

Package ‘statMatchLCM’

May 15, 2026

Type Package

Title Statistical Matching Using Latent Class Models

Version 1.2

Description

Tools for statistical matching based on latent class models. The package implements statistical matching procedures based on latent class models. It allows researchers to perform data integration when no unique identifiers are available by modeling the joint distribution of variables through latent categorical structures. The package supports estimation of latent class models, probabilistic matching between donor and recipient data sets, and generation of synthetic linked data under uncertainty. It is particularly useful in survey research and data fusion applications where combining information from multiple sources is required while preserving statistical properties and accounting for measurement error and missing data mechanisms.

License GPL-3

Encoding UTF-8

RoxygenNote 7.3.3

Imports nnet, StatMatch

Suggests NPBayesImputeCat

Depends R (>= 4.1.0)

LazyData true

NeedsCompilation no

Author Alicja Wolny-Dominiak [aut, cre],
Israa Lewaaelhamd [aut],
Mohammed Ali Ismail [aut]

Maintainer Alicja Wolny-Dominiak <woali@ue.katowice.pl>

Repository CRAN

Date/Publication 2026-05-15 20:20:08 UTC

Contents

datA	2
datAB_to_SM	3

datB	3
fact_to_num	4
num_to_fact	4
sma1_step2	5
sma2_step2	6
sma3_step2	8
sma_step1	10
smc1_step2	11
smc2_step2	12
smc3_step2	14
smc_step1	15
sm_quality	17
Index	18

 datA

Dataset dataA

Description

A simple dataset with categorical variables.

Usage

datA

Format

A data frame with 20 observations and 2 variables:

X Category (e.g., "F", "M")

Y1 Color category (e.g., "blue")

Source

Simulated data

datAB_to_SM	<i>Create stacked A/B dataset for MDC</i>
-------------	---

Description

Creates a joint data set with missing Y1/Z1, required for mass data combination.

Usage

```
datAB_to_SM(datA, datB)
```

Arguments

datA	data.frame A
datB	data.frame B

Value

data.frame with harmonized structure

Examples

```
data(datA)
data(datB)
datAB_to_SM(datA, datB)
```

datB	<i>Dataset datB</i>
------	---------------------

Description

A simple dataset with categorical variables.

Usage

```
datB
```

Format

A data frame with 15 observations and 2 variables:

X Category (e.g., "F", "M")

Z1 Color category (e.g., "red", "green", "blue")

Source

Simulated data

fact_to_num	<i>Convert factor/character variables to numeric with mapping</i>
-------------	---

Description

Converts all factor or character columns in a data.frame to numeric codes and stores the mapping tables.

Usage

```
fact_to_num(df)
```

Arguments

df data.frame with factor or character columns

Value

A list with:

data data.frame with numeric-coded variables

levels list of factor levels

tables list of mapping tables (factor → numeric)

Examples

```
data(datA)
fact_to_num(datA)
```

num_to_fact	<i>Convert numeric codes back to factor</i>
-------------	---

Description

Restores a factor variable from numeric codes using a mapping table created by fact_to_num().

Usage

```
num_to_fact(x, table)
```

Arguments

x numeric vector

table data.frame with columns level_fact and level_num

Value

A factor vector:

levels Defined by `table$level_num`.

labels Defined by `table$level_fact`.

sma1_step2

SMA1 – Nearest-neighbour hot deck on shared variables

Description

Second step of the SMA procedure using nearest-neighbour hot deck matching on shared variables (Y1, Z1) to fuse datasets while preserving their joint distribution.

Usage

```
sma1_step2(step1)
```

Arguments

`step1` list returned by the SMA step 1 procedure

Value

A list with:

datA_fused_2 The final fused dataset A.

out_nndd NNDD matching results.

References

- Lewaa, I., Hafez, M. S., Ismail, M. A. (2023). *Mixed Statistical Matching Approaches Using a Latent Class Model: Simulation Studies*. *Journal of Statistics Applications Probability*, 12(1), 247–265.
- D’Orazio, M., Di Zio, M., and Scanu, M. (2006). *Statistical Matching: Theory and Practice*. John Wiley and Sons.
- D’Orazio, M., Di-Zio, M., and Scanu, M. (2019). Auxiliary variable selection in a statistical matching problem.
- Zhang, L.-C., and Chambers, R. Analysis of integrated data. *CRC/Chapman and Hall*, pp. 101–120.
- Conti, P. L., Marella, D., and Scanu, M. (2016). Statistical matching analysis for complex survey data with applications. *Journal of the American Statistical Association*, 111(516), 1715–1725.

Examples

```

if (requireNamespace("NPBayesImputeCat", quietly = TRUE)) {
  data(datA)
  data(datB)

  datAB <- datAB_to_SM(datA, datB)
  datC <- data.frame(
    X = datB$X[1:4],
    Y1 = datA$Y1[1:4],
    Z1 = datB$Z1[1:4]
  )

  # adding auxiliary information
  datABC <- rbind(datAB, datC)

  # call DPMPM (reduced settings for speed)
  output_list_AII <- NPBayesImputeCat::DPMPM_nozeros_imp(
    X = datABC,
    nrun = 50,
    burn = 10,
    thin = 10,
    K = 20,
    aalpha = 0.25,
    balpha = 0.25,
    m = 2,
    seed = 1234,
    silent = TRUE
  )

  step_first_aii <- sma_step1(datA, datB, datC, output_list_AII)
  step_second_aii1 <- sma_step2(step_first_aii)

  result_aii1 <- step_second_aii1$datA_fused_2
  head(result_aii1)
}

```

sma2_step2

SMA2 – Hot deck on fitted multinomial probabilities

Description

Second step of SMA using multinomial models and nearest-neighbour hot deck on fitted probabilities to improve data fusion accuracy.

Usage

```
sma2_step2(step1)
```

Arguments

step1 list returned by the SMA step 1 procedure

Value

A list with:

datA_fused_2 A data frame containing the final fused (imputed) version of dataset A.

out_nnd A list with results from the nearest-neighbour distance matching procedure.

model_form A formula object specifying the model used.

FitteddatA_imp1 A data frame of fitted values obtained for dataset A.

FitteddatC_imp1 A data frame of fitted values obtained for dataset C.

References

- Lewaa, I., Hafez, M. S., Ismail, M. A. (2023). *Mixed Statistical Matching Approaches Using a Latent Class Model: Simulation Studies*. *Journal of Statistics Applications Probability*, 12(1), 247–265.
- D’Orazio, M., Di Zio, M., and Scanu, M. (2006). *Statistical Matching: Theory and Practice*. John Wiley and Sons.
- D’Orazio, M., Di~Zio, M., and Scanu, M. (2019). Auxiliary variable selection in a statistical matching problem.
- Zhang, L.-C., and Chambers, R. Analysis of integrated data. *CRC/Chapman and Hall*, pp. 101–120.
- Conti, P. L., Marella, D., and Scanu, M. (2016). Statistical matching analysis for complex survey data with applications. *Journal of the American Statistical Association*, 111(516), 1715–1725.

Examples

```
if (requireNamespace("NPBayesImputeCat", quietly = TRUE)) {
  data(dataA)
  data(datB)

  datAB <- datAB_to_SM(dataA, datB)
  datC <- data.frame(
    X = datB$X[1:4],
    Y1 = datA$Y1[1:4],
    Z1 = datB$Z1[1:4]
  )

  # adding auxiliary information
  datABC <- rbind(datAB, datC)

  output_list_AII <- NPBayesImputeCat::DPMPM_nozeros_imp(
    X = datABC,
    nrun = 50,
    burn = 10,
```

```

    thin = 10,
    K = 20,
    aalpha = 0.25,
    balpha = 0.25,
    m = 2,
    seed = 1234,
    silent = TRUE
  )

  step_first_ain <- sma_step1(datA, datB, datC, output_list_AII)
  step_second_ain2 <- sma2_step2(step_first_ain)

  result_ain2 <- step_second_ain2$datA_fused_2
  head(result_ain2)
}

```

sma3_step2

SMA3 – Multinomial simulation approach

Description

Second step of SMA is to generate the target variable using multinomial probability distributions derived from a donor-based model to achieve statistically coherent data fusion

Usage

```
sma3_step2(step1)
```

Arguments

step1 list returned by the SMA step 1 procedure

Value

A list with:

datA_fused_2 A data frame containing the final fused version of dataset A.

beta_BZ A numeric vector of estimated model coefficients.

AA_dummy A matrix or data frame of dummy variables constructed for dataset A.

BB_dummy A matrix or data frame of dummy variables constructed for dataset B.

prob_Z_all A numeric matrix of predicted probabilities for each category of variable Z (rows correspond to observations).

zero_one A binary matrix (one-hot encoded) sampled from prob_Z_all, where each row contains a single 1 indicating the selected category of Z.

References

- Lewaa, I., Hafez, M. S., Ismail, M. A. (2023). *Mixed Statistical Matching Approaches Using a Latent Class Model: Simulation Studies*. *Journal of Statistics Applications Probability*, 12(1), 247–265.
- D’Orazio, M., Di Zio, M., and Scanu, M. (2006). *Statistical Matching: Theory and Practice*. John Wiley and Sons.
- D’Orazio, M., Di~Zio, M., and Scanu, M. (2019). Auxiliary variable selection in a statistical matching problem.
- Zhang, L.-C., and Chambers, R. Analysis of integrated data. *CRC/Chapman and Hall*, pp. 101–120.
- Conti, P. L., Marella, D., and Scanu, M. (2016). Statistical matching analysis for complex survey data with applications. *Journal of the American Statistical Association*, 111(516), 1715–1725.

Examples

```

if (requireNamespace("NPBayesImputeCat", quietly = TRUE)) {
  data(datA)
  data(datB)

  datAB <- datAB_to_SM(datA, datB)
  datC <- data.frame(
    X = datB$X[1:4],
    Y1 = datA$Y1[1:4],
    Z1 = datB$Z1[1:4]
  )

  # adding auxiliary information
  datABC <- rbind(datAB, datC)

  # call DPMPM (reduced settings for speed)
  output_list_AII <- NPBayesImputeCat::DPMPM_nozeros_imp(
    X = datABC,
    nrun = 50,
    burn = 10,
    thin = 10,
    K = 20,
    aalpha = 0.25,
    balpha = 0.25,
    m = 2,
    seed = 1234,
    silent = TRUE
  )

  step_first_ain <- sma_step1(datA, datB, datC, output_list_AII)
  step_second_ain3 <- sma1_step2(step_first_ain)

  result_ain3 <- step_second_ain3$datA_fused_2
  head(result_ain3)
}

```

sma_step1	<i>SMA step 1: selection of best imputed dataset</i>
-----------	--

Description

Selects the imputed dataset minimizing the Hellinger distance between the reference distribution from dataset A and the synthetic distribution from dataset B, in a three-sample statistical matching framework (A, B, C).

Usage

```
sma_step1(datA, datB, datC, output_ll)
```

Arguments

datA	data.frame A
datB	data.frame B
datC	data.frame C
output_ll	list with imputed datasets (impdata)

Value

A list with imputed datasets:

datABC_imp1 The full imputed dataset combining A, B, and C.

datA_imp1 Subset of the imputed data corresponding to dataset A.

datB_imp1 Subset corresponding to dataset B.

datC_imp1 Subset corresponding to dataset C.

References

- Lewaa, I., Hafez, M. S., and Ismail, M. A. (2023). *Mixed Statistical Matching Approaches Using a Latent Class Model: Simulation Studies*. *Journal of Statistics Applications and Probability*, 12(1), 247–265.
- D’Orazio, M., Di Zio, M., and Scanu, M. (2006). *Statistical Matching: Theory and Practice*. John Wiley and Sons.
- D’Orazio, M., Di~Zio, M., and Scanu, M. (2019). Auxiliary variable selection in a statistical matching problem.
- Zhang, L.-C., and Chambers, R. Analysis of integrated data. CRC/Chapman and Hall, pp. 101–120.
- Conti, P. L., Marella, D., and Scanu, M. (2016). Statistical matching analysis for complex survey data with applications. *Journal of the American Statistical Association*, 111(516), 1715–1725.

Examples

```

if (requireNamespace("NPBayesImputeCat", quietly = TRUE)) {
  data(datA)
  data(datB)

  datAB <- datAB_to_SM(datA, datB)
  datC <- data.frame(
    X = datB$X[1:4],
    Y1 = datA$Y1[1:4],
    Z1 = datB$Z1[1:4]
  )

  # adding auxiliary information
  datABC <- rbind(datAB, datC)

  # call DPMPM
  output_list_AII <- NPBayesImputeCat::DPMPM_nozeros_imp(
    X = datABC, nrun = 500, burn = 50, thin = 50,
    K = 80, aalpha = 0.25, balpha = 0.25,
    m = 2, seed = 1234, silent = FALSE
  )

  step_first_ain <- sma_step1(datA, datB, datC, output_list_AII)
  names(step_first_ain)
}

```

smc1_step2

SMC1 - Nearest-neighbour hot deck on original variables

Description

Performs nearest-neighbour hot deck imputation on the selected imputed datasets by matching recipient units in dataset A with donor units in dataset B based on common variables. The procedure transfers the target variable from the nearest donor to construct a statistically coherent fused dataset.

Usage

```
smc1_step2(step1)
```

Arguments

```
step1          output from smc_step1()
```

Value

A list with:

datA_fused_2 The final fused dataset A after step 2 of the procedure.

out_nndd A list containing nearest-neighbour distance matching results.

References

- Lewaa, I., Hafez, M. S., and Ismail, M. A. (2023). *Mixed Statistical Matching Approaches Using a Latent Class Model: Simulation Studies*. *Journal of Statistics Applications and Probability*, 12(1), 247–265.
- D’Orazio, M., Di Zio, M., and Scanu, M. (2006). *Statistical Matching: Theory and Practice*. John Wiley and Sons.
- D’Orazio, M., Di~Zio, M., and Scanu, M. (2019). Auxiliary variable selection in a statistical matching problem.
- Zhang, L.-C., and Chambers, R. Analysis of integrated data. CRC/Chapman and Hall, pp. 101–120.
- Conti, P. L., Marella, D., and Scanu, M. (2016). Statistical matching analysis for complex survey data with applications. *Journal of the American Statistical Association*, 111(516), 1715–1725.

Examples

```

if (requireNamespace("NPBayesImputeCat", quietly = TRUE)) {
  data(datA)
  data(datB)

  datAB <- datAB_to_SM(datA, datB)

  # call DPMPM (reduced settings for speed)
  output_list <- NPBayesImputeCat::DPMPM_nozeros_imp(
    X = datAB,
    nrun = 50,
    burn = 10,
    thin = 10,
    K = 20,
    aalpha = 0.25,
    balpha = 0.25,
    m = 2,
    seed = 1234,
    silent = TRUE
  )

  step_first <- smc_step1(datA, datB, output_list)
  step_second <- smc1_step2(step_first)

  result <- step_second$data_fused_2
  result
}

```

Description

Implements a model-based hot deck data fusion approach by estimating multinomial models on both datasets and matching observations based on their fitted probability distributions

Usage

```
smc2_step2(step1)
```

Arguments

step1 output from smc_step1()

Value

A list with:

datA_fused_2 A data frame containing the final fused (imputed) version of dataset A.

out_nndd A list with results from the nearest-neighbour distance matching procedure.

model_form A formula object used to fit the models.

FitteddatA_imp1 A data frame of fitted values obtained from the model for dataset A.

FitteddatB_imp1 A data frame of fitted values obtained from the model for dataset B.

References

- Lewaa, I., Hafez, M. S., and Ismail, M. A. (2023). *Mixed Statistical Matching Approaches Using a Latent Class Model: Simulation Studies*. *Journal of Statistics Applications and Probability*, 12(1), 247–265.
- D’Orazio, M., Di Zio, M., and Scanu, M. (2006). *Statistical Matching: Theory and Practice*. John Wiley and Sons.
- D’Orazio, M., Di~Zio, M., and Scanu, M. (2019). Auxiliary variable selection in a statistical matching problem.
- Zhang, L.-C., and Chambers, R. Analysis of integrated data. CRC/Chapman and Hall, pp. 101–120.
- Conti, P. L., Marella, D., and Scanu, M. (2016). Statistical matching analysis for complex survey data with applications. *Journal of the American Statistical Association*, 111(516), 1715–1725.

Examples

```
if (requireNamespace("NPBayesImputeCat", quietly = TRUE)) {
  data(datA)
  data(datB)

  datAB <- datAB_to_SM(datA, datB)

  # call DPMPM (reduced settings for speed)
  output_list <- NPBayesImputeCat::DPMPM_nozeros_imp(
    X = datAB,
```

```

    nrun = 50,
    burn = 10,
    thin = 10,
    K = 20,
    aalpha = 0.25,
    balpha = 0.25,
    m = 2,
    seed = 1234,
    silent = TRUE
  )

  step_first <- smc_step1(datA, datB, output_list)
  step_second <- smc2_step2(step_first)

  result <- step_second$datA_fused_2
  result
}

```

smc3_step2

SMC3 – Multinomial simulation approach

Description

Generates the target variable using multinomial simulation based on fitted probability distributions to preserve uncertainty and variability

Usage

```
smc3_step2(step1)
```

Arguments

step1 output from smc_step1()

Value

A list with:

datA_fused_2 A data frame containing the final fused (imputed) version of dataset A.

FitteddatA_imp1 A data frame of fitted values obtained from the model for dataset A.

References

- Lewaa, I., Hafez, M. S., and Ismail, M. A. (2023). *Mixed Statistical Matching Approaches Using a Latent Class Model: Simulation Studies*. *Journal of Statistics Applications and Probability*, 12(1), 247–265.
- D’Orazio, M., Di Zio, M., and Scanu, M. (2006). *Statistical Matching: Theory and Practice*. John Wiley and Sons.

- D’Orazio, M., Di-Zio, M., and Scanu, M. (2019). Auxiliary variable selection in a statistical matching problem.
- Zhang, L.-C., and Chambers, R. Analysis of integrated data. CRC/Chapman and Hall, pp. 101–120.
- Conti, P. L., Marella, D., and Scanu, M. (2016). Statistical matching analysis for complex survey data with applications. *Journal of the American Statistical Association*, 111(516), 1715–1725.

Examples

```

if (requireNamespace("NPBayesImputeCat", quietly = TRUE)) {
  data(dataA)
  data(datB)

  datAB <- datAB_to_SM(dataA, datB)

  # call DPMPM (reduced settings for speed)
  output_list <- NPBayesImputeCat::DPMPM_nozeros_imp(
    X = datAB,
    nrun = 50,
    burn = 10,
    thin = 10,
    K = 20,
    aalpha = 0.25,
    balpha = 0.25,
    m = 2,
    seed = 1234,
    silent = TRUE
  )

  step_first <- smc_step1(dataA, datB, output_list)
  step_second <- smc3_step2(step_first)

  result <- step_second$data_fused_2
  result
}

```

smc_step1

SMC step 1: selection of best imputed dataset

Description

Identifies and selects the imputed dataset that minimizes the Hellinger distance between the reference and synthetic distributions, thereby achieving the highest level of distributional similarity and improving the statistical consistency of the imputation.

Usage

```
smc_step1(datA, datB, output_ll)
```

Arguments

datA	data.frame A
datB	data.frame B
output_ll	list with imputed datasets (impdata)

Value

A list with imputed datasets:

datAB_imp1 The full imputed dataset obtained after combining A and B.

datA_imp1 Subset of datAB_imp1 corresponding to dataset A.

datB_imp1 Subset of datAB_imp1 corresponding to dataset B.

References

- Lewaa, I., Hafez, M. S., Ismail, M. A. (2023). *Mixed Statistical Matching Approaches Using a Latent Class Model: Simulation Studies*. *Journal of Statistics Applications Probability*, 12(1), 247–265.
- D’Orazio, M., Di Zio, M., and Scanu, M. (2006). *Statistical Matching: Theory and Practice*. John Wiley and Sons.
- D’Orazio, M., Di~Zio, M., and Scanu, M. (2019). Auxiliary variable selection in a statistical matching problem.
- Zhang, L.-C., and Chambers, R. Analysis of integrated data. *CRC/Chapman and Hall*, pp. 101–120.
- Conti, P. L., Marella, D., and Scanu, M. (2016). Statistical matching analysis for complex survey data with applications. *Journal of the American Statistical Association*, 111(516), 1715–1725.

Examples

```
if (requireNamespace("NPBayesImputeCat", quietly = TRUE)) {
  data(datA)
  data(datB)

  datAB <- datAB_to_SM(datA, datB)

  # call DPMPM (reduced settings for speed)
  output_list <- NPBayesImputeCat::DPMPM_nozeros_imp(
    X = datAB,
    nrun = 50,
    burn = 10,
    thin = 10,
    K = 20,
    aalpha = 0.25,
    balpha = 0.25,
    m = 2,
    seed = 1234,
    silent = TRUE
  )
}
```



```
)  
  
  step_first <- smc_step1(datA, datB, output_list)  
  str(step_first$datA_imp1)  
  str(step_first$datB_imp1)  
}
```

sm_quality

Quality assessment of synthetic ZI

Description

Evaluates the quality of the synthetic target variable by computing the Hellinger distance between the reference and synthetic distributions. This measure quantifies the degree of similarity between the two distributions, providing an assessment of the accuracy and coherence of the data fusion process

Usage

```
sm_quality(step1, step2)
```

Arguments

step1	output from smc_step1()
step2	output from step 2 method

Value

A list with:

heli_latent A numeric value representing the Hellinger distance between the reference and synthetic distributions.

ref_distr A numeric vector or table representing the reference (original) distribution.

synt_distr A numeric vector or table representing the synthetic (generated) distribution.

Index

* datasets

datA, 2

datB, 3

datA, 2

datAB_to_SM, 3

datB, 3

fact_to_num, 4

num_to_fact, 4

sm_quality, 17

sma1_step2, 5

sma2_step2, 6

sma3_step2, 8

sma_step1, 10

smc1_step2, 11

smc2_step2, 12

smc3_step2, 14

smc_step1, 15