeftPD="runif",parSuppLeftPD=list(min=0,max=1),
+   suppRightPD="runif", parSuppRightPD=list(min=0,max=1),
+   type="trapezoidal")

> EpistemicTest(sample1$value,sample2$value,algorithm = "ms",
+   bootstrapMethod="anti",cutsNumber=100)

[1] 0.1873127
```

An example of the one-sample KS test can be found in the supplementary file.

## 2.5 Real-life dataset

**FuzzySimRes** provides a fuzzy epistemic dataset `controlChartData` concerning electronic circuit thickness, which is one of the most important quality characteristics in the production of the electronic boards for vacuum cleaners (see Faraz and Shapiro (2010) for the relevant source). This dataset is given as a list of 90 TRFNs and contains 30 samples, each of size three. Every observation has its own label `X.y.z`, where `y` is a sample number, and `z` stands for the element number in a sample, e.g.

```
> controlChartData$X.1.2$
```

```
Trapezoidal fuzzy number with:
  support=[70.19,74.15],
  core=[71.4,71.4].
```

is the second value in the first sample.

### 3 Statistical applications with the package

As it was mentioned, the epistemic bootstrap provides real-valued samples generated from the initial fuzzy sample. It enables us to apply many classical statistical methods instead of using procedures specifically designed for fuzzy data (usually underdeveloped in the R environment). In the following, we present some statistical applications of such approaches for both synthetic and real-life datasets.

In the first case, using `SimulateSample`, the respective samples are generated from the probability distributions described in Table 2. Available TPFNs are grouped by their types, wherein the normal distribution with the mean  $\mu$  and standard deviation  $\sigma$  is denoted by  $N(\mu, \sigma)$ , the uniform distribution on the interval  $(a, b)$  – by  $U(a, b)$ , the exponential distribution with the parameter  $\lambda$  – by  $Exp(\lambda)$ , the Weibull distribution with the shape  $k$  and scale  $\lambda$  parameters – by  $Weib(k, \lambda)$ , and the Gamma distribution with the shape  $\alpha$  and rate  $\beta$  parameters – by  $\Gamma(\alpha, \beta)$ , respectively. In the case of the real-life dataset, the data `controlChartData` embedded in `FuzzySimRes` is applied.

In the following, only some of the results are presented in the tables and graphs to reduce the overall length of the paper. All of the outputs can be found in the supplementary script file.

Type	X	$C^l, C^r$	$S^l$	$S^r$
$\bar{F}_{(N,U,U,U)}$	$N(0, 1)$	$U(0, 0.6)$	$U(0, 1)$	$U(0, 1)$
$\bar{F}_{(Weib,Exp,Exp,Exp)}$	$Weib(2, 1)$	$Exp(5)$	$Exp(5)$	$Exp(4)$
$\bar{F}_{(\Gamma,U,U,U)}$	$\Gamma(2, 2)$	$U(0, 0.6)$	$U(0, 0.8)$	$U(0, 0.8)$

**Table 2:** Scenarios for simulating fuzzy random variables.

#### 3.1 Comparison of estimators

We start with a comparison of some estimators of the mean, variance, and median for both epistemic approaches, i.e., the *std* and *anti*. For all types of the TPFNs mentioned in Table 2, the function `EpistemicEstimator` was applied with  $b = 100$   $\alpha$ -cuts. To limit the randomness impact, each numerical experiment was repeated  $m = 1000$  times. Both small ( $n = 10$ ) and moderate ( $n = 100$ ) samples were considered.

Since the function `SimulateSample` produces also the “true values” of the fuzzy samples (i.e., their originals), it gives an opportunity (quite exceptional in real-life applications) to compare the epistemic bootstrap estimators based on fuzzy samples with the results related to these originals. Then, we can calculate the respective error – **Originals Absolute Error** (abbreviated as OAE) – that measures the absolute difference between the epistemic bootstrap estimator  $\hat{\theta}_j^*$  based on the  $j$ -th synthetic sample and its counterpart  $\hat{\theta}_j^o$  obtained from the originals for this  $j$ -th sample, i.e.,

$$\text{OAE} = \frac{1}{m} \sum_{j=1}^m |\hat{\theta}_j^* - \hat{\theta}_j^o|, \quad (5)$$

where  $m$  is the number of simulations.

In general, it seems that the *anti* approach gives better results – the resulting estimates are closer to their “true” values and the respective errors are lower (see Table 3 and the supplementary file). To facilitate the understanding of Table 3, the best outputs (i.e., the estimators that are the closest to the respective true values of the parameters, and the lowest errors in each case) are given in boldface there. Of course, the answers may vary for the different error measures (e.g., sometimes the OAE is slightly lower for the *std* approach). However, the *anti* method clearly provides the significant improvement measured with the SE, slightly less important (but still visible) in the case of the MSE. Taking into account the low additional numerical burden of this approach when it is compared with the *std* method (i.e. generation of two values: from the  $\alpha$ -cut and its  $(1 - \alpha)$  counterpart instead of only a single drawing), the *anti* algorithm should be recommended to users. The above-mentioned conclusions are similar to the ones discussed in Grzegorzewski and Romaniuk (2021, 2022c).

#### 3.2 Detection of the difference in location

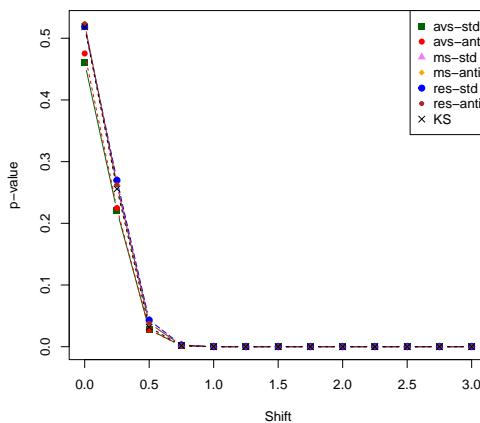
Then, we conducted the power analysis for the two-sample KS test taking into account all the considered epistemic bootstrap approaches. Two independent samples corresponding to the types of the TPFNs from Table 2 were generated and the deterministic shift was added to the second sample. As previously, both the small ( $n = 10$ ) and moderate ( $n = 100$ ) samples were considered, and each

	Mean		Variance		Median	
	std	anti	std	anti	std	anti
$\mathbb{F}_{(N,U,U,U)}, n = 10$						
Value	-0.0055	<b>-0.0053</b>	1.1476	<b>1.0854</b>	<b>-0.0034</b>	-0.0038
SE	0.1089	<b>0.0760</b>	0.2424	<b>0.1595</b>	0.1797	<b>0.1347</b>
MSE	0.1145	<b>0.1078</b>	0.3469	<b>0.2972</b>	0.1602	<b>0.1522</b>
OAЕ	0.0407	<b>0.0405</b>	0.1555	<b>0.1105</b>	0.0874	<b>0.0819</b>
$\mathbb{F}_{(N,U,U,U)}, n = 100$						
Value	0.0016	0.0016	1.1472	<b>1.0850</b>	<b>0.0004</b>	0.0006
SE	0.0345	<b>0.0242</b>	0.0999	<b>0.0498</b>	0.0705	<b>0.0573</b>
MSE	0.0115	<b>0.0110</b>	0.0526	<b>0.0305</b>	0.0179	<b>0.0168</b>
OAЕ	0.0135	<b>0.0134</b>	0.1480	<b>0.0860</b>	0.0405	<b>0.0375</b>
$\mathbb{F}_{(Weib,Exp,Exp,Exp)}, n = 10$						
Value	0.8917	<b>0.8912</b>	0.2884	<b>0.2671</b>	0.8536	<b>0.8517</b>
SE	0.0636	<b>0.0440</b>	0.0728	<b>0.0478</b>	0.0941	<b>0.0706</b>
MSE	0.0272	<b>0.0251</b>	0.0287	<b>0.0222</b>	0.0398	<b>0.0378</b>
OAЕ	<b>0.0409</b>	0.0412	0.0716	<b>0.0553</b>	<b>0.0609</b>	0.0621
$\mathbb{F}_{(Weib,Exp,Exp,Exp)}, n = 100$						
Value	<b>0.8864</b>	0.8865	0.2828	<b>0.2614</b>	0.8417	<b>0.8387</b>
SE	0.0203	<b>0.0141</b>	0.0295	<b>0.0152</b>	0.0359	<b>0.0290</b>
MSE	0.0029	<b>0.0027</b>	0.0068	<b>0.0037</b>	0.0044	<b>0.0041</b>
OAЕ	0.0127	0.0127	0.0676	<b>0.0462</b>	0.0252	<b>0.0247</b>

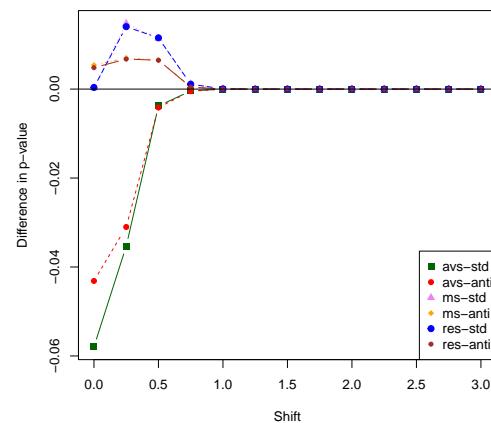
**Table 3:** Numerical comparison of the estimators (for the mean, variance, and median) and their errors (the standard error – SE, the mean squared error – MSE, and the originals absolute error – OAE) based on the epistemic bootstrap (the standard – std, and the antithetic epistemic bootstrap – anti).

numerical experiment was repeated  $m = 1000$  times. Besides the estimation of the null hypothesis rejection percentage for the significance level  $\alpha = 0.05$ , the p-values for the increasing shift were also obtained and aggregated by simple averaging.

Using `SimulateSample` which delivers the originals of the simulated fuzzy sample, we can compare the results of the epistemic bootstrap tests with their “crisp” counterpart, so the results of the classical two-sample KS test serve us as a benchmark. We can see that the estimated p-values (see Fig. 3) and power curves (see Fig. 5) for the moderate sample of TPFNs described by  $\mathbb{F}_{(N,U,U,U)}$  are very close to their respective benchmarks, especially for the shift larger than 0.75. To visualize the results better, the differences in p-values and power curves between the epistemic bootstrap approaches and the classical KS test were also calculated (see Fig. 4 and 6, respectively).

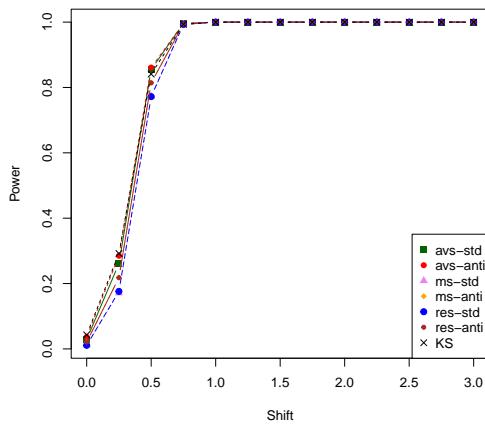


**Figure 3:** Estimated p-values of the two-sample epistemic and “crisp” KS tests for  $\mathbb{F}_{(N,U,U,U)}, n = 100$ , and shift in location.

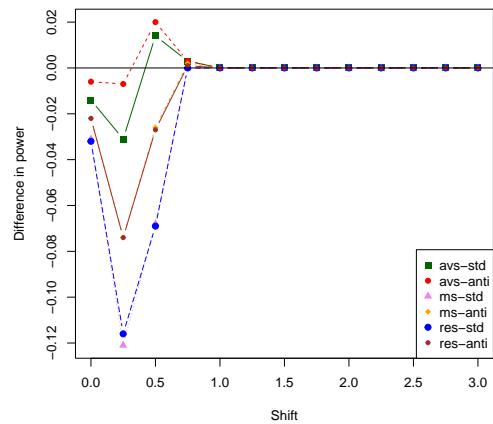


**Figure 4:** Differences in estimated p-values between the two-sample epistemic and “crisp” KS tests for  $\mathbb{F}_{(N,U,U,U)}, n = 100$ , and shift in location.

In general, the estimation error for p-values is lower when the *ms* or *res* approaches are used (especially when they are combined with the *anti* method), and the power curves are closer to the



**Figure 5:** Power curves of the two-sample epistemic and “crisp” KS tests for  $\mathbb{F}_{(N,U,U,U)}$ ,  $n = 100$ , and shift in location.



**Figure 6:** Differences in power curves between the two-sample epistemic and “crisp” KS tests for  $\mathbb{F}_{(N,U,U,U)}$ ,  $n = 100$ , and shift in location.

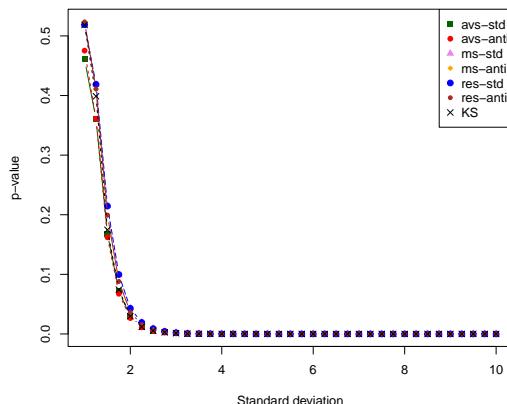
respective benchmarks for the *avs* and *ms* algorithms (the *anti* method has also a beneficial effect). Additional examples can be found in the supplementary file and Grzegorzewski and Romaniuk (2022a).

### 3.3 Detection of the difference in scale

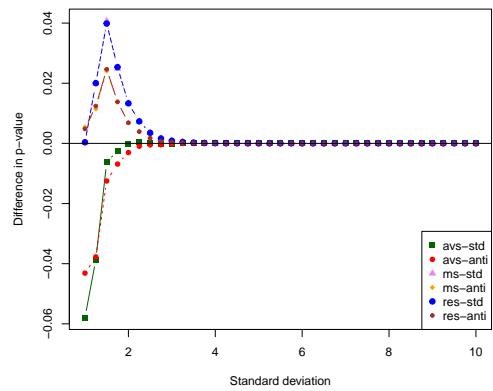
Next, we conducted the power study of tests to detect the difference in dispersion. This case was modeled by gradually increasing the standard deviation of the second sample when the first one is simulated according to  $\mathbb{F}_{(N,U,U,U)}$  type.

As previously, the p-values and power curves (see Fig. 7 and 9) were estimated for the moderate sample and the respective simulation parameters:  $m = 1000$ ,  $\alpha = 0.05$ , and  $b = 100$ . A comparison of the epistemic bootstrap approaches and our benchmark (i.e., the two-sample “crisp” KS test) was also done (see Fig. 8 and 10, respectively). It seems that the estimation error of p-values is lower for the *ms* or *res* approach and the power curves are closer to the respective results of the “crisp” KS test for the *avs* and *ms* algorithms. Thus, the *anti* method again improves the results.

An additional example for the small sample is provided in the supplementary file.



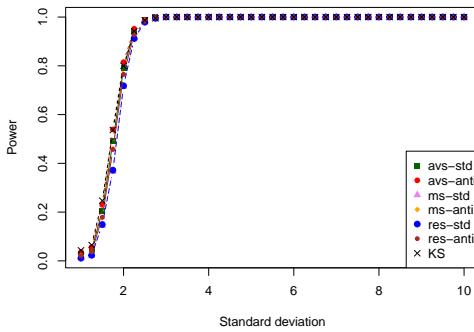
**Figure 7:** Estimated p-values of the two-sample epistemic and “crisp” KS tests for  $\mathbb{F}_{(N,U,U,U)}$ ,  $n = 100$ , difference in scale.



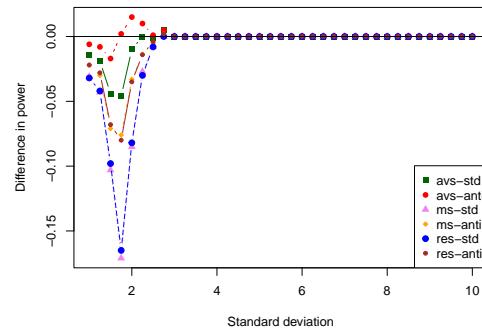
**Figure 8:** Differences in estimated p-values between the two-sample epistemic and “crisp” KS tests for  $\mathbb{F}_{(N,U,U,U)}$ ,  $n = 100$ , difference in scale.

### 3.4 Goodness-of-fit test in quality control

Finally, we applied the KS two-sample test for the manufacturing data embedded in **FuzzySimRes**. These fuzzy data can be used to build the respective control charts to check the behavior of the underlying process (Faraz and Shapiro, 2010). But in our experiment, the sample was divided



**Figure 9:** Power curves of the two-sample epistemic and “crisp” KS tests for  $\mathbb{F}_{(N,U,U,U)}$ ,  $n = 100$ , difference in scale.



**Figure 10:** Differences in power curves between the two-sample epistemic and “crisp” KS tests for  $\mathbb{F}_{(N,U,U,U)}$ ,  $n = 100$ , difference in scale.

randomly into two parts to check if they came from the same distribution (so they were “not statistically different”).

```
> set.seed(5678)
> randomSetsCCD <- sample(length(controlChartData), length(controlChartData)/2)

> EpistemicTest(controlChartData[randomSetsCCD], controlChartData[-randomSetsCCD],
+   algorithm="avs", cutsNumber=1000)

[1] 0.3319477

> EpistemicTest(controlChartData[randomSetsCCD], controlChartData[-randomSetsCCD],
+   algorithm="ms", combineMethod="mean", cutsNumber=1000)

[1] 0.433548

> EpistemicTest(controlChartData[randomSetsCCD], controlChartData[-randomSetsCCD],
+   algorithm="res", combineMethod="mean", cutsNumber=1000, K=200)

[1] 0.4616578
```

As we can see, all of the considered algorithms do not reject the null hypothesis for the KS test, even for high significance levels. These results are consistent with the findings in Faraz and Shapiro (2010). Moreover, as Grzegorzewski and Romaniuk (2024) described, the epistemic KS test clearly indicates the issues caused by the troublesome 21st subsample. It makes the process out of control and results in the lower p-values in the goodness-of-fit tests.

## 4 Conclusions

**FuzzyResampling** package delivers resampling methods developed to overcome some shortcomings of the classical Efron’s bootstrap in the fuzzy environment (see also Romaniuk and Grzegorzewski (2023)). However, this package was intended for the ontic fuzzy data.

Meanwhile, **FuzzySimRes** is a package that has a completely new purpose. The proposed epistemic bootstrap methods allow the generation of real-valued samples from the epistemic fuzzy data, which can then be directly utilized as input values for the various classical statistical procedures (like estimators, tests, etc.). It seems that the proposed methods combined with some well-known statistical techniques can be competitive with available fuzzy procedures which are not too popular among practitioners. Moreover, as was shown in the respective examples, the results of the suggested approaches implemented to imprecise data are comparable with their counterparts – the benchmarks related to the real-valued originals of the fuzzy perceptions.

Of course, further investigations on epistemic bootstrap are still required. They can be aimed both at new resampling epistemic procedures and their applications in statistical inference and machine learning.

## References

- A. Báez-Sánchez, A. Moretti, and M. Rojas-Medar. On polygonal fuzzy sets and numbers. *Fuzzy Sets and Systems*, 209:54–65, 2012. URL <https://doi.org/10.1016/j.fss.2012.04.003>. [p177]
- A. Ban, L. Coroianu, and P. Grzegorzewski. *Fuzzy Numbers: Approximations, Ranking and Applications*. Polish Academy of Sciences, Warsaw, 2015. [p176, 177]
- R. Berkachy and L. Donzé. Testing hypotheses by fuzzy methods: A comparison with the classical approach. In A. Meier, E. Portmann, and L. Terán, editors, *Applying Fuzzy Logic for the Digital Economy and Society*, pages 1–22, Cham, 2019. Springer International Publishing. URL [https://doi.org/10.1007/978-3-030-03368-2\\_1](https://doi.org/10.1007/978-3-030-03368-2_1). [p176]
- R. Berkachy and L. Donzé. *FuzzySTs: Fuzzy Statistical Tools*, 2020. URL <https://CRAN.R-project.org/package=FuzzySTs>. R package version 0.2. [p176]
- A. Colubi. Statistical inference about the means of fuzzy random variables: Applications to the analysis of fuzzy- and real-valued data. *Fuzzy Sets and Systems*, 160(3):344–356, 2009. URL <https://doi.org/10.1016/j.fss.2007.12.019>. [p176]
- L. Coroianu, M. Gagolewski, and P. Grzegorzewski. Nearest piecewise linear approximation of fuzzy numbers. *Fuzzy Sets and Systems*, 233:26–51, 2013. ISSN 0165-0114. URL <https://doi.org/10.1016/j.fss.2013.02.005>. [p176]
- L. Coroianu, M. Gagolewski, and P. Grzegorzewski. Piecewise linear approximation of fuzzy numbers: algorithms, arithmetic operations and stability of characteristics. *Soft Computing*, 23(19):9491–9505, Oct 2019. URL <https://doi.org/10.1007/s00500-019-03800-2>. [p177]
- I. Couso and D. Dubois. Statistical reasoning with set-valued information: Ontic vs. epistemic views. *International Journal of Approximate Reasoning*, 55:1502–1518, 2014. doi: j.ijar.2013.07.002. URL <https://doi.org/10.1016/j.ijar.2013.07.002>. [p175, 177]
- A. C. Davison and D. V. Hinkley. *Bootstrap Methods and their Application*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press, 1997. URL <https://doi.org/10.1017/CBO9780511802843>. [p175]
- B. Efron and R. J. Tibshirani. *An Introduction to the Bootstrap*. Number 57 in Monographs on Statistics and Applied Probability. Chapman & Hall/CRC, Boca Raton, Florida, USA, 1993. [p175, 177]
- A. Faraz and A. F. Shapiro. An application of fuzzy random variables to control charts. *Fuzzy Sets and Systems*, 161(20):2684–2694, 2010. ISSN 0165-0114. URL <https://doi.org/10.1016/j.fss.2010.05.004>. [p178, 183, 186, 187]
- M. Gagolewski and J. Caha. *FuzzyNumbers Package: Tools to Deal with Fuzzy Numbers in R*, 2021. URL <https://github.com/gagolews/FuzzyNumbers/>. [p175]
- J. D. Gibbons and S. Chakraborti. *Nonparametric Statistical Inference*. Chapman and Hall/CRC, 2010. [p178]
- M. Gil, M. Montenegro, G. González-Rodríguez, A. Colubi, and M. Casals. Bootstrap approach to the multi-sample test of means with imprecise data. *Computational Statistics and Data Analysis*, 51: 148–162, 2006. URL <https://doi.org/10.1016/j.csda.2006.04.018>. [p177]
- G. González-Rodríguez, A. Colubi, and W. Trutschnig. Simulation of fuzzy random variables. *Information Sciences*, 179(5):642–653, 2009. ISSN 0020-0255. URL <https://doi.org/10.1016/j.ins.2008.10.018>. [p176]
- P. Grzegorzewski and M. Romaniuk. Epistemic bootstrap for fuzzy data. In *Joint Proceedings of IFSA-EUSFLAT-AGOP 2021 Conferences*, pages 538–545. Atlantis Press, 2021. URL <https://doi.org/10.2991/asum.k.210827.071>. [p175, 177, 178, 180, 184]
- P. Grzegorzewski and M. Romaniuk. Bootstrapped Kolmogorov-Smirnov test for epistemic fuzzy data. In D. Ciucci, I. Couso, J. Medina, D. Ślezak, D. Petturiti, B. Bouchon-Meunier, and R. R. Yager, editors, *Information Processing and Management of Uncertainty in Knowledge-Based Systems*, pages 494–507, Cham, 2022a. Springer International Publishing. URL [https://doi.org/10.1007/978-3-031-08974-9\\_39](https://doi.org/10.1007/978-3-031-08974-9_39). [p175, 178, 179, 180, 181, 182, 186]

- P. Grzegorzewski and M. Romaniuk. Bootstrap methods for fuzzy data. In K. T. Atanassov, V. Atanassova, J. Kacprzyk, A. Kałuszko, M. Krawczak, J. W. Owsiński, S. S. Sotirov, E. Szmidt, and S. Zadrożny, editors, *Uncertainty and Imprecision in Decision Making and Decision Support: New Advances, Challenges, and Perspectives*, pages 28–47, Cham, 2022b. Springer International Publishing. URL [https://doi.org/10.1007/978-3-030-95929-6\\_3](https://doi.org/10.1007/978-3-030-95929-6_3). [p176, 177, 178]
- P. Grzegorzewski and M. Romaniuk. Bootstrap methods for epistemic fuzzy data. *International Journal of Applied Mathematics and Computer Science*, 32(2):285–297, 2022c. URL <https://doi.org/10.34768/amcs-2022-0021>. [p175, 177, 178, 180, 181, 182, 184]
- P. Grzegorzewski and M. Romaniuk. Bootstrapped tests for epistemic fuzzy data. *International Journal of Applied Mathematics and Computer Science*, 34, 2024. (in print). [p177, 178, 180, 187]
- P. Grzegorzewski, O. Hryniewicz, and M. Romaniuk. Flexible resampling for fuzzy data. *International Journal of Applied Mathematics and Computer Science*, 30(2):281–297, 2020. URL <https://doi.org/10.34768/amcs-2020-0022>. [p176, 177]
- S. Heilpern. The expected value of a fuzzy number. *Fuzzy Sets and Systems*, 47(1):81–86, 1992. ISSN 0165-0114. URL <https://www.sciencedirect.com/science/article/pii/0165011492900629>. [p177]
- G. Hesamian and S. Taheri. Linear rank tests for two-sample fuzzy data: a p-value approach. *Journal of Uncertain Systems*, 7(2):129–137, 2013. [p178]
- G. James, D. Witten, T. Hastie, and R. Tibshirani. *An Introduction to Statistical Learning with Applications in R*. Springer New York, New York, 2021. URL <https://doi.org/10.1007/978-1-0716-1418-1>. [p175]
- R. Kruse. The strong law of large numbers for fuzzy random variables. *Information Sciences*, 28:233–241, 1982. URL [https://doi.org/10.1016/0020-0255\(82\)90049-4](https://doi.org/10.1016/0020-0255(82)90049-4). [p177]
- H. Kwakernaak. Fuzzy random variables, part I: Definitions and theorems. *Information Sciences*, 15:1–15, 1978. URL [https://doi.org/10.1016/0020-0255\(78\)90019-1](https://doi.org/10.1016/0020-0255(78)90019-1). [p177]
- A. Lubiano and S. de la Rosa de Saa. *FuzzyStatTra: Statistical Methods for Trapezoidal Fuzzy Numbers*, 2017. URL <https://CRAN.R-project.org/package=FuzzyStatTra>. R package version 0.1. [p176]
- M. A. Lubiano, M. Montenegro, B. Sinova, S. de la Rosa de Sáa, and M. A. Gil. Hypothesis testing for means in connection with fuzzy rating scale-based data: algorithms and applications. *European Journal of Operational Research*, 251:918–929, 2016. URL <https://doi.org/10.1016/j.ejor.2015.11.016>. [p176, 177, 178]
- M. A. Lubiano, A. Salas, C. Carleos, S. de la Rosa de Sáa, and M. A. Gil. Hypothesis testing-based comparative analysis between rating scales for intrinsically imprecise data. *International Journal of Approximate Reasoning*, 88:128–147, 2017. URL <https://doi.org/10.1016/j.ijar.2017.05.007>. [p178]
- M. Montenegro, A. Colubi, M. Rosa Casals, and M. Ángeles Gil. Asymptotic and bootstrap techniques for testing the expected value of a fuzzy random variable. *Metrika*, 59(1):31–49, 2004. URL <https://doi.org/10.1007/s001840300270>. [p176, 177]
- A. Parchami. *Sim.PLFN: Simulation of Piecewise Linear Fuzzy Numbers*, 2017. URL <https://CRAN.R-project.org/package=Sim.PLFN>. R package version 1.0. [p176]
- A. Parchami, P. Grzegorzewski, and M. Romaniuk. Statistical simulations with LR random fuzzy numbers. *Statistical Papers*, 2024. URL <https://doi.org/10.1007/s00362-024-01533-5>. (in print). [p176, 177]
- R Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria, 2023. URL <https://www.R-project.org/>. [p179]
- A. B. Ramos-Guajardo, Á. Blanco-Fernández, and G. González-Rodríguez. *Applying Statistical Methods with Imprecise Data to Quality Control in Cheese Manufacturing*, pages 127–147. Springer International Publishing, Cham, 2019. URL [https://doi.org/10.1007/978-3-030-03201-2\\_8](https://doi.org/10.1007/978-3-030-03201-2_8). [p178]
- A. Rauschenberger, I. Ciocanea-Teodorescu, R. X. Menezes, M. A. Jonker, and M. A. van de Wiel. Sparse classification with paired covariates. *Advances in Data Analysis and Classification*, 14:571–588, 2020. URL <https://doi.org/10.1007/s11634-019-00375-6>. [p183]

- M. Romaniuk. On some applications of simulations in estimation of maintenance costs and in statistical tests for fuzzy settings. In A. Steland, E. Rafajłowicz, and O. Okhrin, editors, *Stochastic Models, Statistics and Their Applications*, pages 437–448. Springer International Publishing, 2019. URL [https://doi.org/10.1007/978-3-030-28665-1\\_33](https://doi.org/10.1007/978-3-030-28665-1_33). [p175]
- M. Romaniuk and P. Grzegorzewski. Resampling fuzzy numbers with statistical applications: FuzzyResampling package. *The R Journal*, 15:271–283, 2023. URL <https://doi.org/10.32614/RJ-2023-036>. [p176, 187]
- M. Romaniuk and O. Hryńiewicz. Interval-based, nonparametric approach for resampling of fuzzy numbers. *Soft Computing*, 23:5883–5903, 2019. URL <https://doi.org/10.1007/s00500-018-3251-5>. [p176]
- M. Romaniuk, P. Grzegorzewski, and O. Hryńiewicz. *FuzzyResampling: Resampling Methods for Triangular and Trapezoidal Fuzzy Numbers*, 2022. URL <https://CRAN.R-project.org/package=FuzzyResampling>. R package version 0.6.0. [p176]
- B. Sinova, M. A. Gil, and S. Van Aelst. M-estimates of location for the robust central tendency of fuzzy data. *IEEE Transactions on Fuzzy Systems*, 24(4):945–956, 2016. URL <https://doi.org/10.1109/TFUZZ.2015.2489245>. [p176]
- W. Trutschnig, M. A. Lubiano, and J. Lastra. SAFD — an R package for statistical analysis of fuzzy data. In C. Borgelt, M. Á. Gil, J. M. Sousa, and M. Verleysen, editors, *Towards Advanced Data Analysis by Combining Soft Computing and Statistics*, pages 107–118. Springer Berlin Heidelberg, Berlin, Heidelberg, 2013. URL [https://doi.org/10.1007/978-3-642-30278-7\\_10](https://doi.org/10.1007/978-3-642-30278-7_10). [p176]

Maciej Romaniuk  
Systems Research Institute Polish Academy of Sciences  
Newelska 6, 01-447 Warsaw  
Poland  
WIT Academy  
Newelska 6, 01-447 Warsaw  
Poland  
(0000-0001-9649-396X)  
[mroman@ibspan.waw.pl](mailto:mroman@ibspan.waw.pl)

Przemysław Grzegorzewski  
Faculty of Mathematics and Information Science, Warsaw University of Technology  
Koszykowa 75, 00-662 Warsaw  
Poland  
Systems Research Institute Polish Academy of Sciences  
Newelska 6, 01-447 Warsaw  
Poland  
(0000-0002-5191-4123)  
[przemyslaw.grzegorzewski@pw.edu.pl](mailto:przemyslaw.grzegorzewski@pw.edu.pl)

Abbas Parchami  
Department of Statistics, Faculty of Mathematics and Computer  
Shahid Bahonar University of Kerman, Kerman  
Iran  
(0000-0002-0593-7324)  
[parchami@uk.ac.ir](mailto:parchami@uk.ac.ir)

# Fast and Flexible Search for Homologous Biological Sequences with DECIPHER v3

by Erik S. Wright

**Abstract** The rapid growth in available biological sequences makes large-scale analyses increasingly challenging. The DECIPHER package was designed with the objective of helping to manage big biological data, which is even more relevant today than when the package was first introduced. Here, I present DECIPHER version 3 with improvements to sequence databases, as well as fast and flexible sequence search. DECIPHER now allows users to find regions of local similarity between sets of DNA, RNA, or amino acid sequences at high speed. I show the power of DECIPHER v3 by (a) comparing against BLAST and MMseqs2 on a protein homology benchmark, (b) locating nucleotide sequences in a genome, (c) finding the nearest sequences in a reference database, and (d) searching for orthologous sequences common to human and zebrafish genomes. Search hits can be quickly aligned, which enables a variety of downstream applications. These new features of DECIPHER v3 make it easier to manage big biological sequence data.

## 1 Introduction

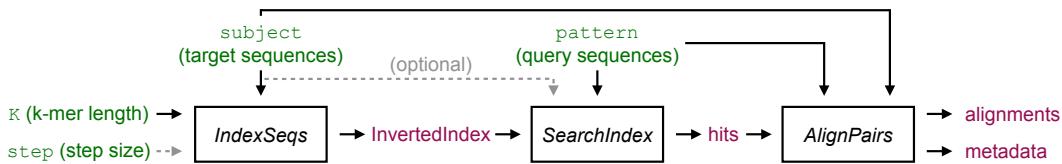
Searching is one of the most common operations in biological sequence analysis with a wide variety of applications. The [Biostrings](#) package includes the Aho-Corasick, Boyer-Moore, and Shift-Or algorithms that are frequently used for matching short sequences with relatively few differences, such as oligonucleotide probes or primers. More difficult queries typically rely on calling programs outside of R. For example, error-prone long reads can be rapidly mapped with programs that make use of seed-and-extend approaches to identify matching regions between longer sequences ([Sahlin et al., 2023](#)). More distantly related homologs can be found with BLAST ([Altschul et al., 1997](#)) or quicker alternatives DIAMOND ([Buchfink et al., 2021](#)), LAST ([Kielbasa et al., 2011](#)), and MMseqs2 ([Steinegger and Soding, 2017](#)). The ubiquity of these tools in sequence analyses attests to the importance of searching through biological sequences.

BLAST is perhaps the most commonly used bioinformatics tool. With over 30 years since its first release, BLAST still represents the gold standard of sensitivity and specificity for sequence search, and competitors typically compare themselves to BLAST's speed and accuracy. Both DIAMOND version 2 and MMseqs2 offer modes with accuracy rivaling BLAST and several fold its speed for practical search tasks. There is a clear trade-off between search sensitivity and speed, with different programs specializing at opposite ends of the spectrum. Improved sensitivity can be obtained by converting initial search hits into a profile and iteratively repeating the search, which is the basis of programs such as HHblits ([Remmert et al., 2011](#)), PSI-BLAST ([Altschul et al., 1997](#)), and PHMMER ([Eddy, 2011](#)). Greater speed is achievable when sensitivity is less of a concern, such as when mapping nearly identical short reads to a genome. For example, the [Rsubread](#) package maps and quantifies short reads from R using the fast subread aligner ([Liao et al., 2019](#)).

To my knowledge, there is currently no R-based implementation of a generalized homology search program, although there exist interfaces to external programs that can be called from within R. The [DECIPHER](#) package is a natural fit for large-scale homology search because of its functionality for managing big biological data ([Wright, 2016](#)). With this in mind, here I introduce [DECIPHER](#) v3 with the addition of general-purpose search and improvements to sequence databases. [DECIPHER](#) can now find locally homologous regions between two sets of nucleotide or protein sequences, report their scores and, if desired, align the search hits. In the examples below, I show the application of [DECIPHER](#)'s search functions to common search problems. [DECIPHER](#) v3 is available from the Bioconductor package repository.

## 2 Design of the sequence search algorithm implemented in DECIPHER

The development of search functionality within [DECIPHER](#) was inspired by user experiences with BLAST. First, BLAST is very sensitive but relatively slow for large sets of query and target sequences. Alternative search programs showed it was possible to exceed BLAST's speed, potentially at the expense of sensitivity. Second, BLAST results are limited by multiple parameters that have led to misinterpretation ([Gonzalez-Pech et al., 2019; Shah et al., 2019; Madden et al., 2019](#)). Third, reliance on BLAST requires running software outside of the R environment, and it would be preferable to have similar functionality available within R. These issues guided the development of fast and versatile



**Figure 1:** Overview of the main parameters and functions comprising homology search in DECIPHER.

search and pairwise alignment functions in DECIPHER.

As shown in Figure 1, homology search in DECIPHER v3 is separated into three functions: indexing the target sequences with IndexSeqs, searching a query (pattern) for significant target (subject) hits with SearchIndex, and, optionally, aligning those hits with AlignPairs. The statistical basis for DECIPHER’s heuristic search algorithm is related to that of BLAT (Kent, 2002). Assuming equal frequencies of letters from an alphabet of size  $A$ , the probability of finding at least one match of length  $k$  among  $T$  target k-mers is:  $\text{probability} = 1 - (1 - (1/A)^k)^T$ . When this probability is near zero, as is typically the case when  $4 \geq A \leq 20$  and  $k > 1$ , it can be approximated as  $T/A^k$ . It is convenient to score using the negative log-likelihood of probability:  $-\log_e(\text{probability}) \approx k * \log_e(A) - \log_e(T)$ . The  $k * \log_e(A)$  term can be replaced with the negative  $\log_e$  of the normalized frequency ( $f$ ) expected for a specific k-mer when  $k$  is variable in length or alphabet composition is non-uniform (i.e.,  $f \approx A^{-k}$ ).

This scoring approach can be extended to approximate the probability of finding multiple k-mer matches within a region of a given size (Kent, 2002). However, the existence of insertions or deletions (i.e., gaps) between matches complicates the application of an analytical scoring approach. Hence, DECIPHER uses a search formula with empirical costs for the number of positions ( $sep$ ) and gaps ( $gap$ ) between chained matches. The final scoring formula for a hit between a query and target sequence is:

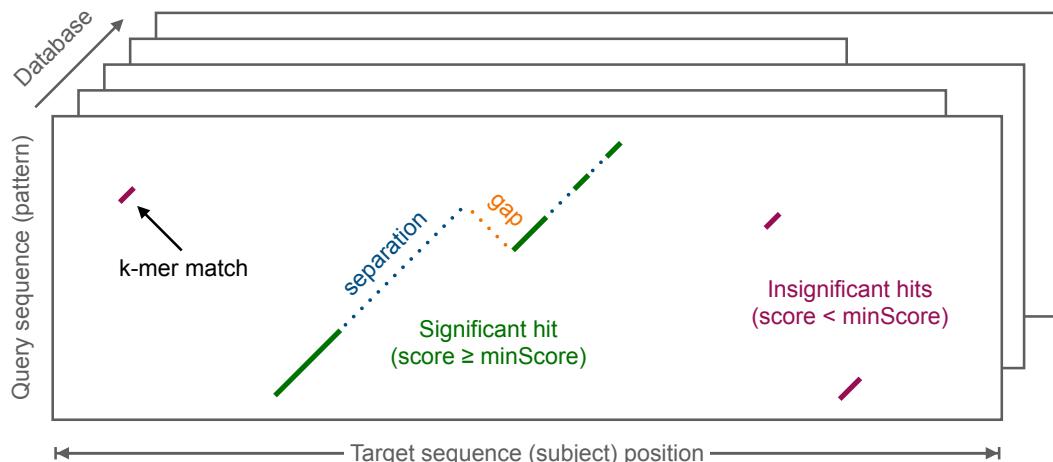
$$\text{score} = \sum_{i=1}^n (-\log_e(f)) + \sum_{i=1}^{n-1} (\text{sepCost} * \sqrt{\text{sep}}) + \sum_{i=1}^{n-1} (\text{gapCost} * \sqrt{\text{gap}}) - \log_e(Q) - \log_e(T/\text{step})$$

Where  $f$  is the expected frequency of the  $i^{th}$  k-mer match,  $sep$  is the number of positions separating neighboring (i.e.,  $i$  and  $i+1$ ) matches,  $gap$  is the minimum length of insertions or deletions (indels) between neighboring matches,  $Q$  is the number of k-mers queried, and  $T$  is the number of k-mers in the target sequence, and  $step$  is the staggering between adjacent k-mers recorded in the inverted index.

Affine (i.e., constant plus linear) gap penalties are popular for sequence alignment, although it is known that the distribution of indel lengths is approximated by a Zipfian distribution (Wygoda et al., 2024). In the absence of alignment it is only feasible to estimate the total number of gaps between k-mers rather than their individual lengths. For this reason, the square root function was selected to penalize distances between adjacent k-mer matches in SearchIndex. The square root functional form gradually decreases the marginal penalty for longer distances, which reflects the diminishing probability of observing longer gaps. Since the probability of gaps relative to mismatches is unknown, both  $sepCost$  and  $gapCost$  are empirically optimized parameters.

The score for each hit represents the significance of the match, with scores above  $\text{minScore}$  defined as significant. To correct for multiple testing, the default minimum score for reporting hits is  $\log_e(D - T + 1)$ , where  $D$  is the number of unmasked positions in the database. By default, masking is applied to low complexity, ambiguous, and repeat regions of both the query and target sequences to avoid false positive hits within regions that are not due to common descent (i.e., homology). Masking is also applied to target k-mers that are far more numerous than expected in order to accelerate the search process. Since all remaining positions are fully indexed, DECIPHER’s search functions guarantee that all matches of at least length  $k$  are scored. Protein matches are found in a reduced amino acid alphabet with residue groups: {A}, {C}, {D, E}, {F, W, Y}, {G}, {H}, {I, L, M, V}, {N}, {P}, {Q}, {R, K}, and {S, T}.

The SearchIndex algorithm is written in C code and begins by calculating the normalized expected frequency ( $f$ ) for every k-mer. Query (pattern) sequences are processed in parallel if DECIPHER was compiled with OpenMP enabled and processors is not 1. After masking each query, the locations of all matching unmasked target (subject) k-mers are extracted from the inverted index and sorted by their location within each target sequence. Next, k-mers with starting positions separated by exactly the step size are collapsed into contiguous matches by combining their scores. If subject sequences are provided to SearchIndex then the matches are extended to the left and right until encountering their nearest neighbor or their score decreases. Matches between query and target sequences are chained together using dynamic programming (Abouelhoda and Ohlebusch, 2003). The minimum score is applied to the set of top scoring chains to generate the candidate hit set (Fig. 2). Top scoring hits are returned after imposing any user-specified limits on the number of hits per target (perSubjectLimit)



**Figure 2:** Matches between a query sequence and each target sequence in a database are chained into high scoring hits. Hits are scored based on the number of positions separating matches and the minimum number of implied gaps between matches.

and query (`perPatternLimit`) for each sequence queried.

The selection of an appropriate value for k-mer length (i.e.,  $k$ ) is critically important to balance sensitivity and search speed. An advantage of **DECIPHER**'s scoring formulation is that it is possible to approximate *sensitivity* at a value of  $k$  for the objective of finding homologous matches of length  $Q$  with a given percent identity ( $PID$ ). If  $k$  is unknown, the `IndexSeqs` function allows users to provide their goal *sensitivity*,  $Q$ , and  $PID$  to automatically estimate a reasonable value for  $k$ . Relatedly, a step size ( $s$ ) between target k-mers of greater than one position can be provided to make k-mers partly overlapping or non-overlapping, which reduces the amount of memory required for the inverted index at the expense of decreased *sensitivity*.

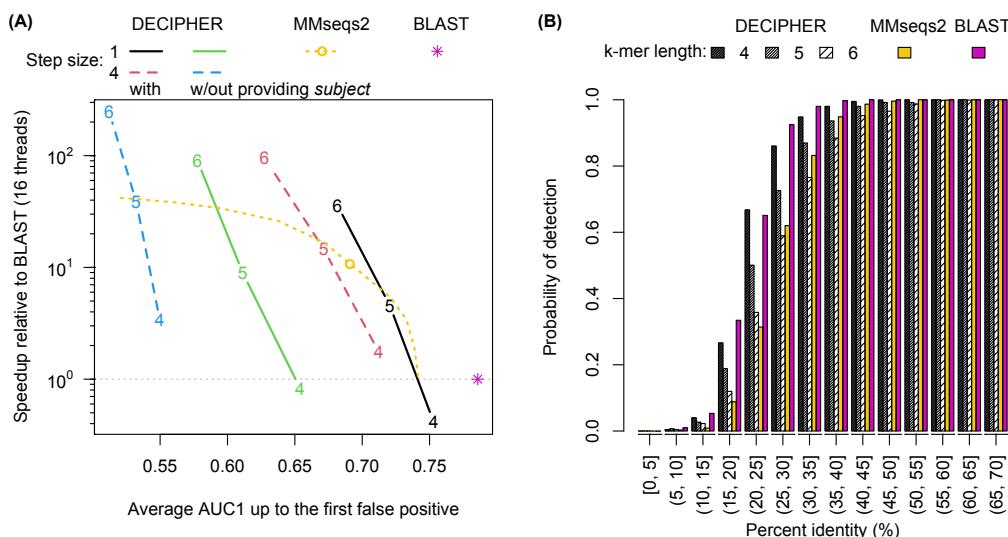
### 3 Comparing programs on a protein search benchmark

Homology search programs are typically compared on benchmarks composed of sequences assigned to protein families based on their structures (Fox et al., 2014). To this end, I downloaded the set of 93,153 non-redundant amino acid sequences from the SCOPe (v2.08) database, which contains representative sequences from 5,047 protein families. Families may contain homologous sequences with very high or low similarity, representing a wide gamut of search cases. To create an independent test set, I withheld a single protein from each of the 4,041 protein families with more than one representative sequence. I constructed a matching decoy set by reversing the held-out proteins, as this is a common benchmarking approach since proteins do not evolve by reversal (Steinegger and Soding, 2017). Therefore, the final query set consisted of 8,082 sequences and the target set contained 89,112 sequences.

Ideally, all sequences in the same family as the query sequence would be found with higher scores than any sequences in the decoy set. For each search, accuracy can be approximated as the fraction of a protein's family that was given higher scores than the highest scoring false positive hit to the matching decoy protein, an accuracy measure commonly known as AUC1 (Buchfink et al., 2021). Figure 3 shows that providing `SearchIndex` with subject (target) sequences for k-mer extension improved accuracy at the expense of approximately 2- to 3-fold slower search speed. As expected, lower values of k-mer length and step size improved accuracy and decreased search speed. Impressively, **DECIPHER** showed comparable accuracy and speed to MMseqs2 (release 15-6f452) at some values of k-mer length and step size. Both programs were less sensitive than protein BLAST (v2.15.0), but able to achieve much higher speeds under some parameterizations. All programs accurately detected the vast majority of true positives above 40% identity (Fig. 3) despite the short lengths of many SCOPe sequences.

### 4 Example 1: Locating nucleotide sequences in a genome

A common application of sequence search is to find the location of nucleotide sequences in a genome. This may include query sequences that are long or short, contiguous or discontiguous (e.g., exons), and similar or distant to the genome. Consider the case of locating promoter elements in a bacterial genome. Here, the query sequences are relatively short but may contain mismatches or indels at



**Figure 3:** The SCOPe database of amino acid sequences assigned to protein families was used to compare search programs. (A) Speed is shown relative to protein BLAST (horizontal gray line) for easier comparison. All programs were configured to use the maximum available processor threads. DECIPHER achieved greater accuracy (average AUC1) with inclusion of subject sequences to allow extension of k-mer matches. The value of k-mer length ( $k$ ) is shown on each DECIPHER point. Lower values of  $k$  and step size improved accuracy at the expense of speed. MMseqs2 results are shown across a range of user-specified sensitivities from  $-s 1$  to  $-s 10$ , with the default of 5.7 shown as a point. BLAST was run under its default settings, which reports hits with E-values less than 10. (B) The average fraction of true positives given higher scores than those of their corresponding decoy false positive increases with percent identity between the true positives and query sequences. DECIPHER results are shown at different k-mer lengths for a step size of 1 when supplying the subject sequences (i.e., the black line in (A)). MMseqs2 and BLAST results are shown at their default settings.

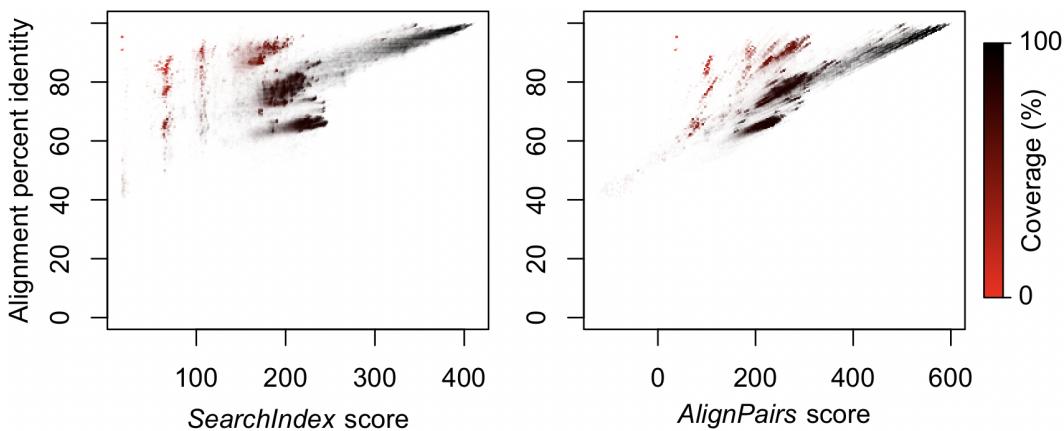
some sites relative to the genome. As a case study, I downloaded the set of 1,440 confirmed and strong confidence promoters in the *E. coli* genome from RegulonDB (Salgado et al., 2024). Given that promoters can map to the forward or reverse strand, it is necessary to build an index from both strands using IndexSeqs. Since the optimal value for  $k$  is unknown, it is possible to instead specify the goal of 90% sensitivity for query sequences of length 60 nucleotides and 80% identity to the genome, resulting in an estimated  $k$  of 5. Then, SearchIndex can be used to obtain a data frame containing hits meeting the minimum score:

```
query <- readDNAStringSet("<<path to promoters.fas>>")
target <- readDNAStringSet("<<path to genome.fas>>")
target <- c(target, reverseComplement(target)) # search both stands
index <- IndexSeqs(target, sensitivity=0.9, percentIdentity=80, patternLength=60)
hits <- SearchIndex(query, index, perSubjectLimit=0, processors=NULL) # use all CPUs
```

Performing this search for 1,440 promoter sequences of length 60 takes 4.1 seconds on a machine with 16 processor cores and 45 seconds using a single processor. In contrast, searching for a single promoter using the **Biostrings** function matchPattern takes 6.7 seconds when allowing for up to 12 mismatches or gaps (i.e.,  $\geq 80\%$  identity). This equates to a speedup of over 2,000-fold with 16 processors or 200-fold with a single processor using **DECIPHER**. Alternatively, the **Biostrings** function matchPDict can search for all patterns without indels in 22 seconds. The SearchIndex function finds 1,593 significant hits, while the matchPattern function finds 1,478 matches and the matchPDict function finds 1,455 matches. All matchPDict matches are found by matchPattern, and all but one of the matchPattern matches overlap with a hit found by SearchIndex.

## 5 Example 2: Aligning to the nearest neighbor in a reference database

Perhaps the most common use of biological sequence search is to find homologous sequences in a reference set. For example, the online BLAST search tool allows selection from a wide variety of databases, including many that are taxon or gene specific. To test the use of **DECIPHER** for this purpose, I downloaded the set of 17,103 sequences that are part of the Fungi Internal Transcribed



**Figure 4:** The scores output by *SearchIndex* and *AlignPairs* for each hit (point) show similar correlations with percent identity. Lower read (query) coverage often resulted in higher percent identities for the same score, implying the hit corresponded to a smaller aligned region with higher *PID*.

Spacer (ITS) project (BioProject PRJNA177353), which serves as a marker region for fungal taxonomy (Schoch et al., 2014) and is searchable from BLAST’s online tool. As a query, I downloaded a set of 4,823 unique ITS reads (SRA accession SRR5098782) obtained from the gut of a healthy human subject as part of the Human Microbiome Project (Nash et al., 2017). Since every read is expected to be homologous to every reference sequence, it is useful to limit the number of hits per read by specifying a *perPatternLimit*.

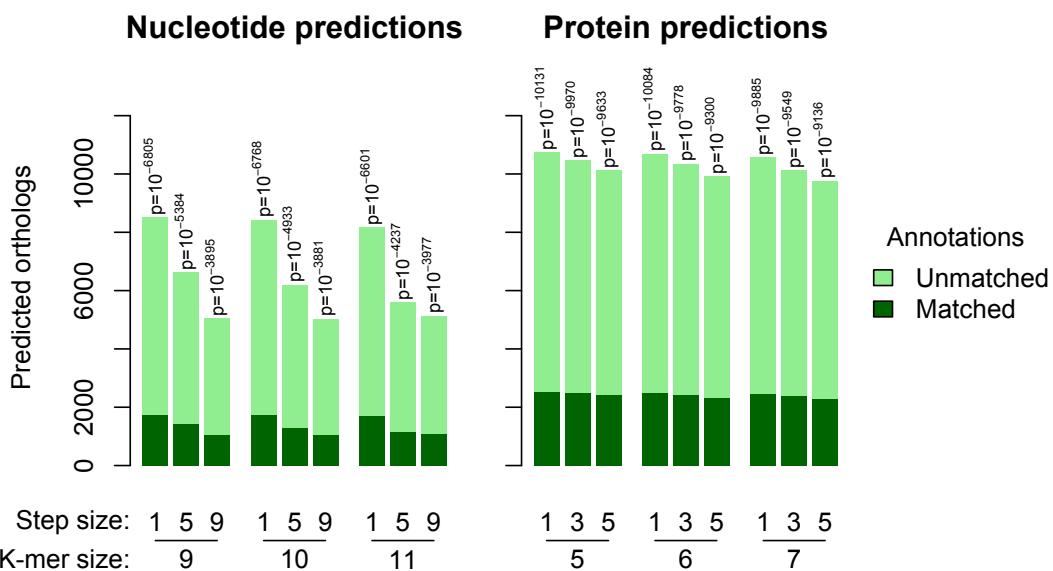
```
index <- IndexSeqs(reference, K=8L, processors=NULL)
hits <- SearchIndex(reads, index, perPatternLimit=100, processors=NULL)
aligned <- AlignPairs(reads, reference, hits, processors=NULL)
PID <- 100*aligned$Matches/aligned$AlignmentLength
coverage <- 100*(aligned$Matches + aligned$Mismatches)/width(reads)[aligned$Pattern]
```

Searching the inverted index required 14 seconds using 16 processors and a k-mer length of 8. This resulted in 475,443 search hits, which were aligned by *AlignPairs* in 3.7 seconds using 16 processors. Alignment provides a means to easily calculate the *PID* between each query/target hit, making it straightforward to determine the nearest reference sequence to each read using R code. There are multiple ways to formulate *PID* (Raghava and Barton, 2006), all of which are possible to calculate from the metadata output by *AlignPairs*. Here, we will compute *PID* of the alignment, which is correlated with the scores output by *SearchIndex* and *AlignPairs* (Fig. 4). This example illustrates the power and speed offered by the new search functions in **DECIPHER** v3.

## 6 Example 3: Searching for homology between genomes

Determination of orthology represents a challenging application of sequence search, as establishing homology is difficult in the presence of substantial sequence divergence. As an example, I used **DECIPHER**’s search functions to find orthologous proteins in humans and zebrafish (*Danio rerio*), which share a common ancestor approximately 450 million years ago (Bi et al., 2021). Many studies use reciprocal best hits to define orthologs by intersecting the best hit for each sequence from a first (query) genome to a second (target) genome with the best hit in the other direction (Hernandez-Salmeron and Moreno-Hagelsieb, 2020). Writing a *ReciprocalSearch* function in R is straightforward, given the two sets of sequences, k-mer length (*k*), and step size (*s*):

```
ReciprocalSearch <- function(seqs1, seqs2, k, s, CPUs=NULL) { # bidirectional searching
  index1 <- IndexSeqs(seqs2, K=k, step=s, verbose=FALSE, processors=CPUs)
  hits1 <- SearchIndex(seqs1, index1, perPatternLimit=1, verbose=FALSE, processors=CPUs)
  index2 <- IndexSeqs(seqs1, K=k, step=s, verbose=FALSE, processors=CPUs)
  hits2 <- SearchIndex(seqs2, index2, perPatternLimit=1, verbose=FALSE, processors=CPUs)
  m <- match(hits1$Pattern, hits2$Subject) # intersect the best hits
  w <- which(!is.na(m))
  w <- w[hits1$Subject[w] == hits2$Pattern[m[w]]]
  hits1[w,]
}
```



**Figure 5:** More orthologous pairs were predicted when searching protein (amino acid) sequences rather than their corresponding coding (nucleotide) sequences. Lower values of k-mer length and, especially, step size permitted the identification of a greater number of putative orthologs with or without matching protein annotations. The approximate p-value (one-sided binomial test) for observing a given number of matched annotations by chance is shown above each bar.

Since many animal proteins are isoforms, I reduced each set to the unique protein sequences having distinct protein names and distinct gene names, maintaining the longest isoform per gene. This process resulted in 20,169 human proteins (BioProject PRJNA559484) and 23,005 zebrafish proteins (BioProject PRJNA11776). Applying the reciprocal search strategy predicted substantially more orthologs when searching with protein rather than their corresponding coding (nucleotide) sequences (Fig. 5). Consistent with the previous results, smaller k-mer lengths generally provided more orthologous hits, while increasing step size had a more detrimental effect than increasing k-mer length.

Additional information is required to determine whether the greater number of reciprocal best hits is due to more correct or incorrect predictions. It is possible to use annotation agreement as a proxy for the quality of orthology predictions (Altenhoff et al., 2016). In total, 3,481 distinct protein names are common to human and zebrafish genomes, and the remaining proteins either lack a protein name, are without an ortholog in the other species, or follow alternative nomenclature. The total number of predicted orthologous pairs was correlated with the number of matching annotations, with putative protein orthologs having far more matched annotations than predicted nucleotide orthologs. Furthermore, increases in matched annotations were statistically significant (Fig. 5). These results confirm the expectations that amino acid search provides greater sensitivity than nucleotide search and lower values of step size improve sensitivity.

## 7 Coupling search to improved sequence databases in DECIPHER v3

The requirement for in-memory indexing is a downside of DECIPHER's search functions relative to BLAST for large target databases. For example, loading all chromosomes of the complete human genome into memory requires 2.9 GB per strand before even building an inverted index. To mitigate this downside, users can construct workflows that process queries in subsets, such as searching individual chromosomes one at a time. For example, it is possible to search a database connection (`dbConn`) in batches of a given number of sequences (i.e., `batchSize` in the example below):

```
n <- 0
batchSize <- 1000
repeat {
  seqs <- SearchDB(dbConn, limit=paste(n, batchSize, sep=","))
  hits <- SearchIndex(seqs, index)
  # do something with hits here before the next batch
  if (length(seqs) < batchSize)
    break
  n <- n + batchSize
}
```

}

Since inception, **DECIPHER** has supported the curation of large amounts of biological sequences through SQLite databases. **DECIPHER** v3 now supports multiple SQL database engines. This is accomplished by relying on generic SQL syntax in conjunction with an interface to the **DBI** package such as **RMariaDB**, **RPostgres**, or **RSQlite**. Supporting multiple database drivers allows users to tailor their database configuration to their needs. While SQLite is portable and straightforward to configure, other database types allow multiple users and distributed systems. This empowers users to construct more complex workflows by accessing or updating subsets of a large sequence database.

## 8 Discussion

The new search functions in **DECIPHER** v3 have a wide variety of potential applications. Here, I showed how users can control the trade-off between speed and accuracy by adjusting k-mer length and step size. The search functions described above are parallelized with OpenMP, which provides a substantial speed-up when multiple processors are enabled. Since all searches are in-memory, they can be used in conjunction with **DECIPHER**'s database functionality to create adaptable and distributable workflows. Unlike application-specific search tools, I designed **DECIPHER**'s search functions to be highly flexible and empower users with more control. My hope is that these search functions will be useful for creating customized workflows within the R environment.

A nice feature of BLAST is its general purpose use, which is mirrored by **DECIPHER**'s search functions. Some alternative search methods make assumptions about the relative length of the query versus target sequences and the degree to which they overlap. For example, short read alignment methods typically expect the query (i.e., read) to be much shorter than the target (i.e., genome). Similarly, protein search can be accelerated by assuming the query and target sequences have similar lengths, because it is possible to constrain the search space to matches between similar positions in the query and target sequences. This assumption fails to hold if query or target sequences are incomplete or truncated. In contrast, **DECIPHER**'s search functions are intended to find all hits sharing a significant set of (unmasked) k-mers.

As was used here, most benchmarks for homology inference are created to gauge how well a program can identify structurally similar protein domains at very low percent identity (Saripella et al., 2016). Structure-based search methods, such as DeepBLAST (Hamamsy et al., 2023) and Foldseek (van Kempen et al., 2023), outperform sequence-based search methods on these benchmarks. This result has led some to question whether protein BLAST is now obsolete (Al-Fatlawi et al., 2023). However, domain-based benchmarks only represent a fraction of the situations where search is commonly used. Many use cases involve finding full-length proteins in a proteome (Hernandez-Salmeron and Moreno-Hagelsieb, 2020) or short nucleotide sequences in a genome (Marrama et al., 2023). **DECIPHER**'s search functionality was designed to provide dependable search for these common use cases, as well as many others.

The main strength of **DECIPHER**'s search functions is their versatility. The main weakness is also their versatility, because the functions are not tailored for any particular purpose. Specific applications, such as paired-end read mapping, would require customized R code to process the results. For example, **DECIPHER** does not automatically output BAM or SAM files that are common in mapping applications, although similar information is contained within the functions' outputs. In this sense, **DECIPHER** is best suited for R users who have the inclination, time, and ability to develop code around its functions. I anticipate that the new features in **DECIPHER** v3 will further encourage and empower users to perform bioinformatics in the R environment.

## 9 Acknowledgements

This study was funded by the NIAID at the NIH (grant number 1U01AI176418-01). This work used computer resources provided by the Open Science Grid.

## References

- M. I. Abouelhoda and E. Ohlebusch. A local chaining algorithm and its applications in comparative genomics. *Algorithms in Bioinformatics*, pages 1–16. Springer Berlin Heidelberg, 2003. ISBN 978-3-540-39763-2. URL [https://link.springer.com/chapter/10.1007/978-3-540-39763-2\\_1](https://link.springer.com/chapter/10.1007/978-3-540-39763-2_1). [p192]

- A. Al-Fatlawi, M. Menzel, and M. Schroeder. Is protein BLAST a thing of the past? *Nat Commun*, 14(1):8195, 2023. ISSN 2041-1723. doi: 10.1038/s41467-023-44082-5. URL <https://www.ncbi.nlm.nih.gov/pubmed/38081865>. [p197]
- A. M. Altenhoff, B. Boeckmann, S. Capella-Gutiérrez, D. A. Dalquen, T. DeLuca, K. Forslund, J. Huerta-Cepas, B. Linard, C. Pereira, L. P. Pryszcz, F. Schreiber, A. S. da Silva, D. Szklarczyk, C.-M. Train, P. Bork, O. Lecompte, C. von Mering, I. Xenarios, K. Sjölander, L. J. Jensen, M. J. Martin, M. Muffato, c. Quest for Orthologs, T. Gabaldón, S. E. Lewis, P. D. Thomas, E. Sonnhammer, and C. Dessimoz. Standardized benchmarking in the quest for orthologs. *Nature Methods*, 13(5):425–430, 2016. doi: 10.1038/nmeth.3830. URL <https://pubmed.ncbi.nlm.nih.gov/27043882/>. [p196]
- S. F. Altschul, T. L. Madden, A. A. Schäffer, J. Zhang, Z. Zhang, W. Miller, and D. J. Lipman. Gapped BLAST and PSI-BLAST: a new generation of protein database search programs. *Nucleic Acids Research*, 25(17):3389–3402, 1997. doi: 10.1093/nar/25.17.3389. URL <https://pubmed.ncbi.nlm.nih.gov/9254694/>. [p191]
- X. Bi, K. Wang, L. Yang, H. Pan, H. Jiang, Q. Wei, M. Fang, H. Yu, C. Zhu, Y. Cai, Y. He, X. Gan, H. Zeng, D. Yu, Y. Zhu, H. Jiang, Q. Qiu, H. Yang, Y. E. Zhang, W. Wang, M. Zhu, S. He, and G. Zhang. Tracing the genetic footprints of vertebrate landing in non-teleost ray-finned fishes. *Cell*, 184(5):1377–1391 e14, 2021. ISSN 0092-8674. doi: 10.1016/j.cell.2021.01.046. URL <https://www.ncbi.nlm.nih.gov/pubmed/33545088>. [p195]
- B. Buchfink, K. Reuter, and H. G. Drost. Sensitive protein alignments at tree-of-life scale using DIAMOND. *Nat Methods*, 18(4):366–368, 2021. ISSN 1548-7091. doi: 10.1038/s41592-021-01101-x. URL <https://www.ncbi.nlm.nih.gov/pubmed/33828273>. [p191, 193]
- S. R. Eddy. Accelerated profile HMM searches. *PLoS Comput Biol*, 7(10):e1002195, 2011. ISSN 1553-734X. doi: 10.1371/journal.pcbi.1002195. URL <https://www.ncbi.nlm.nih.gov/pubmed/22039361>. [p191]
- N. K. Fox, S. E. Brenner, and J. M. Chandonia. SCOPe: Structural classification of proteins—extended, integrating SCOP and ASTRAL data and classification of new structures. *Nucleic Acids Res*, 42 (Database issue):D304–9, 2014. ISSN 0305-1048. doi: 10.1093/nar/gkt1240. URL <https://www.ncbi.nlm.nih.gov/pubmed/24304899>. [p193]
- R. A. Gonzalez-Pech, T. G. Stephens, and C. X. Chan. Commonly misunderstood parameters of NCBI BLAST and important considerations for users. *Bioinformatics*, 35(15):2697–2698, 2019. ISSN 1367-4803. doi: 10.1093/bioinformatics/bty1018. URL <https://www.ncbi.nlm.nih.gov/pubmed/30541060>. [p191]
- T. Hamamsy, J. T. Morton, R. Blackwell, D. Berenberg, N. Carriero, V. Gligorijevic, C. E. M. Strauss, J. K. Leman, K. Cho, and R. Bonneau. Protein remote homology detection and structural alignment using deep learning. *Nat Biotechnol*, 2023. ISSN 1087-0156. doi: 10.1038/s41587-023-01917-2. URL <https://www.ncbi.nlm.nih.gov/pubmed/37679542>. [p197]
- J. E. Hernandez-Salmeron and G. Moreno-Hagelsieb. Progress in quickly finding orthologs as reciprocal best hits: comparing blast, last, diamond and mmseqs2. *BMC Genomics*, 21(1):741, 2020. ISSN 1471-2164. doi: 10.1186/s12864-020-07132-6. URL <https://www.ncbi.nlm.nih.gov/pubmed/33099302>. [p195, 197]
- W. J. Kent. BLAT—the BLAST-like alignment tool. *Genome research*, 2002. doi: DOI:10.1101/gr.229202. URL <https://pubmed.ncbi.nlm.nih.gov/11932250/>. [p192]
- S. M. Kielbasa, R. Wan, K. Sato, P. Horton, and M. C. Frith. Adaptive seeds tame genomic sequence comparison. *Genome Research*, 21(3):487–93, 2011. ISSN 1088-9051. doi: 10.1101/gr.113985.110. URL <https://www.ncbi.nlm.nih.gov/pubmed/21209072>. [p191]
- Y. Liao, G. K. Smyth, and W. Shi. The R package Rsubread is easier, faster, cheaper and better for alignment and quantification of RNA sequencing reads. *Nucleic Acids Res*, 47(8):e47, 2019. ISSN 0305-1048. doi: 10.1093/nar/gkz114. URL <https://www.ncbi.nlm.nih.gov/pubmed/30783653>. [p191]
- T. L. Madden, B. Busby, and J. Ye. Reply to the paper: Misunderstood parameters of NCBI BLAST impacts the correctness of bioinformatics workflows. *Bioinformatics*, 35(15):2699–2700, 2019. ISSN 1367-4803. doi: 10.1093/bioinformatics/bty1026. URL <https://www.ncbi.nlm.nih.gov/pubmed/30590429>. [p191]
- D. Marrama, W. D. Chronister, L. Westernberg, R. Vita, Z. Kosaloglu-Yalcin, A. Sette, M. Nielsen, J. A. Greenbaum, and B. Peters. PEPMatch: a tool to identify short peptide sequence matches in large sets of proteins. *BMC Bioinformatics*, 24(1):485, 2023. ISSN 1471-2105. doi: 10.1186/s12859-023-05606-4. URL <https://www.ncbi.nlm.nih.gov/pubmed/38110863>. [p197]

- A. K. Nash, T. A. Auchtung, M. C. Wong, D. P. Smith, J. R. Gesell, M. C. Ross, C. J. Stewart, G. A. Metcalf, D. M. Muzny, R. A. Gibbs, N. J. Ajami, and J. F. Petrosino. The gut mycobiome of the Human Microbiome Project healthy cohort. *Microbiome*, 5(1):153, 2017. ISSN 2049-2618 (Electronic) 2049-2618 (Linking). doi: 10.1186/s40168-017-0373-4. URL <https://www.ncbi.nlm.nih.gov/pubmed/29178920>. [p195]
- G. P. Raghava and G. J. Barton. Quantification of the variation in percentage identity for protein sequence alignments. *BMC Bioinformatics*, 7:415, 2006. ISSN 1471-2105 (Electronic) 1471-2105 (Linking). doi: 10.1186/1471-2105-7-415. URL <https://www.ncbi.nlm.nih.gov/pubmed/16984632>. [p195]
- M. Remmert, A. Biegert, A. Hauser, and J. Soding. HHblits: lightning-fast iterative protein sequence searching by HMM-HMM alignment. *Nat Methods*, 9(2):173–5, 2011. ISSN 1548-7091. doi: 10.1038/nmeth.1818. URL <https://www.ncbi.nlm.nih.gov/pubmed/22198341>. [p191]
- K. Sahlin, T. Baudeau, B. Cazaux, and C. Marchet. A survey of mapping algorithms in the long-reads era. *Genome Biol.*, 24(1):133, 2023. ISSN 1474-7596. doi: 10.1186/s13059-023-02972-3. URL <https://www.ncbi.nlm.nih.gov/pubmed/37264447>. [p191]
- H. Salgado, S. Gama-Castro, P. Lara, C. Mejia-Almonte, G. Alarcon-Carranza, A. G. Lopez-Almazo, F. Betancourt-Figueroa, P. Pena-Loredo, S. Alquicira-Hernandez, D. Ledezma-Tejeida, L. Arizmendi-Zagal, F. Mendez-Hernandez, A. K. Diaz-Gomez, E. Ochoa-Praxedis, L. J. Muniz-Rascado, J. S. Garcia-Sotelo, F. A. Flores-Gallegos, L. Gomez, C. Bonavides-Martinez, V. M. Del Moral-Chavez, A. J. Hernandez-Alvarez, A. Santos-Zavaleta, S. Capella-Gutierrez, J. L. Gelpi, and J. Collado-Vides. RegulonDB v12.0: a comprehensive resource of transcriptional regulation in *E. coli* K-12. *Nucleic Acids Res.*, 52(D1):D255–D264, 2024. ISSN 1362-4962 (Electronic) 0305-1048 (Print) 0305-1048 (Linking). doi: 10.1093/nar/gkad1072. URL <https://www.ncbi.nlm.nih.gov/pubmed/37971353>. [p194]
- G. V. Saripella, E. L. Sonnhammer, and K. Forslund. Benchmarking the next generation of homology inference tools. *Bioinformatics*, 32(17):2636–41, 2016. ISSN 1367-4803. doi: 10.1093/bioinformatics/btw305. URL <https://www.ncbi.nlm.nih.gov/pubmed/27256311>. [p197]
- C. L. Schoch, B. Robbertse, V. Robert, D. Vu, G. Cardinali, L. Irinyi, W. Meyer, R. H. Nilsson, K. Hughes, A. N. Miller, P. M. Kirk, K. Abarenkov, M. C. Aime, H. A. Ariyawansa, M. Bidartondo, T. Boekhout, B. Buyck, Q. Cai, J. Chen, A. Crespo, P. W. Crous, U. Damm, Z. W. De Beer, B. T. Dentinger, P. K. Divakar, M. Duenas, N. Feau, K. Fliegerova, M. A. Garcia, Z. W. Ge, G. W. Griffith, J. Z. Groenewald, M. Groenewald, M. Grube, M. Gryzenhout, C. Gueidan, L. Guo, S. Hambleton, R. Hamelin, K. Hansen, V. Hofstetter, S. B. Hong, J. Houbraken, K. D. Hyde, P. Inderbitzin, P. R. Johnston, S. C. Karunarathna, U. Koljalg, G. M. Kovacs, E. Kraichak, K. Krizsan, C. P. Kurtzman, K. H. Larsson, S. Leavitt, P. M. Letcher, K. Liimatainen, J. K. Liu, D. J. Lodge, J. J. Luangsa-ard, H. T. Lumbsch, S. S. Maharachchikumbura, D. Manamgoda, M. P. Martin, A. M. Minnis, J. M. Moncalvo, G. Mule, K. K. Nakasone, T. Niskanen, I. Olariaga, T. Papp, T. Petkovits, R. Pino-Bodas, M. J. Powell, H. A. Raja, D. Redecker, J. M. Sarmiento-Ramirez, K. A. Seifert, B. Shrestha, S. Stenroos, B. Stielow, S. O. Suh, K. Tanaka, L. Tedersoo, M. T. Telleria, D. Udayanga, W. A. Untereiner, J. Dieguez Uribeondo, K. V. Subbarao, C. Vagvolgyi, C. Visagie, K. Voigt, D. M. Walker, B. S. Weir, M. Weiss, N. N. Wijayawardene, M. J. Wingfield, J. P. Xu, Z. L. Yang, N. Zhang, W. Y. Zhuang, and S. Federhen. Finding needles in haystacks: linking scientific names, reference specimens and molecular data for Fungi. *Database (Oxford)*, 2014, 2014. ISSN 1758-0463 (Electronic) 1758-0463 (Linking). doi: 10.1093/database/bau061. URL <https://www.ncbi.nlm.nih.gov/pubmed/24980130>. [p195]
- N. Shah, M. G. Nute, T. Warnow, and M. Pop. Misunderstood parameter of NCBI BLAST impacts the correctness of bioinformatics workflows. *Bioinformatics*, 35(9):1613–1614, 2019. ISSN 1367-4803. doi: 10.1093/bioinformatics/bty833. URL <https://www.ncbi.nlm.nih.gov/pubmed/30247621>. [p191]
- M. Steinegger and J. Soding. MMseqs2 enables sensitive protein sequence searching for the analysis of massive data sets. *Nat Biotechnol.*, 35(11):1026–1028, 2017. ISSN 1087-0156. doi: 10.1038/nbt.3988. URL <https://www.ncbi.nlm.nih.gov/pubmed/29035372>. [p191, 193]
- M. van Kempen, S. S. Kim, C. Tumescheit, M. Mirdita, J. Lee, C. L. M. Gilchrist, J. Soding, and M. Steinegger. Fast and accurate protein structure search with Foldseek. *Nat Biotechnol.*, 2023. ISSN 1087-0156. doi: 10.1038/s41587-023-01773-0. URL <https://www.ncbi.nlm.nih.gov/pubmed/37156916>. [p197]
- E. S. Wright. Using DECIIPHER v2.0 to analyze big biological sequence data in R. *The R Journal*, 8(1):352–359, 2016. doi: 10.32614/RJ-2016-025. URL <https://doi.org/10.32614/RJ-2016-025>. [p191]

E. Wygoda, G. Loewenthal, A. Moshe, M. Alburquerque, I. Mayrose, and T. Pupko. Statistical framework to determine indel-length distribution. *Bioinformatics*, 40(2), 2024. ISSN 1367-4803. doi: 10.1093/bioinformatics/btae043. URL <https://www.ncbi.nlm.nih.gov/pubmed/38269647>. [p192]

*Erik S. Wright*

*University of Pittsburgh*

*Department of Biomedical Informatics*

*Pittsburgh, PA, USA*

*ORCID: 0000-0002-1457-4019*

*eswright@pitt.edu*

# Bioconductor Notes, June 2024

by Maria Doyle, Bioconductor Community Manager, and Bioconductor Core Developer Team

**Abstract** We discuss general project news.

## 1 Introduction

Bioconductor provides tools for the analysis and comprehension of high-throughput genomic data. The project has entered its twentieth year, with funding for core development and infrastructure maintenance secured through 2025 (NIH NHGRI 2U24HG004059). Additional support is provided by NIH NCI, Chan-Zuckerberg Initiative, National Science Foundation, Microsoft, and Amazon. In this news report, we give some updates on core team and project activities.

## 2 Software

In May 2024, Bioconductor 3.19 was released\*. It is compatible with R 4.4 and includes 2289 software packages, 431 experiment data packages, 928 up-to-date annotation packages, 30 workflows, and 5 books. Books are built regularly from source, ensuring full reproducibility; an example is the community-developed [Orchestrating Single-Cell Analysis with Bioconductor](#).

\*Note: Bioconductor 3.20 was subsequently released in October 2024. For details on the latest release, visit the [Bioconductor website](#).

## 3 Community and Impact

### 3.1 CZI EOSS6 Grants Announcement

In June 2024, the Chan Zuckerberg Initiative (CZI) announced funding for five Bioconductor projects as part of the EOSS Cycle 6. These projects focus on advancing training accessibility, developer support, GPU computing, resource tagging with ontologies, and the expansion of Tidyomics tools. For more details on each project and its collaborators, see the blog post [here](#).

### 3.2 CZI Open Science Meeting Highlights

On June 12–14, 2024, Bioconductor members attended the Chan Zuckerberg Initiative (CZI) Open Science Meeting in Boston. This gathering provided an opportunity for current CZI Open Science grantees to connect and share their progress. Bioconductor was recognised as one of the top open-source software projects in the CZ Software Mentions dataset. For more details, see the blog post [here](#).

### 3.3 CSAMA 2024: Biological Data Science Summer School

The 20th edition of [CSAMA \(Computational Statistics for Biological Data Analysis\)](#) Summer School took place in Bressanone-Brixen, Italy, from June 23–28, 2024. This intensive course covered advanced analysis of molecular data using the R/Bioconductor ecosystem. Faculty included Robert Gentleman, co-founder of R and Bioconductor, and Vince Carey, current Bioconductor lead, highlighting Bioconductor's role in open-source training.

## 4 Conferences

### 4.1 BioC2024 Reminder

The annual Bioconductor Conference (BioC2024) is fast approaching! This year's event will be held in Grand Rapids, Michigan, from July 24–26, 2024. Featuring keynote talks, interactive workshops, and ample opportunities for community engagement, BioC2024 is a must-attend event for Bioconductor users and developers alike. Visit the [conference website](#) to register and learn more.

## 4.2 EuroBioC2024 Reminder

Planning is well underway for EuroBioC2024, the European Bioconductor Conference, scheduled for September 4–6, 2024, in Oxford, UK. Like BioC2024, EuroBioC2024 will include keynote presentations, hands-on workshops, and plenty of opportunities for networking with members of the Bioconductor community. For details and updates, visit the [conference website](#).

## 5 Using Bioconductor

Start using Bioconductor by installing the most recent version of R and evaluating the commands

```
if (!requireNamespace("BiocManager", quietly = TRUE))
  install.packages("BiocManager")
BiocManager::install()
```

Install additional packages and dependencies, e.g., [SingleCellExperiment](#), with

```
BiocManager::install("SingleCellExperiment")
```

[Docker](#) images provides a very effective on-ramp for power users to rapidly obtain access to standardized and scalable computing environments. Key resources include:

- [bioconductor.org](#) to install, learn, use, and develop Bioconductor packages.
- A list of [available software](#) linking to pages describing each package.
- A question-and-answer style [user support site](#) and developer-oriented [mailing list](#).
- A community slack workspace ([sign up](#)) for extended technical discussion.
- The [F1000Research Bioconductor gateway](#) for peer-reviewed Bioconductor workflows as well as conference contributions.
- The [Bioconductor YouTube](#) channel includes recordings of keynote and talks from recent conferences, in addition to video recordings of training courses.
- Our [package submission](#) repository for open technical review of new packages.

Upcoming and recently completed events are browsable at our [events page](#).

The [Technical](#) and [Community](#) Advisory Boards provide guidance to ensure that the project addresses leading-edge biological problems with advanced technical approaches, and adopts practices (such as a project-wide [Code of Conduct](#)) that encourages all to participate. We look forward to welcoming you!

We welcome your feedback on these updates and invite you to connect with us through the [Bioconductor Slack](#) workspace or by emailing [community@bioconductor.org](mailto:community@bioconductor.org).

*Maria Doyle, Bioconductor Community Manager  
University of Limerick*

*Bioconductor Core Developer Team  
Dana-Farber Cancer Institute, Roswell Park Comprehensive Cancer Center, City University of New York, Fred Hutchinson Cancer Research Center, Mass General Brigham*

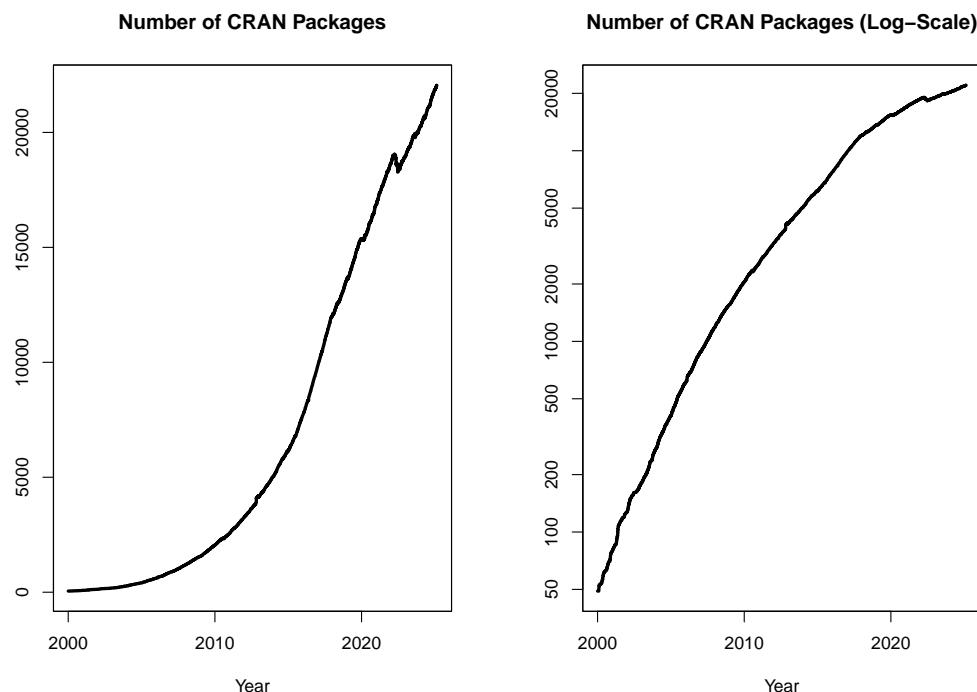
# Changes on CRAN

2024-07-01 to 2024-09-30

by Kurt Hornik, Uwe Ligges, and Achim Zeileis

## 1 CRAN growth

In the past 3 months, 473 new packages were added to the CRAN package repository. 144 packages were unarchived, 207 were archived and 1 had to be removed. The following shows the growth of the number of active packages in the CRAN package repository:



On 2024-09-30, the number of active packages was around 21414.

## 2 Changes in the CRAN Repository Policy

Similar to providing ORCID identifiers for individual authors and contributors in R packages, the policy now encourages to provide ROR identifiers (from the Research Organization Registry, <https://ror.org/>) in a package's DESCRIPTION file in case there are organizations among the authors and contributors.

The policy now points out explicitly that package names on CRAN are persistent and in general it is not permitted to change a package's name.

## 3 CRAN package submissions

From July 2024 to September 2024 CRAN received 6705 package submissions. For these, 10875 actions took place of which 7759 (71%) were auto processed actions and 3116 (29%) manual actions.

Minus some special cases, a summary of the auto-processed and manually triggered actions follows:

	archive	inspect	newbies	pending	pretest	publish	recheck	waiting
auto	2184	724	1459	219	0	2001	736	364
manual	1255	22	24	49	63	1282	351	58

These include the final decisions for the submissions which were

	archive	publish
auto	2046 (31.2%)	1717 (26.2%)
manual	1235 (18.8%)	1564 (23.8%)

where we only count those as *auto* processed whose publication or rejection happened automatically in all steps.

## 4 CRAN mirror security

Currently, there are 94 official CRAN mirrors, 73 of which provide both secure downloads via ‘`https`’ and use secure mirroring from the CRAN master (via `rsync` through ssh tunnels). Since the R 3.4.0 release, `chooseCRANmirror()` offers these mirrors in preference to the others which are not fully secured (yet).

## 5 CRAN Task View Initiative

There is one new task view:

- **Dynamic Visualizations and Interactive Graphics**: Maintained by Sherry Zhang, Diane Cook, Ian Lytle.

Currently there are 46 task views (see <https://CRAN.R-project.org/web/views/>), with median and mean numbers of CRAN packages covered 110 and 123, respectively. Overall, these task views cover 4741 CRAN packages, which is about 21% of all active CRAN packages.

*Kurt Hornik*  
WU Wirtschaftsuniversität Wien  
Austria  
ORCID: 0000-0003-4198-9911  
[Kurt.Hornik@R-project.org](mailto:Kurt.Hornik@R-project.org)

*Uwe Ligges*  
TU Dortmund  
Germany  
ORCID: 0000-0001-5875-6167  
[Uwe.Ligges@R-project.org](mailto:Uwe.Ligges@R-project.org)

*Achim Zeileis*  
Universität Innsbruck  
Austria  
ORCID: 0000-0003-0918-3766  
[Achim.Zeileis@R-project.org](mailto:Achim.Zeileis@R-project.org)

# Changes in R

by Tomas Kalibera and Sebastian Meyer

**Abstract** We present selected changes in the development version of R, which is referred to as R-devel and is to become R 4.5. We also provide some statistics on bug tracking activities in 2024, which covers 4 months before the release of R 4.4 and 8 months of work on R-devel.

## 1 Selected changes in R-devel

R 4.5.0 is due to be released around April 2025. The following gives a selection of changes in R-devel, which are likely to appear in the new release.

For a more complete list of changes, run `news(is.na(Version))` in R-devel or inspect the nightly rendered version of the NEWS.Rd source file at <https://CRAN.R-project.org/doc/manuals/r-devel/NEWS.html> or its RSS feed at <https://developer.R-project.org/RSSfeeds.html>. The bundled *R News* give concise descriptions of a large number of changes, including many smaller ones.

### 1.1 Selected user-visible changes

- R packages for installation are now downloaded in parallel, which improves the download speed, on some systems by a large factor.

Existing support for simultaneous download in `download.files()` has been improved and functions `install.packages()` and `download.packages()` now use it with the default download method (`libcurl`). Downloading a single file requires a number of handshakes between the client and the server, which may be expensive especially on connections with large latencies. With simultaneous download, existing connections to the same server may be re-used and handshakes be made in parallel with each other and with other transfers. The actual speedup depends on network characteristics, but significant speedups have been observed both with some very good as well as with some rather poor network connections.

The speedups are bigger when HTTP version 2 is used, which depends on the server (many CRAN mirrors support it) and when the curl library used on the client supports it (on Unix and macOS it typically does, on Windows the upcoming Rtools45 will support it).

More details about this work can be found in a blog post at <https://blog.r-project.org/2024/12/02/faster-downloads>. The previous implementation of simultaneous download has been improved to be more careful about available resources, particularly available file handles.

The implementation of R internet timeout has been improved to allow parallel downloads of many files, though there is a semantic gap between an absolute-time timeout in R and simultaneous download, where transfers may be delayed by ongoing transfers and various connection assets may be re-used between transfers.

The status reporting for simultaneous transfers in `download.file()` has been improved so that the caller can easily find out which individual transfers have succeeded. The simultaneous download can be used directly also for other files, not only R packages.

- One can now enter and edit arbitrarily long input lines in the R console on Linux and macOS (when the Readline library is used, which is normally the case) and in Rterm on Windows. Arbitrarily long input lines can also be used in Rgui on Windows, but only the last segment can be edited. In previous versions of R, there was a hard-coded limit of about 4K bytes, which worked fine for input entered manually, but some users ran into problems when pasting generated code. On some systems the input was truncated, with Rgui on Windows it could also be corrupted.

More details about this work can be found in a blog post at <https://blog.r-project.org/2024/08/30/long-input-lines>. The R parser itself is always invoked on a fixed-length prefix of the input. When the prefix turns out too short to make any progress parsing, the parser is re-started with a longer prefix and the length is potentially unlimited. This iterative mechanism is exposed in Rgui on Windows, when the user can enter additional segments of the line, but then cannot edit the previous segment anymore. Still, several bugs leading to input corruption were found during this work and fixed.

Readline, used on Linux and macOS in R, allows to edit lines of arbitrary length, so an intermediate layer has been implemented which provides the input to the parser using that iterative mechanism. Rterm on Windows uses a rewrite of getline library, which has been extended to support editing lines of arbitrary length, and again an intermediate layer has been implemented to provide that to the parser.

- Support for the SHA-256 hashing algorithm has been added to R's **tools** package via function `sha256sum()`. It allows hashing files on disk and raw vectors of bytes in memory. The underlying implementation used is public domain code written by Ulrich Drepper. SHA-256 from the SHA-2 family of hash functions is considered more secure than MD5 (`md5sum()`) and hence is more appropriate whenever the aim is to detect not only accidental data corruption, but also malicious modification. SHA-256 is generally slower to compute than MD5 and the checksum is twice as long (MD5 is 128-bit). `md5sum()` has been extended to also support computation of a hash of raw vector of bytes in memory, so that it has the same interface as `sha256sum()`.
- R now supports zstd compression. Function `zstdfile()` has been added which allows to create R connections to read from and write to zstd-compressed files. One can use zstd compression with serialization and there is also some support for zstd compression in the `tar()` and `untar()` functions. The zstd compression support is currently optional in R builds: it is only included when the zstd library is available at build time on Unix. On Windows, the zstd library is part of Rtools so zstd support is always available.

With the default tuning options used currently in R, zstd offers typically slightly worse compression than xz, but is faster. The tradeoffs, however, become different with other tuning options.

- There is a new wrapper function `grepv(...)`, short hand for `grep(..., value = TRUE)`, to return the matching elements themselves rather than their indices.
- A PDF document may include metadata in the so-called *document information dictionary*. For instance, PDF documents produced by R's `pdf()` device always set "R" as Creator and by default use "R Graphics Output" as Title. An additional argument `author` has now been added to set the Author entry (omitted by default), possibly via `pdf.options()` in an R profile or init script. Furthermore, the new logical arguments `timestamp` (setting `CreationDate` and `ModDate`) and `producer` ("R <major>.<minor>") can be set to FALSE to disable the corresponding metadata fields. Disabling automatic timestamps simplifies tracking uncompressed PDF output in version control systems.
- R CMD check gains an option `--run-demo` to check the R scripts in the demo directory similar to those in tests. This means demos are run, potentially compared to reference output in corresponding `.Rout.save` files, and it is checked that required R packages are declared in the `DESCRIPTION` file. Whereas package tests are not installed by default, demos are and are thus exposed to users, so package maintainers should ensure they work. The new check option provides an alternative to including `demos()` calls in examples or test scripts. Demos may be interactive (e.g., some await a browser response, require authentication, or contain instructions for special data or software setups) and are sometimes used for computationally intensive examples (e.g., replication scripts from associated publications), so R CMD check does not run demo scripts by default

(similar to `\donttest` examples, which are only checked with option `--run-donttest`). Package maintainers could add an explicit condition such as `if (!interactive()) quit()` to demo scripts that require interaction, so (occasional) checks with `--run-demo` can be used to cover the others.

At the time of writing, 603 (3%) of 22133 packages on CRAN contained a demo directory. Of these packages, 117 used undeclared packages in demo scripts. Experiments showed that dozens indeed required an interactive user or a special setup, and others did not complete within 90 minutes (the chosen timeout on the test system). Of the remaining 504 packages, 346 (69%) produced noteworthy errors: a considerable amount of 275 packages failed a demo with “could not find function” or “object not found” errors, often simply because the script forgot to attach the own package, and sometimes because functionality has been removed in the meantime. Remaining issues include coding errors that are caught in recent versions of R (e.g., `length > 1` conditions), dead URLs, or breakage that is likely caused by changes in dependencies.

## 1.2 Selected low-level changes

A number of changes in R-devel — perhaps more than usual — are low-level and invisible to R users, but important for the reliability and maintainability of R and R packages now and in the future. Some changes of this kind are listed here.

- Progress has been made on improving the C API for R packages and embedding applications. The interface has evolved organically over the years, sometimes exposing more than necessary or helpful. With some functions it is too easy to make errors leading to segfaults or incorrect results. Some functions expose too much internal structure or behavior, which needs to be changed occasionally to be able to maintain and improve R itself.

The interface is defined in the ‘Writing R Extensions’ manual, see `RShowDoc("R-exts")` in your version of R or online at <https://CRAN.R-project.org/manuals.html>: what is mentioned there and is available in installed header files (in directory `R.home("include")`) is in the API. Most internal functions not in the API are “hidden” in the dynamic libraries, which prevents their accidental use. In R-devel they are now hidden also on macOS, previously only on Windows and Linux.

Some non-API functions currently have to remain unhidden, because they are needed by base R packages that are part of the R distribution and are maintained with R. The cleanup involved hiding some functions that were unnecessarily exposed, replacing some functions with safer alternatives, and adding some to the API by documenting them. This required cooperation of package authors and was tested and tuned on CRAN and Bioconductor packages.

The ‘Writing R Extensions’ manual has been improved so that functions/symbols in the API are explicitly flagged to be in the (regular) API, Fortran API (for Fortran code), experimental API (subject to change, so users must be prepared to adapt much more than with other categories), or embedding API (only for applications embedding R, including front-ends). There is now internal functionality in R that can query if a given symbol is in the API, which allows stricter checking.

- R includes a number of workarounds for bugs in the `iconv` library shipped with recent versions of macOS. The library is used for encoding conversions, e.g., when converting between Latin-1 and UTF-8. Despite most R installations today use UTF-8 as the native encoding, such conversions are still happening, e.g., for plot labels and when importing historical data from legacy file formats.

The `iconv` implementation that appeared in macOS 14.0 changed the behavior with characters not representable in the target encoding, but also caused crashes with legitimate use and caused incorrect conversions (garbage in output), which has been

detected by R developers and users. The bugs lead to incorrect behavior with conversion state reset, with incremental conversion of an input stream, and with treating byte-order marks. The details can be found in a blog post at <https://blog.r-project.org/2024/12/11/problems-with-iconv-on-macos>.

The current CRAN builds of R for macOS use a static build of the iconv library that existed in earlier versions of macOS (still truly GNU libiconv 1.11), so are not yet exposed to these problems. However, users building R from source on macOS would run into them if not applying workarounds. R 4.4 already included some of these workarounds, but R-devel includes more and the previous ones have been improved. Also, R-devel fixes a problem of interaction between one of the older workarounds and an iconv bug causing a problem even in the CRAN builds.

As iconv would usually be linked dynamically when building R from source, the workarounds are based on runtime tests that are executed on first use of encoding conversion functions. Hopefully these bugs will eventually be fixed upstream in macOS, so that these workarounds can be removed. They represent a surprisingly large amount of code.

- The detection of invalid memory accesses caused by attempts to access elements of empty vectors has been enabled by default and improved to be more robust to different inlining decisions by the compiler. Such an access now causes an immediate crash of R and can be easily located using a debugger. This is achieved via “poisoning”: a data pointer to an empty R vector is intentionally returned invalid.

Invalid memory accesses can otherwise lead to crashes later in the computation or to incorrect results. As a result of this change, a large number of memory access bugs have been found in CRAN packages and reported to CRAN maintainers. If packages from other repositories, where regular checking against R-devel is not in place, start crashing with R 4.5, this is one of the possible causes.

The amount of related changes in R-devel is small, but the effort debugging packages and providing fixes or advice to package maintainers has been significant.

- R-devel on Windows has been updated so that it can be built also with some alternative MsyS2 toolchains, not only with Rtools. This required changes in the make files and is based on `pkg-config`, which is used for figuring out library dependencies and C preprocessor options. It simplifies testing of the Windows-specific code in R with newer compilers (newer than currently in Rtools) and it allows using sanitizers on this code. Also, newer compilers with better diagnostics may find problems relevant also for the current compilers in Rtools. R-devel has been tested on Windows with some of the MsyS2 toolchains and sanitizers. In previous years, some alternative toolchains have also been tested, but it required ad-hoc manual modification of the make files in the R sources.

This support can also be used for testing package code with alternative toolchains, which however requires some Windows-specific knowledge. It is essential that all code in use is built by the same toolchain, so for these experiments, one needs to start with an empty library. Also, MsyS2 uses dynamic linking, so one can only run (test) the result with the corresponding MsyS2 toolchain environment. More details are available in a blog at <https://blog.r-project.org/2025/01/28/alternative-toolchains-on-windows>.

- C23 is a new C standard that brings several new features to the language, including `bool`, `true` and `false` keywords. It is the default standard used by GCC 15, which is expected to be released as stable around the release of R 4.5. R-devel, as well as the upcoming third patch release of R 4.4, has been made installable with the current pre-release version of GCC 15 and also with older compilers when the C23 standard is selected explicitly.

R packages affected by this change may either choose to explicitly require some older C standard (possibly C17), or better be updated to work also with C23. R-devel has

been switched to use C23 always when the compiler supports it, even when it is not the default standard of the compiler. This change has been tested on CRAN and Bioconductor packages and fixes were provided to the maintainers.

In addition to the general need of updating code to newer standards so that it can be maintained in the future (say, when other components, e.g., libraries, would use newer standards), some of the C23 features may be useful in R, such as the `bool` type. See ‘Writing R Extensions’, Section 1.2.5 ([online version](#)), for more details on the use of C standards with R.

- Further progress has been made on getting rid of symbol remapping, i.e., automated renaming of C symbols like `allocVector()` to `Rf_allocVector()` done by R headers. The remapping causes problems when it accidentally modifies code included via C preprocessor headers. The risks can be minimized by including headers in a certain order, but it is an additional constraint package authors need to think about, and in practice, the conflicts are a common cause of package compilation errors arising after updates of the toolchain, OS, or a library.

The remapping in R-devel has been disabled for C++ code in packages, because most of the conflicts happened to be found with C++ code. The remapping still remains available and is done by default for C code, but it is recommended not to use the remapping in new C code. The ‘Writing R Extensions’ manual has been updated to use full names of the functions including the prefixes; see Section 6 (“The R API”, [online version](#)), for more details on the re-mapping situation in R-devel.

Similarly, “strict R headers” are now the default, which means that legacy definitions of `Calloc`, `Realloc`, `Free` and `PI` have been removed. One can still use these allocation functions with the `R_` prefix, and `M_PI` is a standard replacement for `PI`.

Like with other similar tasks, this required only relatively small changes in R itself, but most of the effort has been spent on testing with CRAN and Bioconductor packages and helping the package maintainers with the necessary updates.

## 2 Bug statistics for 2024

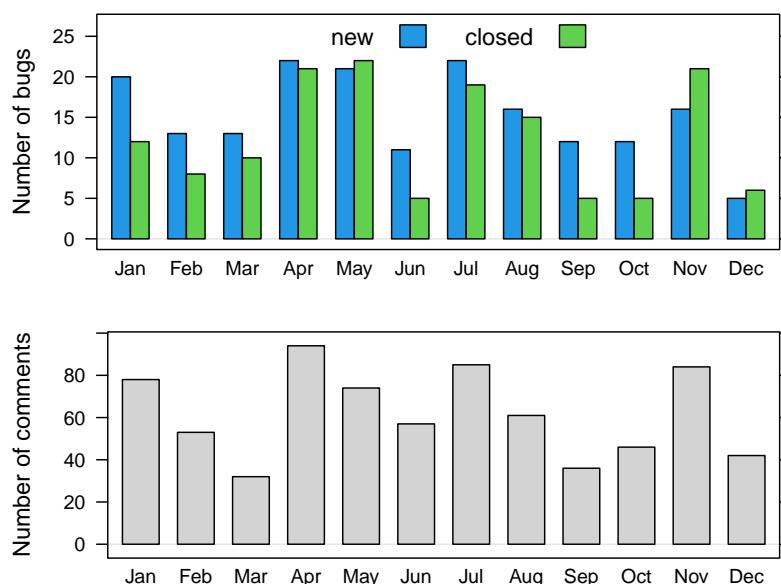
Summaries of bug-related activities over the past year were derived from the database underlying [R’s Bugzilla system](#). Overall, 183 new bugs or feature requests (15%) were submitted in 2024, very much like in the previous two years. The top 5 selected components for these reports were “Documentation”, “Low-level”, “Misc”, “Language”, and “Graphics”, where the first and last increased their rank compared to 2023.

A total of 107 contributors added 742 comments on 250 different reports; 149 reports were closed. Whereas the number of new submissions remained at the same level, open issues were handled at a slower rate compared to 2023, with a decrease in the numbers of closures (−27%), comments (−21%), and contributors (−11%). Compared to 2022, there is also a decrease in bug closures and comments, but it is smaller (−13% closures, −15% comments, −13% contributors).

The numbers of comments and different contributors in the bug tracking system may be influenced by ‘R Dev Days’ organized by the [R Contribution Working Group](#): the bugs are pre-discussed at those events mostly outside the core system, which later receives summarized comments. An ‘R Dev Day’ organized just after UseR 2024 had 8 R Core Team members present, so some “comments” were exchanged in person.

The monthly numbers shown in Figure 1 suggest that more issues were closed in April/May, July/August and November. April/May was around the release of R 4.4, but also one of the ‘R Dev Days’ took place. July/August covers the UseR! 2024 conference with an attached R Core Team meeting and ‘R Dev Day’.

As every year, not all bug reports go through R Bugzilla. Some are picked up from the R-devel, R-pkg-devel or R-help mailing lists. Some come via private communication by e-mail or in person. A number of bugs are always discovered by the R Core Team when



**Figure 1:** Bug tracking activity by month in 2024.

testing changes in R or packages; these bugs are usually not reported by any channel, but are fixed directly.

*Tomas Kalibera  
R Core Team  
Prague, Czechia  
ORCID: 0000-0002-7435-734X  
Tomas.Kalibera@R-project.org*

*Sebastian Meyer  
Friedrich-Alexander-Universität Erlangen-Nürnberg  
Erlangen, Germany  
ORCID: 0000-0002-1791-9449  
Sebastian.Meyer@R-project.org*

# R Foundation News

by *Torsten Hothorn*

## 1 Donations and members

Membership fees and donations received between 2024-12-06 and 2025-01-29.

### 1.1 Donations

James Curran (New Zealand)  
Joseph Luchman (United States)  
Riccardo Martinelli (Italy)  
The University of Auckland, Statistics Department (New Zealand)  
Fergus Reig Gracia (Spain)  
Thomas Stadler (Switzerland)

### 1.2 Supporting benefactors

b-data GmbH, Winterthur (Switzerland)

### 1.3 Supporting institutions

oikostat GmbH, Ettiswil (Switzerland)

### 1.4 Supporting members

Richard Abdill (United States)  
Gilberto Camara (Brazil)  
Ivan Maria Castellani (Italy)  
Keith Chamberlain (United States)  
Cédric Chambru (Switzerland)  
Michael Chirico (United States)  
Dan Dediu (Spain)  
Kevin DeMaio (United States)  
Anna Doizy (Réunion)  
Fraser Edwards (United Kingdom)  
Isaac Florence (United Kingdom)  
Neil Frazer (United States)  
Sven Garbade (Germany)  
Jan Marvin Garbuszus (Germany)

Eduardo García Galea (Spain)  
Brian Gramberg (Netherlands)  
Spencer Graves (United States)  
Krushi Gurudu (United States)  
Frank Hafner (United States)  
Joe Harwood (United Kingdom)  
Philippe Heymans Smith (Costa Rica)  
Adam Hill (United States)  
JUNE KEE KIM (Korea, Republic of)  
Caleb Lareau (United States)  
Thomas Levine (United States)  
Baoxiao Liu (Netherlands)  
Keon-Woong Moon (Korea, Republic of)  
yoshinobu nakahashi (Japan)  
Antonio Paez (Canada)  
Sermet Pekin (Turkey)  
Ingo Ruczinski (United States)  
Choonghyun Ryu (Korea, Republic of)  
Kai Streicher (Switzerland)  
Petr Waldauf (Czechia)  
Fredrik Wartenberg (Sweden)  
Agnieszka Zawiejska (Poland)  
Vaidotas Zemlys-Balevičius (Lithuania)  
广宇曾 (China)

*Torsten Hothorn*  
*Universität Zürich*  
*Switzerland*  
*ORCID: 0000-0001-8301-0471*  
*Torsten.Hothorn@R-project.org*