

Package ‘penetrance’

March 3, 2025

Type Package

Title Methods for Penetrance Estimation in Family-Based Studies

Version 0.1.0

Depends R (>= 3.5.0)

Imports clipp, stats, parallel, MASS, kinship2 (>= 1.8.5), graphics,
grDevices

Description Implements statistical methods for estimating disease penetrance in family-based studies. Penetrance refers to the probability of disease manifestation in individuals carrying specific genetic variants. The package provides tools for age-specific penetrance estimation, handling missing data, and accounting for ascertainment bias in family studies.

Cite as: Kubista, N., Braun, D. & Parmigiani, G. (2024) <[doi:10.48550/arXiv.2411.18816](https://doi.org/10.48550/arXiv.2411.18816)>.

License GPL-3

URL <https://github.com/nicokubi/penetrance>,
<https://nicokubi.github.io/penetrance/>

Encoding UTF-8

LazyData true

LazyDataCompression xz

RoxygenNote 7.3.2

VignetteBuilder knitr

SystemRequirements C++11

Suggests knitr, rmarkdown, scales, ggplot2, roxygen2, testthat (>= 3.0.0)

Config/testthat/edition 3

NeedsCompilation no

Author Nicolas Kubista [aut, cre],
BayesMendel Lab [aut]

Maintainer Nicolas Kubista <bmendel@jimmy.harvard.edu>

Repository CRAN

Date/Publication 2025-03-03 11:40:05 UTC

Contents

absValue	3
apply_burn_in	3
apply_thinning	4
baseline_data_default	4
calculateBaseline	5
calculateEmpiricalDensity	5
calculateNCPen	6
calculate_weibull_parameters	7
combine_chains	8
combine_chains_noSex	8
distribution_data_default	9
drawBaseline	9
drawEmpirical	10
generate_density_plots	10
generate_summary	11
generate_summary_noSex	12
imputeAges	12
imputeAgesInit	15
imputeUnaffectedAges	16
lik.fn	17
lik_noSex	18
makePriors	19
mhChain	20
mhLogLikelihood_clipp	23
mhLogLikelihood_clipp_noSex	24
out_sim	26
penetrance	27
plot_acf	31
plot_loglikelihood	32
plot_pdf	32
plot_penetrance	33
plot_trace	34
printRejectionRates	34
prior_params_default	35
risk_proportion_default	36
simulated_families	36
test_fam2	37
transformDF	38
validate_weibull_parameters	39

Index

41

absValue	<i>Function to return absolute values</i>
----------	---

Description

Function to return absolute values

Usage

```
absValue(x)
```

Arguments

x	Numeric, the input value.
---	---------------------------

Value

Numeric, the absolute value of the input.

apply_burn_in	<i>Apply Burn-In</i>
---------------	----------------------

Description

Apply Burn-In

Usage

```
apply_burn_in(results, burn_in)
```

Arguments

results	A list of MCMC chain results.
burn_in	The fraction roportion of results to discard as burn-in (0 to 1). The default is no burn-in, burn_in=0.

Value

A list of results with burn-in applied.

apply_thinning	<i>Apply Thinning</i>
----------------	-----------------------

Description

Apply Thinning

Usage

```
apply_thinning(results, thinning_factor)
```

Arguments

results	A list of MCMC chain results.
thinning_factor	The factor by which to thin the results (positive integer). The default thinning factor is 1, which implies no thinning.

Value

A list of results with thinning applied.

baseline_data_default	<i>Default Baseline Data</i>
-----------------------	------------------------------

Description

This dataset contains age-specific cancer penetrance rates for both females and males. As an example, the data is derived from the SEER program for colorectal cancer females and males.

Usage

```
baseline_data_default
```

```
baseline_data_default
```

Format

A data frame with 94 rows and 3 variables:

Age Age in years

Female Female baseline risk

Male Male baseline risk

A data frame with 94 rows and 3 variables:

Age Age in years (1 to 94)
Female Penetrance rate for females
Male Penetrance rate for males

calculateBaseline *Calculate Baseline Risk*

Description

This function extracts the penetrance data for a specified cancer type, gene, race, and penetrance type from the provided database.

Usage

```
calculateBaseline(cancer_type, gene, race, type, db)
```

Arguments

cancer_type	The type of cancer for which the risk is being calculated.
gene	The gene of interest for which the risk is being calculated.
race	The race of the individual.
type	The type of penetrance calculation.
db	The dataset used for the calculation, containing penetrance data.

Value

A matrix of penetrance data for the specified parameters.

calculateEmpiricalDensity
Calculate Empirical Age Density

Description

Calculates empirical age density distributions for different subgroups in the data, separated by sex and genetic testing status.

Usage

```
calculateEmpiricalDensity(  
  data,  
  aff_column = "aff",  
  age_column = "age",  
  sex_column = "sex",  
  geno_column = "geno",  
  n_points = 10000,  
  sex_specific = TRUE  
)
```

Arguments

data	A data frame containing the family data
aff_column	Name of the affection status column
age_column	Name of the age column
sex_column	Name of the sex column
geno_column	Name of the genotype column
n_points	Number of points to use in density estimation
sex_specific	Logical; whether to calculate sex-specific densities

Value

A list of density objects for different subgroups (tested/untested, male/female)

calculateNCPen	<i>Calculate Age-Specific Non-Carrier Penetrance</i>
----------------	--

Description

This function calculates the age-specific non-carrier penetrance based on SEER baseline data, penetrances for carriers, and allele frequencies. It adjusts penetrance estimates for genetic testing by incorporating the genetic risk attributable to specified alleles.

Usage

```
calculateNCPen(SEER_baseline, alpha, beta, delta, gamma, prev, max_age)
```

Arguments

SEER_baseline	Numeric, the baseline penetrance derived from SEER data for the general population without considering genetic risk factors.
alpha	Numeric, shape parameter for the Weibull distribution used to model carrier risk.
beta	Numeric, scale parameter for the Weibull distribution used to model carrier risk.
delta	Numeric, location parameter for the Weibull distribution used to model carrier risk.
gamma	Numeric, scaling factor applied to the Weibull distribution to adjust carrier risk.
prev	Numeric, the prevalence of the risk allele in the population.
max_age	Integer, the maximum age up to which the calculations are performed.

Value

A list containing:

weightedCarrierRisk

Numeric vector, the weighted risk for carriers at each age based on prevalence.

yearlyProb

Numeric vector, the yearly probability of not getting the disease at each age.

cumulativeProb

Numeric vector, the cumulative probability of not getting the disease up to each age.

calculate_weibull_parameters

Calculate Weibull Parameters

Description

This function calculates the shape (alpha) and scale (beta) parameters of a Weibull distribution given the median, first quartile, and delta values.

Usage

```
calculate_weibull_parameters(given_median, given_first_quartile, delta)
```

Arguments

given_median The median of the data.

given_first_quartile

The first quartile of the data.

delta

A constant offset value.

Value

A list containing the calculated Weibull parameters:

alpha The shape parameter of the Weibull distribution

beta The scale parameter of the Weibull distribution

Examples

```
# Calculate Weibull parameters
params <- calculate_weibull_parameters(
  given_median = 50,
  given_first_quartile = 30,
  delta = 15
)
print(params)
```

combine_chains	<i>Combine Chains Function to combine the posterior samples from the multiple chains.</i>
----------------	---

Description

Combine Chains Function to combine the posterior samples from the multiple chains.

Usage

```
combine_chains(results)
```

Arguments

results A list of MCMC chain results.

Value

A list with combined results, including median, threshold, first quartile, and asymptote values.

combine_chains_noSex	<i>Combine Chains for Non-Sex-Specific Estimation</i>
----------------------	---

Description

Combines the posterior samples from multiple MCMC chains for non-sex-specific estimations.

Usage

```
combine_chains_noSex(results)
```

Arguments

results A list of MCMC chain results, where each element contains posterior samples of parameters.

Value

A list with combined results, including samples for median, threshold, first quartile, asymptote values, log-likelihoods, and log-acceptance ratios.

distribution_data_default
Default Distribution Data

Description

Default data frame structure with row names for use in the makePriors function.

Usage

```
distribution_data_default
```

```
distribution_data_default
```

Format

A data frame for prior distribution parameters

A data frame with the following columns:

age Age values (NA for default).

at_risk Proportion of people at risk (NA for default).

drawBaseline *Draw Ages Using the Inverse CDF Method from the baseline data*

Description

This function draws ages using the inverse CDF method from baseline data.

Usage

```
drawBaseline(baseline_data)
```

Arguments

baseline_data A data frame containing baseline data with columns 'cum_prob' and 'age'.

Value

A single age value drawn from the baseline data.

drawEmpirical	<i>Draw Ages Using the Inverse CDF Method from Empirical Density</i>
---------------	--

Description

This function draws ages using the inverse CDF method from empirical density data, based on sex and whether the individual was tested.

Usage

```
drawEmpirical(empirical_density, sex, tested, sex_specific = TRUE)
```

Arguments

empirical_density	A list of density objects containing the empirical density of ages for different groups.
sex	Numeric, the sex of the individual (1 for male, 2 for female).
tested	Logical, indicating whether the individual was tested (has a non-NA 'geno' value).
sex_specific	Logical, indicating whether the imputation should be sex-specific. Default is TRUE.

Value

A single age value drawn from the appropriate empirical density data.

generate_density_plots	<i>Generate Posterior Density Plots</i>
------------------------	---

Description

Generates histograms of the posterior samples for the different parameters

Usage

```
generate_density_plots(data)
```

Arguments

data	A list with combined results.
------	-------------------------------

Value

No return value, called for side effects. Creates density plots for each parameter.

Examples

```
# Create example data
data <- list(
  median_male_results = rnorm(1000, 50, 5),
  median_female_results = rnorm(1000, 45, 5),
  threshold_male_results = runif(1000, 20, 30),
  threshold_female_results = runif(1000, 25, 35),
  asymptote_male_results = rbeta(1000, 2, 2),
  asymptote_female_results = rbeta(1000, 2, 2)
)

# Generate density plots
old_par <- par(no.readonly = TRUE) # Save old par settings
generate_density_plots(data)
par(old_par) # Restore old par settings
```

generate_summary	<i>Generate Summary</i>
------------------	-------------------------

Description

Function to generate summary statistics

Usage

```
generate_summary(data, verbose = FALSE)
```

Arguments

data	A list with combined results.
verbose	Logical, whether to print summary to console. Default is FALSE.

Value

A data.frame containing summary statistics (min, 1st quartile, median, mean, 3rd quartile, max) for each parameter.

`generate_summary_noSex`*Generate Summary for Non-Sex-Specific Estimation*

Description

Generates summary statistics for the combined MCMC results for non-sex-specific estimations.

Usage

```
generate_summary_noSex(data, verbose = FALSE)
```

Arguments

<code>data</code>	A list containing combined results of MCMC chains, typically the output of <code>combine_chains_noSex</code> .
<code>verbose</code>	Logical, whether to print summary to console. Default is FALSE.

Value

A data.frame containing summary statistics (min, 1st quartile, median, mean, 3rd quartile, max) for median, threshold, first quartile, and asymptote values.

`imputeAges`*Impute Missing Ages in Family-Based Data*

Description

Imputes missing ages in family-based data using a combination of Weibull distributions for affected individuals and empirical distributions for unaffected individuals. The function can perform both sex-specific and non-sex-specific imputations.

Usage

```
imputeAges(  
  data,  
  na_indices,  
  baseline_male = NULL,  
  baseline_female = NULL,  
  alpha_male = NULL,  
  beta_male = NULL,  
  delta_male = NULL,  
  alpha_female = NULL,  
  beta_female = NULL,  
  delta_female = NULL,
```

```

baseline = NULL,
alpha = NULL,
beta = NULL,
delta = NULL,
max_age,
sex_specific = TRUE,
max_attempts = 100,
geno_freq,
trans,
lik
)

```

Arguments

<code>data</code>	A data frame containing family-based data with columns: family, individual, father, mother, sex, aff, age, geno, and isProband
<code>na_indices</code>	Vector of indices where ages need to be imputed
<code>baseline_male, baseline_female</code>	Data frames containing baseline age distributions for males/females
<code>alpha_male, alpha_female</code>	Shape parameters for male/female Weibull distributions
<code>beta_male, beta_female</code>	Scale parameters for male/female Weibull distributions
<code>delta_male, delta_female</code>	Location parameters for male/female Weibull distributions
<code>baseline</code>	Data frame containing overall baseline age distribution (non-sex-specific)
<code>alpha, beta, delta</code>	Overall Weibull parameters (non-sex-specific)
<code>max_age</code>	Maximum allowable age
<code>sex_specific</code>	Logical; whether to use sex-specific parameters
<code>max_attempts</code>	Maximum number of attempts for generating valid ages
<code>geno_freq</code>	Vector of genotype frequencies
<code>trans</code>	Transmission probabilities
<code>lik</code>	Likelihood matrix

Value

A data frame with the following modifications:

<code>age</code>	Updated with imputed ages for previously missing values
------------------	---

The rest of the data frame remains unchanged.

Examples

```

# Create sample data with the same structure as used in mhChain
data <- data.frame(
  family = rep(1:2, each=5),
  individual = rep(1:5, 2),
  father = c(NA,1,1,1,1, NA,6,6,6,6),
  mother = c(NA,2,2,2,2, NA,7,7,7,7),
  sex = c(1,2,1,2,1, 1,2,1,2,1),
  aff = c(1,0,1,0,NA, 1,0,1,0,NA),
  age = c(45,NA,25,NA,20, 50,NA,30,NA,22),
  geno = c("1/2",NA,"1/2",NA,NA, "1/2",NA,"1/2",NA,NA),
  isProband = c(1,0,0,0,0, 1,0,0,0,0)
)

# Initialize parameters
na_indices <- which(is.na(data$age))
geno_freq <- c(0.999, 0.001) # Frequency of normal and risk alleles
trans <- matrix(c(1,0,0.5,0.5), nrow=2) # Transmission matrix
lik <- matrix(1, nrow=nrow(data), ncol=2) # Likelihood matrix

# Create baseline data for both sex-specific and non-sex-specific cases
age_range <- 20:94
n_ages <- length(age_range)

# Sex-specific baseline data
baseline_male <- data.frame(
  age = age_range,
  cum_prob = (1:n_ages)/n_ages * 0.8 # Male cumulative probabilities
)

baseline_female <- data.frame(
  age = age_range,
  cum_prob = (1:n_ages)/n_ages * 0.9 # Female cumulative probabilities
)

# Non-sex-specific baseline data
baseline <- data.frame(
  age = age_range,
  cum_prob = (1:n_ages)/n_ages * 0.85 # Overall cumulative probabilities
)

# Example with sex-specific imputation
imputed_data_sex <- imputeAges(
  data = data,
  na_indices = na_indices,
  baseline_male = baseline_male,
  baseline_female = baseline_female,
  alpha_male = 3.5,
  beta_male = 20,
  delta_male = 20,
  alpha_female = 3.2,
  beta_female = 18,

```

```

    delta_female = 18,
    max_age = 94,
    sex_specific = TRUE,
    geno_freq = geno_freq,
    trans = trans,
    lik = lik
  )

# Example with non-sex-specific imputation
imputed_data_nosex <- imputeAges(
  data = data,
  na_indices = na_indices,
  baseline = baseline,
  alpha = 3.3,
  beta = 19,
  delta = 19,
  max_age = 94,
  sex_specific = FALSE,
  geno_freq = geno_freq,
  trans = trans,
  lik = lik
)

```

imputeAgesInit	<i>Initialize Age Imputation</i>
----------------	----------------------------------

Description

Initializes the age imputation process by filling missing ages with random values between a threshold and maximum age.

Usage

```
imputeAgesInit(data, threshold, max_age)
```

Arguments

data	A data frame containing family-based data
threshold	Minimum age value for initialization
max_age	Maximum age value for initialization

Value

A list containing:

data	The data frame with initialized ages
na_indices	Indices of missing age values

Examples

```
# Create sample data
data <- data.frame(
  family = c(1, 1),
  individual = c(1, 2),
  father = c(NA, 1),
  mother = c(NA, NA),
  sex = c(1, 2),
  aff = c(1, 0),
  age = c(NA, NA),
  geno = c("1/2", NA),
  isProband = c(1, 0)
)

# Initialize ages with random values between 20 and 94
result <- imputeAgesInit(data, threshold = 20, max_age = 94)

# Access the results
imputed_data <- result$data
missing_indices <- result$na_indices
```

imputeUnaffectedAges *Impute Ages for Unaffected Individuals*

Description

This function imputes ages for unaffected individuals in a dataset based on their sex and whether they were tested, using empirical age distributions.

Usage

```
imputeUnaffectedAges(data, na_indices, empirical_density, max_age)
```

Arguments

data	A data frame containing the individual data, including columns for age, sex, and geno.
na_indices	A vector of indices indicating the rows in the data where ages need to be imputed.
empirical_density	A list of density objects containing the empirical density of ages for different groups.
max_age	Integer, the maximum age considered in the analysis.

Value

The data frame with imputed ages for unaffected individuals.

lik.fn

*Penetrance Function***Description**

Calculates the penetrance for an individual based on Weibull distribution parameters. This function estimates the probability of developing cancer given the individual's genetic and demographic information.

Usage

```
lik.fn(
  i,
  data,
  alpha_male,
  alpha_female,
  beta_male,
  beta_female,
  delta_male,
  delta_female,
  gamma_male,
  gamma_female,
  max_age,
  baselineRisk,
  BaselineNC,
  prev
)
```

Arguments

i	Integer, index of the individual in the data set.
data	Data frame, containing individual demographic and genetic information. Must include columns for 'sex', 'age', 'aff' (affection status), and 'geno' (genotype).
alpha_male	Numeric, Weibull distribution shape parameter for males.
alpha_female	Numeric, Weibull distribution shape parameter for females.
beta_male	Numeric, Weibull distribution scale parameter for males.
beta_female	Numeric, Weibull distribution scale parameter for females.
delta_male	Numeric, shift parameter for the Weibull function for males.
delta_female	Numeric, shift parameter for the Weibull function for females.
gamma_male	Numeric, asymptote parameter for males (only scales the entire distribution).
gamma_female	Numeric, asymptote parameter for females (only scales the entire distribution).
max_age	Integer, maximum age considered in the analysis.
baselineRisk	Numeric matrix, baseline risk for each age by sex. Rows correspond to sex (1 for male, 2 for female) and columns to age.
BaselineNC	Logical, indicates if non-carrier penetrance should be based on SEER data.
prev	Numeric, prevalence of the risk allele in the population.

Value

Numeric vector, containing penetrance values for unaffected and affected individuals.

 lik_noSex

Likelihood Calculation without Sex Differentiation

Description

This function calculates the likelihood for an individual based on Weibull distribution parameters without considering sex differentiation.

Usage

```
lik_noSex(
  i,
  data,
  alpha,
  beta,
  delta,
  gamma,
  max_age,
  baselineRisk,
  BaselineNC,
  prev
)
```

Arguments

i	Integer, index of the individual in the data set.
data	Data frame, containing individual demographic and genetic information. Must include columns for 'age', 'aff' (affection status), and 'geno' (genotype).
alpha	Numeric, Weibull distribution shape parameter.
beta	Numeric, Weibull distribution scale parameter.
delta	Numeric, shift parameter for the Weibull function.
gamma	Numeric, asymptote parameter (only scales the entire distribution).
max_age	Integer, maximum age considered in the analysis.
baselineRisk	Numeric vector, baseline risk for each age.
BaselineNC	Logical, indicates if non-carrier penetrance should be based on SEER data or the calculated non-carrier penetrance.
prev	Numeric, prevalence of the risk allele in the population.

Value

Numeric vector, containing likelihood values for unaffected and affected individuals.

`makePriors`*Make Priors*

Description

This function generates prior distributions based on user input or default parameters. It is designed to aid in the statistical analysis of risk proportions in populations, particularly in the context of cancer research. The distributions are calculated for various statistical metrics such as asymptote, threshold, median, and first quartile.

Usage

```
makePriors(  
  data,  
  sample_size,  
  ratio,  
  prior_params,  
  risk_proportion,  
  baseline_data  
)
```

Arguments

<code>data</code>	A data frame containing age and risk data. If NULL or contains NA values, default parameters are used.
<code>sample_size</code>	Numeric, the total sample size used for risk proportion calculations.
<code>ratio</code>	Numeric, the odds ratio (OR) or relative risk (RR) used in asymptote parameter calculations.
<code>prior_params</code>	List, containing prior parameters for the beta distributions. If NULL, default parameters are used.
<code>risk_proportion</code>	Data frame, with default proportions of people at risk.
<code>baseline_data</code>	Data frame with the baseline risk data.

Details

The function includes internal helper functions for normalizing median and first quartile values, and for computing beta distribution parameters. The function handles various settings: using default parameters, applying user inputs, and calculating parameters based on sample size and risk proportions.

If the OR/RR ratio is provided, the asymptote parameters are computed based on this ratio, overriding other inputs for the asymptote.

The function returns a list of distribution functions for the asymptote, threshold, median, and first quartile, which can be used for further statistical analysis.

Value

A list of functions representing the prior distributions for asymptote, threshold, median, and first quartile.

See Also

[qbeta](#), [runif](#)

mhChain	<i>Execution of a Single Chain in Metropolis-Hastings for Cancer Risk Estimation</i>
---------	--

Description

Performs a single chain execution in the Metropolis-Hastings algorithm for Bayesian inference, specifically tailored for cancer risk estimation. This function can handle both sex-specific and non-sex-specific scenarios.

Usage

```
mhChain(  
  seed,  
  n_iter,  
  burn_in,  
  chain_id,  
  ncores,  
  data,  
  twins,  
  max_age,  
  baseline_data,  
  prior_distributions,  
  prev,  
  median_max,  
  BaselineNC,  
  var,  
  age_imputation,  
  imp_interval,  
  remove_proband,  
  sex_specific  
)
```

Arguments

seed	Integer, the seed for the random number generator to ensure reproducibility.
n_iter	Integer, the number of iterations to perform in the Metropolis-Hastings algorithm.

burn_in	Integer, the number of initial iterations to discard (burn-in period).
chain_id	Integer, the identifier for the chain being executed.
ncores	Integer, the number of cores to use for parallel computation.
data	Data frame, containing family and genetic information used in the analysis.
twins	Information on monozygous twins or triplets in the pedigrees.
max_age	Integer, the maximum age considered in the analysis.
baseline_data	Numeric matrix or vector, containing baseline risk estimates for different ages and sexes.
prior_distributions	List, containing prior distributions for the parameters being estimated.
prev	Numeric, the prevalence of the risk allele in the population.
median_max	Logical, indicates if the maximum median age should be used for the Weibull distribution.
BaselineNC	Logical, indicates if non-carrier penetrance should be based on SEER data.
var	Numeric, the variance for the proposal distribution in the Metropolis-Hastings algorithm.
age_imputation	Logical, indicates if age imputation should be performed.
imp_interval	Integer, the interval at which age imputation should be performed when age_imputation = TRUE.
remove_proband	Logical, indicates if the proband should be removed from the analysis.
sex_specific	Logical, indicates if the analysis should differentiate by sex.

Value

A list containing samples, log likelihoods, log-acceptance ratio, and rejection rate for each iteration.

Examples

```
# Create sample data in PanelPRO format
data <- data.frame(
  ID = 1:10,
  PedigreeID = rep(1, 10),
  Sex = c(0, 1, 0, 1, 0, 1, 0, 1, 0, 1), # 0=female, 1=male
  MotherID = c(NA, NA, 1, 1, 3, 3, 5, 5, 7, 7),
  FatherID = c(NA, NA, 2, 2, 4, 4, 6, 6, 8, 8),
  isProband = c(1, rep(0, 9)),
  CurAge = c(45, 35, 55, 40, 50, 45, 60, 38, 52, 42),
  isAff = c(1, 0, 1, 0, 1, 0, 1, 0, 1, 0),
  Age = c(40, NA, 50, NA, 45, NA, 55, NA, 48, NA),
  Geno = c(1, NA, 1, 0, 1, 0, NA, NA, 1, NA)
)

# Transform data into required format
data <- transformDF(data)

# Set parameters for the chain
```

```

seed <- 123
n_iter <- 10
burn_in <- 0.1 # 10% burn-in
chain_id <- 1
ncores <- 1
max_age <- 100

# Create baseline data (simplified example)
baseline_data <- matrix(
  c(rep(0.005, max_age), rep(0.008, max_age)), # Increased baseline risks
  ncol = 2,
  dimnames = list(NULL, c("Male", "Female")))
)

# Set prior distributions with carefully chosen bounds
prior_distributions <- list(
  prior_params = list(
    asymptote = list(g1 = 2, g2 = 3), # Mode around 0.4
    threshold = list(min = 20, max = 30), # Narrower range for threshold
    median = list(m1 = 3, m2 = 2), # Mode around 0.6
    first_quartile = list(q1 = 2, q2 = 3) # Mode around 0.4
  )
)

# Create variance vector for all 8 parameters in sex-specific case
# Using very small variances for initial stability
var <- c(0.005, 0.005, # asymptotes (smaller variance since between 0-1)
        1, 1, # thresholds
        1, 1, # medians
        1, 1) # first quartiles

# Run the chain
results <- mhChain(
  seed = seed,
  n_iter = n_iter,
  burn_in = burn_in,
  chain_id = chain_id,
  ncores = ncores,
  data = data,
  twins = NULL,
  max_age = max_age,
  baseline_data = baseline_data,
  prior_distributions = prior_distributions,
  prev = 0.05, # Increased prevalence
  median_max = FALSE, # Changed to FALSE for simpler median constraints
  BaselineNC = TRUE,
  var = var,
  age_imputation = FALSE,
  imp_interval = 10,
  remove_proband = TRUE,
  sex_specific = TRUE
)

```

mhLogLikelihood_clipp *Calculate Log Likelihood using clipp Package*

Description

Calculate Log Likelihood using clipp Package

Usage

```
mhLogLikelihood_clipp(  
  paras,  
  families,  
  twins,  
  max_age,  
  baseline_data,  
  prev,  
  geno_freq,  
  trans,  
  BaselineNC,  
  ncores  
)
```

Arguments

paras	Numeric vector of parameters
families	Data frame of pedigree information
twins	Information on monozygous twins
max_age	Integer, maximum age
baseline_data	Numeric matrix of baseline risk data
prev	Numeric, prevalence
geno_freq	Numeric vector of frequencies
trans	Numeric matrix of transmission probabilities
BaselineNC	Logical for baseline choice
ncores	Integer for parallel computation

Value

Numeric value representing the calculated log likelihood.

Examples

```

# Create example parameters and data
paras <- c(0.8, 0.7, 20, 25, 50, 45, 30, 35) # Example parameters

# Create sample data in PanelPRO format
families <- data.frame(
  ID = 1:10,
  PedigreeID = rep(1, 10),
  Sex = c(0, 1, 0, 1, 0, 1, 0, 1, 0, 1), # 0=female, 1=male
  MotherID = c(NA, NA, 1, 1, 3, 3, 5, 5, 7, 7),
  FatherID = c(NA, NA, 2, 2, 4, 4, 6, 6, 8, 8),
  isProband = c(1, rep(0, 9)),
  CurAge = c(45, 35, 55, 40, 50, 45, 60, 38, 52, 42),
  isAff = c(1, 0, 1, 0, 1, 0, 1, 0, 1, 0),
  Age = c(40, NA, 50, NA, 45, NA, 55, NA, 48, NA),
  Geno = c(1, NA, 1, 0, 1, 0, NA, NA, 1, NA)
)

# Transform data into required format
families <- transformDF(families)

trans <- matrix(
  c(
    1, 0, # both parents are wild type
    0.5, 0.5, # mother is wildtype and father is a heterozygous carrier
    0.5, 0.5, # father is wildtype and mother is a heterozygous carrier
    1 / 3, 2 / 3 # both parents are heterozygous carriers
  ),
  nrow = 4, ncol = 2, byrow = TRUE
)

# Calculate log likelihood
loglik <- mhLogLikelihood_clipp(
  paras = paras,
  families = families,
  twins = NULL,
  max_age = 94,
  baseline_data = baseline_data_default,
  prev = 0.001,
  geno_freq = c(0.999, 0.001),
  trans = trans,
  BaselineNC = TRUE,
  ncores = 1
)

```


Description

This function calculates the log likelihood for a set of parameters and data without considering sex differentiation using the clipp package.

Usage

```
mhLogLikelihood_clipp_noSex(  
  paras,  
  families,  
  twins,  
  max_age,  
  baseline_data,  
  prev,  
  geno_freq,  
  trans,  
  BaselineNC,  
  ncores  
)
```

Arguments

paras	Numeric vector, the parameters for the Weibull distribution and scaling factors. Should contain in order: gamma, delta, given_median, given_first_quartile.
families	Data frame, containing pedigree information with columns for 'age', 'aff' (affection status), and 'geno' (genotype).
twins	Information on monozygous twins or triplets in the pedigrees.
max_age	Integer, maximum age considered in the analysis.
baseline_data	Numeric vector, baseline risk data for each age.
prev	Numeric, prevalence of the risk allele in the population.
geno_freq	Numeric vector, represents the frequency of the risk type and its complement in the population.
trans	Numeric matrix, transition matrix that defines the probabilities of allele transmission from parents to offspring.
BaselineNC	Logical, indicates if non-carrier penetrance should be based on the baseline data or the calculated non-carrier penetrance.
ncores	Integer, number of cores to use for parallel computation.

Value

Numeric, the calculated log likelihood.

References

Details about the clipp package and methods can be found in the package documentation.

out_sim	<i>Simulated Output Data</i>
---------	------------------------------

Description

This dataset contains the simulated output data for the penetrance package.

Usage

```
data(out_sim)
```

Format

A list with the following components:

summary_stats A data frame with 18000 observations of 8 variables:

- Median_Male** numeric, Median value for males
- Median_Female** numeric, Median value for females
- Threshold_Male** numeric, Threshold value for males
- Threshold_Female** numeric, Threshold value for females
- First_Quartile_Male** numeric, First quartile value for males
- First_Quartile_Female** numeric, First quartile value for females
- Asymptote_Male** numeric, Asymptote value for males
- Asymptote_Female** numeric, Asymptote value for females

density_plots A list of 1 element, mfrow: integer vector of length 2

trace_plots A list of 1 element, mfrow: integer vector of length 2

penetrance_plot A list of 2 elements: rect and text

pdf_plots A list of 2 elements: rect and text

combined_chains A list of 19 numeric vectors with 18000 elements each

results A list of 1 element which is a list of 24 elements, each with 18000 elements

data A data frame with 4727 observations of 9 variables:

- individual** integer, Individual ID
- isProband** numeric, Indicator if the individual is a proband
- family** integer, Family ID
- mother** numeric, Mother's ID
- father** numeric, Father's ID
- aff** numeric, Affected status
- sex** numeric, Sex of the individual
- age** numeric, Age of the individual
- geno** character, Genotype

Examples

```
data(out_sim)
head(out_sim$summary_stats)
```

Description

A comprehensive package for penetrance estimation in family-based studies. This package implements Bayesian methods using Metropolis-Hastings algorithm for estimating age-specific penetrance of genetic variants. It supports both sex-specific and non-sex-specific analyses, and provides various visualization tools for examining MCMC results.

This function implements the Independent Metropolis-Hastings algorithm for Bayesian penetrance estimation of cancer risk. It utilizes parallel computing to run multiple chains and provides various options for analyzing and visualizing the results.

Usage

```
penetrance(  
  pedigree,  
  twins = NULL,  
  n_chains = 1,  
  n_iter_per_chain = 10000,  
  ncores = 6,  
  max_age = 94,  
  baseline_data = baseline_data_default,  
  remove_proband = FALSE,  
  age_imputation = FALSE,  
  median_max = TRUE,  
  BaselineNC = TRUE,  
  var = c(0.1, 0.1, 2, 2, 5, 5, 5, 5),  
  burn_in = 0,  
  thinning_factor = 1,  
  imp_interval = 100,  
  distribution_data = distribution_data_default,  
  prev = 1e-04,  
  sample_size = NULL,  
  ratio = NULL,  
  prior_params = prior_params_default,  
  risk_proportion = risk_proportion_default,  
  summary_stats = TRUE,  
  rejection_rates = TRUE,  
  density_plots = TRUE,  
  plot_trace = TRUE,  
  penetrance_plot = TRUE,  
  penetrance_plot_pdf = TRUE,  
  plot_loglikelihood = TRUE,  
  plot_acf = TRUE,  
  probCI = 0.95,
```

```
sex_specific = TRUE
)
```

Arguments

pedigree	<p>A data frame containing the pedigree data in the required format. It should include the following columns:</p> <ul style="list-style-type: none"> • PedigreeID: A numeric value representing the unique identifier for each family. There should be no duplicated entries. • ID: A numeric value representing the unique identifier for each individual. There should be no duplicated entries. • Sex: A numeric value where 0 indicates female and 1 indicates male. Unknown sex needs to be coded as NA. • MotherID: A numeric value representing the unique identifier for an individual's mother. • FatherID: A numeric value representing the unique identifier for an individual's father. • isProband: A numeric value where 1 indicates the individual is a proband and 0 otherwise. • CurAge: A numeric value indicating the age of censoring (current age if the person is alive or age at death if the person is deceased). Allowed ages range from 1 to 94. Unknown ages can be left empty or coded as NA. • isAff: A numeric value indicating the affection status of cancer, with 1 for diagnosed individuals, 0 for unaffected individuals, and NA for unknown status. • Age: A numeric value indicating the age of cancer diagnosis, encoded as NA if the individual was not diagnosed. Allowed ages range from 1 to 94. Unknown ages can be left empty or coded as NA. • geno: A column for germline testing or tumor marker testing results. Positive results should be coded as 1, negative results as 0, and unknown results as NA or left empty.
twins	<p>A list specifying identical twins or triplets in the family. For example, to indicate that "ora024" and "ora027" are identical twins, and "aey063" and "aey064" are identical twins, use the following format: <code>twins <- list(c("ora024", "ora027"), c("aey063", "aey064"))</code>.</p>
n_chains	Integer, the number of chains for parallel computation. Default is 1.
n_iter_per_chain	Integer, the number of iterations for each chain. Default is 10000.
ncores	Integer, the number of cores for parallel computation. Default is 6.
max_age	Integer, the maximum age considered for analysis. Default is 94.
baseline_data	Data for the baseline risk estimates (probability of developing cancer), such as population-level risk from a cancer registry. Default data, for exemplary purposes, is for Colorectal cancer from the SEER database.
remove_proband	Logical, indicating whether to remove probands from the analysis. Default is FALSE.

age_imputation	Logical, indicating whether to perform age imputation. Default is FALSE.
median_max	Logical, indicating whether to use the baseline median age or max_age as an upper bound for the median proposal. Default is TRUE.
BaselineNC	Logical, indicating that the non-carrier penetrance is assumed to be the baseline penetrance. Default is TRUE.
var	Numeric vector, variances for the proposal distribution in the Metropolis-Hastings algorithm. Default is $c(0.1, 0.1, 2, 2, 5, 5, 5, 5)$.
burn_in	Numeric, the fraction of results to discard as burn-in (0 to 1). Default is 0 (no burn-in).
thinning_factor	Integer, the factor by which to thin the results. Default is 1 (no thinning).
imp_interval	Integer, the interval at which age imputation should be performed when age_imputation = TRUE.
distribution_data	Data for generating prior distributions.
prev	Numeric, prevalence of the carrier status. Default is 0.0001.
sample_size	Optional numeric, sample size for distribution generation.
ratio	Optional numeric, ratio parameter for distribution generation.
prior_params	List, parameters for prior distributions.
risk_proportion	Numeric, proportion of risk for distribution generation.
summary_stats	Logical, indicating whether to include summary statistics in the output. Default is TRUE.
rejection_rates	Logical, indicating whether to include rejection rates in the output. Default is TRUE.
density_plots	Logical, indicating whether to include density plots in the output. Default is TRUE.
plot_trace	Logical, indicating whether to include trace plots in the output. Default is TRUE.
penetrance_plot	Logical, indicating whether to include penetrance plots in the output. Default is TRUE.
penetrance_plot_pdf	Logical, indicating whether to include PDF plots in the output. Default is TRUE.
plot_loglikelihood	Logical, indicating whether to include log-likelihood plots in the output. Default is TRUE.
plot_acf	Logical, indicating whether to include autocorrelation function (ACF) plots for posterior samples. Default is TRUE.
probCI	Numeric, probability level for credible intervals in penetrance plots. Must be between 0 and 1. Default is 0.95.
sex_specific	Logical, indicating whether to use sex-specific parameters in the analysis. Default is TRUE.

Details

Key features:

- Bayesian estimation of penetrance using family-based data
- Support for sex-specific and non-sex-specific analyses
- Age imputation for missing data
- Visualization tools for MCMC diagnostics
- Integration with the `clipp` package for likelihood calculations

Value

A list containing combined results from all chains, including optional statistics and plots.

Author(s)

Maintainer: Nicolas Kubista <bmendel@jimmy.harvard.edu>

Authors:

- BayesMendel Lab

See Also

Useful links:

- <https://github.com/nicokubi/penetrance>
- <https://nicokubi.github.io/penetrance/>

Examples

```
# Create example baseline data (simplified for demonstration)
baseline_data_default <- data.frame(
  Age = 1:94,
  Female = rep(0.01, 94),
  Male = rep(0.01, 94)
)

# Create example distribution data
distribution_data_default <- data.frame(
  Age = 1:94,
  Risk = rep(0.01, 94)
)

# Create example prior parameters
prior_params_default <- list(
  shape = 2,
  scale = 50
)

# Create example risk proportion
risk_proportion_default <- 0.5
```

```

# Create a simple example pedigree
example_pedigree <- data.frame(
  PedigreeID = rep(1, 4),
  ID = 1:4,
  Sex = c(1, 0, 1, 0), # 1 for male, 0 for female
  MotherID = c(NA, NA, 2, 2),
  FatherID = c(NA, NA, 1, 1),
  isProband = c(0, 0, 1, 0),
  CurAge = c(70, 68, 45, 42),
  isAff = c(0, 0, 1, 0),
  Age = c(NA, NA, 40, NA),
  geno = c(NA, NA, 1, NA)
)

# Basic usage with minimal iterations
result <- penetrance(
  pedigree = list(example_pedigree),
  n_chains = 1,
  n_iter_per_chain = 10, # Very small number for example
  ncores = 1,           # Single core for example
  summary_stats = TRUE,
  plot_trace = FALSE,  # Disable plots for quick example
  density_plots = FALSE,
  penetrance_plot = FALSE,
  penetrance_plot_pdf = FALSE,
  plot_loglikelihood = FALSE,
  plot_acf = FALSE
)

# View basic results
head(result$summary_stats)

```

plot_acf

Plot Autocorrelation for Multiple MCMC Chains (Posterior Samples)

Description

This function plots the autocorrelation for sex-specific or non-sex-specific posterior samples across multiple MCMC chains. It defaults to key parameters like `asymptote_male_samples`, `asymptote_female_samples`, etc.

Usage

```
plot_acf(results, n_chains, max_lag = 50)
```

Arguments

results	A list of MCMC chain results.
n_chains	The number of chains.
max_lag	Integer, the maximum lag to be considered for the autocorrelation plot. Default is 50.

Value

A series of autocorrelation plots for each chain.

plot_loglikelihood *Plot Log-Likelihood for Multiple MCMC Chains*

Description

This function plots the log-likelihood values across iterations for multiple MCMC chains. It helps visualize the convergence of the chains based on the log-likelihood values.

Usage

```
plot_loglikelihood(results, n_chains)
```

Arguments

results	A list of MCMC chain results, each containing the loglikelihood_current values.
n_chains	The number of chains.

Value

A series of log-likelihood plots for each chain.

plot_pdf *Plot Weibull Probability Density Function with Credible Intervals*

Description

This function plots the Weibull PDF with credible intervals for the given MCMC results. It allows for visualization of density curves based on the posterior samples.

Usage

```
plot_pdf(combined_chains, prob, max_age, sex = "NA")
```


Arguments

combined_chains	List of combined MCMC chain results containing posterior samples for penetrance parameters.
prob	Numeric, probability level for confidence intervals (between 0 and 1).
max_age	Integer, maximum age to plot.
sex	Character, one of "Male", "Female", or "NA" for non-sex-specific. Default is "NA".

Value

A plot showing the Weibull PDF with credible intervals.

plot_penetrance	<i>Plot Weibull Distribution with Credible Intervals</i>
-----------------	--

Description

This function plots the Weibull distribution with credible intervals for the given MCMC results. It allows for visualization of penetrance curves based on the posterior samples.

Usage

```
plot_penetrance(combined_chains, prob, max_age, sex = "NA")
```

Arguments

combined_chains	List of combined MCMC chain results containing posterior samples for penetrance parameters.
prob	Numeric, probability level for confidence intervals (between 0 and 1).
max_age	Integer, maximum age to plot.
sex	Character, one of "Male", "Female", or "NA" for non-sex-specific. Default is "NA".

Value

A plot showing the Weibull distribution with credible intervals.

plot_trace	<i>Plot MCMC Trace Plots</i>
------------	------------------------------

Description

Plot MCMC Trace Plots

Usage

```
plot_trace(results, n_chains, verbose = FALSE)
```

Arguments

results	A list of MCMC chain results.
n_chains	Integer, the number of chains.
verbose	Logical, whether to print progress messages. Default is FALSE.

Value

No return value, called for side effects. Creates trace plots for each parameter.

Examples

```
# Create example results list
results <- list(
  list(
    median_samples = rnorm(100),
    threshold_samples = runif(100),
    first_quartile_samples = rgamma(100, 2, 2),
    asymptote_samples = rbeta(100, 2, 2)
  )
)

# Generate trace plots
plot_trace(results, n_chains = 1)
```

printRejectionRates	<i>Print MCMC Rejection Rates</i>
---------------------	-----------------------------------

Description

Print MCMC Rejection Rates

Usage

```
printRejectionRates(results, verbose = TRUE)
```

Arguments

`results` A list of MCMC chain results.
`verbose` Logical, whether to print rates to console. Default is TRUE.

Details

Extracts and prints the rejection rates from MCMC chain results.

Value

A named numeric vector containing the rejection rate (between 0 and 1) for each MCMC chain. Names are of the form "Chain X" where X is the chain number.

Examples

```
# Create example results list with two chains
results <- list(
  list(rejection_rate = 0.3),
  list(rejection_rate = 0.4)
)

# Get rejection rates without printing
rates <- printRejectionRates(results, verbose = FALSE)

# Print rejection rates
rates <- printRejectionRates(results)
```

prior_params_default *Default Prior Parameters*

Description

Default parameters for the prior distributions used in the `makePriors` function.

Usage

```
prior_params_default
prior_params_default
```

Format

A list of prior distribution parameters

A list with the following components:

asymptote A list with components `g1` and `g2`, default values for the asymptote parameters.

threshold A list with components `min` and `max`, default values for the threshold parameters.

median A list with components m1 and m2, default values for the median parameters.

first_quartile A list with components q1 and q2, default values for the first quartile parameters.

risk_proportion_default

Default Risk Proportions

Description

Default proportions of people at risk used in the `makePriors` function.

Usage

```
risk_proportion_default
```

```
risk_proportion_default
```

Format

A data frame of risk proportions

A data frame with the following columns:

median Proportion of people at risk at the median age.

first_quartile Proportion of people at risk at the first quartile age.

max_age Proportion of people at risk at the maximum age.

simulated_families *Processed Family Data*

Description

A dataset containing processed information about the first simulated 130 families. These families are referenced in the vignette `simulation_study.Rmd` and `using_penetrance.Rmd`. The user must specify the `pedigree` argument as a data frame which contains the family data (see `test_fam`). The family data must be in the correct format with the following columns:

Usage

```
simulated_families
```

Format

A list of processed family data.

Details

- ID A numeric value representing the unique identifier for each individual. There should be no duplicated entries.
- Sex A numeric value where 0 indicates female and 1 indicates male. Missing entries are not currently supported.
- MotherID A numeric value representing the unique identifier for an individual's mother.
- FatherID A numeric value representing the unique identifier for an individual's father.
- isProband A numeric value where 1 indicates the individual is a proband and 0 otherwise.
- CurAge A numeric value indicating the age of censoring (current age if the person is alive or age at death if the person is deceased). Allowed ages range from 1 to 94.
- isAff A numeric value indicating the affection status of cancer, with 1 for diagnosed individuals and 0 otherwise. Missing entries are not supported.
- Age A numeric value indicating the age of cancer diagnosis, encoded as NA if the individual was not diagnosed. Allowed ages range from 1 to 94.
- Geno A column for germline testing or tumor marker testing results. Positive results should be coded as 1, negative results as 0, and unknown results as NA or left empty.

Source

Generated for package example

test_fam2

Processed Family Data

Description

A dataset containing processed information about the first simulated 130 families. These families are referenced in the vignette `simulation_study_real.Rmd`. The user must specify the pedigree argument as a data frame which contains the family data (see `test_fam`). The family data must be in the correct format with the following columns:

Usage

`test_fam2`

Format

A list of processed family data.

Details

- ID** A numeric value representing the unique identifier for each individual. There should be no duplicated entries.
- Sex** A numeric value where 0 indicates female and 1 indicates male. Missing entries are not currently supported.
- MotherID** A numeric value representing the unique identifier for an individual's mother.
- FatherID** A numeric value representing the unique identifier for an individual's father.
- isProband** A numeric value where 1 indicates the individual is a proband and 0 otherwise.
- CurAge** A numeric value indicating the age of censoring (current age if the person is alive or age at death if the person is deceased). Allowed ages range from 1 to 94.
- isAff** A numeric value indicating the affection status of cancer, with 1 for diagnosed individuals and 0 otherwise. Missing entries are not supported.
- Age** A numeric value indicating the age of cancer diagnosis, encoded as NA if the individual was not diagnosed. Allowed ages range from 1 to 94.
- Geno** A column for germline testing or tumor marker testing results. Positive results should be coded as 1, negative results as 0, and unknown results as NA or left empty.

Source

Generated for package example

transformDF

Transform Data Frame

Description

This function transforms a data frame from the standard format used in PanelPRO into the required format which conforms to the requirements of penetrance (and clipp).

Usage

```
transformDF(df)
```

Arguments

df The input data frame in the usual PanelPRO format.

Value

A data frame in the format required for clipp with the following columns:

individual	ID of the individual
isProband	Indicator if the individual is a proband
family	Family ID

mother	Mother's ID
father	Father's ID
aff	Affection status
sex	Sex (2 for female, 1 for male)
age	Age at diagnosis or current age
geno	Genotype information

Examples

```
# Create example data frame
df <- data.frame(
  ID = 1:2,
  PedigreeID = c(1,1),
  Sex = c(0,1),
  MotherID = c(NA,1),
  FatherID = c(NA,NA),
  isProband = c(1,0),
  CurAge = c(45,20),
  isAff = c(1,0),
  Age = c(40,NA),
  Geno = c(1,0)
)

# Transform the data frame
transformed_df <- transformDF(df)
```

```
validate_weibull_parameters
      Validate Weibull Parameters
```

Description

This function validates the given parameters for calculating Weibull distribution.

Usage

```
validate_weibull_parameters(
  given_first_quartile,
  given_median,
  threshold,
  asymptote
)
```

Arguments

given_first_quartile	The first quartile of the data.
given_median	The median of the data.
threshold	A constant threshold value.
asymptote	A constant asymptote value (gamma).

Value

Logical value indicating whether the parameters are valid (TRUE) or not (FALSE)

Examples

```
# Validate parameters
is_valid <- validate_weibull_parameters(
  given_first_quartile = 30,
  given_median = 50,
  threshold = 15,
  asymptote = 0.8
)
print(is_valid)
```


Index

* datasets

- baseline_data_default, 4
- distribution_data_default, 9
- out_sim, 26
- prior_params_default, 35
- risk_proportion_default, 36
- simulated_families, 36
- test_fam2, 37

- absValue, 3
- apply_burn_in, 3
- apply_thinning, 4

- baseline_data_default, 4

- calculate_weibull_parameters, 7
- calculateBaseline, 5
- calculateEmpiricalDensity, 5
- calculateNCPen, 6
- combine_chains, 8
- combine_chains_noSex, 8

- distribution_data_default, 9
- drawBaseline, 9
- drawEmpirical, 10

- generate_density_plots, 10
- generate_summary, 11
- generate_summary_noSex, 12

- imputeAges, 12
- imputeAgesInit, 15
- imputeUnaffectedAges, 16

- lik.fn, 17
- lik_noSex, 18

- makePriors, 19
- mhChain, 20
- mhLogLikelihood_clipp, 23
- mhLogLikelihood_clipp_noSex, 24

- out_sim, 26

- penetrance, 27
- penetrance-package (penetrance), 27
- plot_acf, 31
- plot_loglikelihood, 32
- plot_pdf, 32
- plot_penetrance, 33
- plot_trace, 34
- printRejectionRates, 34
- prior_params_default, 35

- qbeta, 20

- risk_proportion_default, 36
- runif, 20

- simulated_families, 36

- test_fam2, 37
- transformDF, 38

- validate_weibull_parameters, 39