

# Package ‘mlm4omics’

June 14, 2019

**Version** 1.2.0

**Title** Multilevel Model for Multivariate Responses with Missing Values

**Description** To conduct Bayesian inference regression for responses with multilevel explanatory variables and missing values; It uses function from 'Stan', a software to implement posterior sampling using Hamiltonian MC and its variation Non-U-Turn algorithms. It implements the posterior sampling of regression coefficients from the multilevel regression models.

The package has two main functions to handle not-missing-at-random missing responses and left-censored with not-missing-at random responses.

The purpose is to provide a similar format as the other R regression functions but using 'Stan' models.

**Maintainer** Irene SL Zeng <i.zeng@auckland.ac.nz>

**License** GPL-3

**Encoding** UTF-8

**LazyData** true

**ByteCompile** true

**Depends** R (>= 3.5.0), Rcpp (>= 0.12.17), methods, stats

**Imports** rstan (>= 2.17.3),rstantools (>= 1.5.0),MASS,Matrix,stats4,ggplot2

**LinkingTo** StanHeaders (>= 2.17.2), rstan (>= 2.17.3), BH (>= 1.66.0-1), Rcpp (>= 0.12.17), RcppEigen (>= 0.3.3.4.0)

**Suggests** testthat, BiocStyle, knitr, rmarkdown, roxygen2 (>= 5.0.0)

**URL** <https://doi.org/10.1101/153049>

**RoxygenNote** 6.0.1

**biocViews** ImmunoOncology,  
Bayesian,CopyNumberVariation,Classification,Regression,MassSpectrometry,Proteomics,Software

**bugReport** <https://github.com/ireneslzeng/mlmm/issues>

**VignetteBuilder** knitr

**SystemRequirements** GNU make

**NeedsCompilation** yes

**git\_url** <https://git.bioconductor.org/packages/mlm4omics>

**git\_branch** RELEASE\_3\_9

**git\_last\_commit** 8c2d5ac

**git\_last\_commit\_date** 2019-05-02

**Date/Publication** 2019-06-13

**Author** Irene Zeng [aut, cre],  
Thomas Lumley [ctb]

## R topics documented:

mlm4omics-package . . . . .	2
mlmc . . . . .	2
mlmm . . . . .	4
pdata . . . . .	6
setinitvalues . . . . .	7
<b>Index</b>	<b>9</b>

---

mlm4omics-package	<i>The 'mlm4omics' package.</i>
-------------------	---------------------------------

---

### Description

To conduct Bayesian inference regression for responses with multilevel explanatory variables and missing values

### References

Stan Development Team (2018). RStan: the R interface to Stan. R package version 2.17.3. <http://mc-stan.org>

---

mlmc	<i>The multilevel function for missing and censored dependents: mlmc().</i>
------	---

---

### Description

mlmc() handles Bayesian multilevel model with response variable that has left-censored values, and missing values that depends on the response value itself. Apart from the response value, the missingness is also known to associate with the other variables. The method is created for analyzing mass-spectrometry data when it has abundance-dependant missing and censored values, and there are prior information available for the associations between the probability of missing and the known variables. The imputed values for the censored response are outputted as part of the parameters.

### Usage

```
mlmc(formula_completed, formula_missing, formula_censor = NULL,
      formula_subject, pdata, respond_dep_missing = TRUE,
      response_censorlim = NULL, pidname, sidname, prec_prior = NULL,
      alpha_prior = NULL, iterno = 100, chains = 3, thin = 1, seed = 125,
      algorithm = "NUTS", warmup = floor(iterno/2), adapt_delta_value = 0.9,
      savefile = FALSE)
```

**Arguments**

formula_completed	The main regression model formula; It has the same formula format as lmr() and it is used to define the first level response and its explanatory variables.
formula_missing	The logistic regression model formula; It has the same formula as formula_completed.
formula_censor	The formula used in the program to define the observations with censored values.
formula_subject	The second level formula in the multilevel model which is used to define responses such as subject and its explanatory variables.
pdata	The dataset contains response and predictors in a long format. Response is a vector with an indicator variable to define the corresponding unit. The data needs to have the following rudimental variables: the indicator variable for first level response, second level indicator variable for subject such as subject id or a sampling unit, an indicator for missingness and indicator of censoring. Missingness and censored are two different classifications, when these two variables are tabulated, there must not have any observation defined as censored and missing. Data structure can be referred from the example and vignette.
respond_dep_missing	A logical variable to indicate whether response value is missing-dependant.
response_censorlim	The detectable limit for the response value, i.e. 1 mg per Liter for intensity value.
pidname	Variable name to define the multilevel response unit, i.e. protein name or gene name.
sidname	Variable name to define the subject unit, i.e. patient id or sampling id.
prec_prior	prior precision matrix of the explanatory variables at the first level unit in the multilevel model, for example, the variables to predict the ion intensity. The dimension will be $q \times q$ , where $q$ is the number of explanatory variables at the right-hand side of formula_completed. The default is a matrix with diagonal values of 0.01 and off-diagonal values of 0.005.
alpha_prior	prior for coefficients of predictors to missing probability in the logistic regression. Its length will be equal to the number of variables at the right-hand side of the formula_missing. Default is a vector of zeros.
iterno	Number of iterations for the posterior samplings.
chains	rstan parameter to define number of chains of posterior samplings.
thin	rstan parameter to define the frequency of iterations saved.
seed	random seed for rstan function.
algorithm	rstan parameter which has three options NUTS, HMC, Fixed param.
warmup	Number of iterations for burn-out in stan.
adapt_delta_value	Adaptive delta value is an adaptation parameters for sampling algorithms, default is 0.85, value between 0-1.
savefile	A logical variable to indicate if the sampling files are to be saved.

**Value**

Return of the function is the result fitted by stan. It will have the summarized parameters from all chains and summary results for each chain.

**Examples**

```
## Not run:
set.seed(150)
library(MASS)
var2 <- abs(rnorm(800,0,1)); treatment <- c(rep(0,400), rep(1,400));
var1 <- (1/0.85)*var2+2*treatment;
geneid <- rep(seq_len(50),16);
sid <- c(rep(seq_len(50),8), rep(seq_len(50)+50,8))
cov1 <- rWishart(1,df=50, Sigma <- diag(rep(1,50)))
u <- rnorm(50,0,1);mu <- mvrnorm(n=1, mu=u, cov1[, ,1])
sdd <- rgamma(1, shape=1, scale=1/10);
for (i in seq_len(800)) {var1[i] <- var1[i]+rnorm(1, mu[geneid[i]], sdd)}
miss_logit <- var2*(-0.9)+var1*(0.001);
miss <- rbinom(800, 1, exp(miss_logit)/(exp(miss_logit)+1));
censor <- rep(0,800)
for (i in seq_len(800)) {if (var1[i]<0.002) censor[i]=1}
pdata <- data.frame(var1, var2, treatment, miss, censor, geneid, sid);
for ( i in seq_len(800))
{if ((pdata$miss[i]==1) & (pdata$censor[i]==1)) pdata$miss[i]=0};
for ( i in seq_len(800)) {
if (pdata$miss[i]==1) pdata$var1[i]=NA;
if (pdata$censor[i]==1) pdata$var1[i]=0.002};
pidname="geneid";sidname="sid";
#copy and paste the following formulas to the mlmm() function respectively
formula_completed=var1~var2+treatment;
formula_missing=miss~var2;
formula_censor=censor~1;
formula_subject=~treatment;
response_censorlim=0.002;

model1 <- mlmc(formula_completed=var1~var2+treatment,
formula_missing=miss~var2,
formula_censor=censor~1,
formula_subject=~treatment,
pdata=pdata,
response_censorlim=0.002,
respond_dep_missing=TRUE,
pidname="geneid",sidname="sid",
iterno=50,
chains=2,
savefile=FALSE)

## End(Not run)
```

## Description

mlmm() handles Bayesian multilevel model with response variable that has missing values that depends on the response value itself. Apart from the response value, the missingness is also known to associate with the other variables. The method is created for analyzing mass-spectrometry data when it has abundance-dependant missing and censored values, and there are prior information available for the associations between the probability of missing and the known variables. The function mlmm is written for response variable has no censored values while mlmc function include imputing censored values.

## Usage

```
mlmm(formula_completed, formula_missing, formula_subject, pdata,
      respond_dep_missing = TRUE, pidname, sidname, prec_prior = NULL,
      alpha_prior = NULL, iterno = 100, chains = 3, thin = 1, seed = 125,
      algorithm = "NUTS", warmup = floor(iterno/2), adapt_delta_value = 0.9,
      savefile = FALSE)
```

## Arguments

formula_completed	The main regression model formula. It has the same formula format as lmr() and it is used to define the first level response and its explanatory variables.
formula_missing	The logistic regression model formula. It has the same formula as formula_completed.
formula_subject	The second level formula in the multilevel model which is used to define second level unit such as subject and explanatory variables.
pdata	The dataset contains response and predictors in a long format. Response is a vector with an indicator variable to define the corresponding unit. The data needs to have the following rudimental variables: the indicator variable for first level response, the indicator variable for second level unit such as subject or a sampling unit, an indicator for missingness and indicator of censoring. Missingness and censored are two different classification, there should not have any overlap between missingness and censored. Data structure can be referenced from the example and reference papers.
respond_dep_missing	An indicator of whether response value is missing-dependant.
pidname	Variable name to define the multilevel response unit, i.e. protein name or gene name
sidname	Vriable name to define the subject unit, i.e. patient id or sampling id
prec_prior	prior precision matrix of the explanatory variables at the first level unit in the multilevel model, for example, the variables to predict the ion intensity. The dimension will be $q \times q$ , where $q$ is the number of explanatory variables at the right-hand side of formula_completed. The default is a matrix with diagonal value of 0.01 and off-diagonal value of 0.005.
alpha_prior	prior for coefficients of predictors to missing probability in the logistic regression. Its length will be equal to the number of variables at the right-hand side of the formula_missing. Default is a vector of zeros.
iterno	Number of iterations for the posterior samplings
chains	rstan() parameter to define number of chains of posterior samplings.

thin	rstan() parameter to define the frequency of iterations saved.
seed	random seed for rstan() function
algorithm	rstan() parameter which has three options c(NUTS,HMC, Fixed_param).
warmup	Number of iterations for burn-out in stan.
adapt_delta_value	Adaptive delta value is an adaptation parameters for sampling algorithms,default is 0.85, value between 0-1.
savefile	A logical variable to indicate if the sampling files are to be saved.

### Value

Return of the function is the result fitted by stan(). It will have the summarized parameters from all chains and summary results for each chain. Plot() function will return the visualization of the mean and parameters.

### Examples

```
library(MASS)
set.seed(150)
var2 <- abs(rnorm(1000,0,1)); treatment <- c(rep(0,500),rep(1,500))
geneid <- rep(seq_len(20),50);
sid <- c(rep(seq_len(25),20),rep(seq_len(25)+25,20))
cov1 <- rWishart(1,df=100,Sigma <- diag(rep(1,100)))
u <- rnorm(100,0,1)
mu <- mvrnorm(n=1,mu=u,cov1[, ,1])
sdd <- rgamma(1,shape=1,scale=1/10)
var1=(1/0.85)*var2+2*treatment
for (i in seq_len(1000)) {var1[i]=var1[i]+rnorm(1,mu[geneid[i]],sdd)}
miss_logit <- var2*(-0.9)+var1*(0.01)
probmiss <- exp(miss_logit)/(exp(miss_logit)+1)
miss <- rbinom(1000,1,probmiss); table(miss)
pdata <- data.frame(var1,var2,treatment,miss,geneid,sid)
for ( i in seq_len(1000)) if (pdata$miss[i]==1) pdata$var1[i]=NA;
pidname="geneid"; sidname="sid";
#copy and paste the following formulas to the mmlm() function respectively
formula_completed=var1~var2+treatment
formula_missing=miss~var2
formula_censor=censor~1
formula_subject=~treatment
model3 <- mmlm(formula_completed=var1~var2+treatment,
formula_missing=miss~var2,
formula_subject=~treatment, pdata=pdata, respond_dep_missing=TRUE,
pidname="geneid", sidname="sid", iterno=10, chains=2,
savefile=FALSE)
```

---

pdata

*pdata for examples and testthat() pdata has 7 variables and var1 is the response variable, var2 is a continuous explanatory variable, treatment is another explanatory variable, miss and censor are indicator for missing and censored, geneid and sid represents gene id and subject id respectively.*

---

**Description**

pdata for examples and testthat() pdata has 7 variables and var1 is the response variable, var2 is a continuous explanatory variable, treatment is another explanatory variable, miss and censor are indicator for missing and censored, geneid and sid represents gene id and subject id respectively.

**Usage**

```
data(pdata)
```

**Format**

An object of class `data.frame` with 100 rows and 7 columns.

---

<code>setinitvalues</code>	<i>The function to set initial values for parameters: <code>setinitvalues()</code>.</i>
----------------------------	---

---

**Description**

Generate initial values for parameters

**Usage**

```
setinitvalues(npred, np, npred_miss, npred_sub, nmiss, nsid,
  censor_lim_upp = 0.008, ita_a = 1, ita_b = 1/10, g_mu = 0,
  g_sig = 1, alpha_mu_u = 0, alpha_mu_s = 1, alpha_theta_a = 1,
  alpha_theta_b = 1/10, beta2_theta_a = 1, beta2_theta_b = 1/10)
```

**Arguments**

<code>npred</code>	number of predictors for the regression model
<code>np</code>	number of protein/metabolite units comprised of the response values (i.e. which represents peptides' ion-intensities used to construct protein/metabolite's abundance)
<code>npred_miss</code>	number of predictors for missingness
<code>npred_sub</code>	number of predictors for the second level units such as subjects
<code>nmiss</code>	number of observations with missing responses values
<code>nsid</code>	number of second level units i.e. subjects
<code>censor_lim_upp</code>	upper-limit of censored value of responses. The default value is 0.001 according to an experiment device. User can change it according to the data.
<code>ita_a</code>	shape parameter for gamma distributed prior-ita (std of response value). The default is set to 1.
<code>ita_b</code>	rate parameter for gamma distributed prior-ita. The default is set to 1/10. The default values of shape and rate parameters provide a reasonable wide range of initial value for ita. Users can change it accordingly.
<code>g_mu</code>	mean of normal distributed location parameter $g$ for re-parameterising $U$ (regression coefficient of model in the completed data). The default value is set to 0 for the mean of standard normal distribution.

<code>g_sig</code>	std of normal distributed location parameter <code>g</code> . The default is set to 1 for the std of normal distribution.
<code>alpha_mu_u</code>	mean of normal distributed location parameter <code>alpha_mu</code> for re-parameterising <code>alpha</code> (regression coefficient of logistic regression model for missing prob). The default is set to 0.
<code>alpha_mu_s</code>	std of normal distributed location parameter <code>alpha_mu</code>
<code>alpha_theta_a</code>	shape parameter of gamma-distributed dispersion parameter <code>alpha_theta</code> for re-parameterising <code>alpha</code> . Default value is set to 1 as a natural starting value.
<code>alpha_theta_b</code>	rate parameter of gamma-distributed <code>alpha_theta</code> . Default value uses 1/10 as for <code>ita_b</code> . Both default values of <code>shape(_a)</code> and <code>rate(_b)</code> of <code>alpha_theta</code> can be changed to give a wider range ( <code>_b=1/10</code> ) or a narrower range ( <code>_b=0.5</code> ).
<code>beta2_theta_a</code>	shape parameter of gamma-distributed <code>beta2_theta</code> for re-parameterizing <code>beta2</code> (regression coefficient for second level units, i.e. subject). Default value uses 1.
<code>beta2_theta_b</code>	rate parameter of gamma distributed dispersion parameter <code>beta2_theta</code> . Default value used 1/10, same as for <code>ita_b</code> .

### Value

`pVAR` precision matrix for predictors in completed data model

`U_latent` standardized multinormal distributed latent variable to re-parameterise regression coefficient `U`.

`g` location parameter to re-parameterise `U`.

`alpha_mu` mean value for `alpha` (regression coefficient of model for missing probability).

`alpha_latent` standardized normal distributed latent variable to re-parameterize `alpha`.

`beta2_latent` standardized multivariate normal distributed latent variable to re-parameterising `beta2`.

`beta2_mu` mean of the multivariate normal distributed `beta2`

`y_m_latent` standardized normal distributed latent variable to re-parameterise response variable.

### Examples

```
testexmp <- setinitvalues(npred=2,np=3,npred_miss=3,npred_sub=2,nmiss=10,
nsid=30)
```



# Index

## \*Topic **datasets**

pdata, [6](#)

mlm4omics (mlm4omics-package), [2](#)

mlm4omics-package, [2](#)

mlmc, [2](#)

mlmm, [4](#)

pdata, [6](#)

setinitvalues, [7](#)