

---

# **PC method to impute missing values for mixed data**

*Release 2024*

**Ambre Adjevi and Laura Fuentes**

**Mar 23, 2024**



**CONTENTS:**

<b>1</b>	<b>algorithms module</b>	<b>1</b>
<b>2</b>	<b>utils module</b>	<b>3</b>
<b>3</b>	<b>metrics_FAMD module</b>	<b>5</b>
<b>4</b>	<b>Indices and tables</b>	<b>7</b>
	<b>Python Module Index</b>	<b>9</b>
	<b>Index</b>	<b>11</b>



## ALGORITHMS MODULE

```
class algorithms.FAMD(data, k1, k2, nb_values_per_cat, n_components=2)
```

Bases: object

**DM()**

Function to define D and M according to current values

**data\_concat()**

Redefines the whole dataset (df) after updating df\_C0 (continuous variables) and df\_categ (categorical variables) with new imputation

**ponderation\_gsbs()**

Weighting step specialized in the gsbs dataset: This function updates: - the standard deviation values (sj) for continuous variables - proportion for categorical variables (pj) according to a new imputation update.

**run\_famd**(verbose=False)

Function that runs the whole FAMD algorithm

**Returns**

reconstructed version of (self.XD\_moins\_sqrt - self.M) according to self.n\_components

**Return type**

Z\_p (pd.DataFrame)

**step3()**

Performs the third step from FAMD algorithm: - Update D, M - Computes the SVD on self.XD\_moins\_sqrt - self.M

**Returns**

reconstructed version of (self.XD\_moins\_sqrt - self.M) according to self.n\_components

**Return type**

Z\_p (pd.DataFrame)

```
class algorithms.IterativeFAMDImputer(data, k1, k2, nb_values_per_cat, n_components=2)
```

Bases: *FAMD*

**impute**(n\_it, tol=0.0001, verbose=False)

Runs the whole imputation process

**Parameters**

- **n\_it** (*int*) – number of iterations to run the algorithm
- **tol** (*float*, *optional*) – Threshold to define early stopping criteria. Defaults to 1e-4.
- **verbose** (*bool*, *optional*) – Defaults to False.

**Returns**

Imputed dataset after algorithm convergence

**Return type**

self.df (pd.DataFrame)

**inital\_impute()**

Initial imputation: - continuous variables: mean/variable - categories: proportion/category

**ponderation\_gsbs()**

Weighting function specialized on the gsbs dataset This function updates: - the standard deviation values (sj) for continous variables, - proportion for categorical variables (pj) according to a new imputation update.

## UTILS MODULE

`utils.compute_metrics(df, cat_idx, n_it, n_components, proba_non_missing)`

Computes the falsely classified rate and nrmse over a synthetic dataset for different probabilities of missingness.

### Parameters

- **df** (*pd.DataFrame*) – dataframe to impute
- **n\_it** (*int*) – maximum number of iterations for iFAMD convergence
- **n\_components** (*int*) – number of principal components for reconstruction
- **proba\_non\_missing** (*list of float*) – List of probabilities that a value is not missing

### Returns

Masked dataframe hence containing missing values

### Return type

`data_missing_raw` (*pd.DataFrame*)

`utils.create_dataset(n, S, K, cat, cat_idx, nb_of_cat_per_var, SNR=1.0)`

Generate dataset with continuous (float variables) and categorical variables (string variables) based on parameters. :param n: Sample size. :type n: int :param S: Underlying dimension. :type S: int :param K: K[s] = number of times the variable s (s in {1,...,S}) is duplicated in the dataset :type K: list of ints :param cat: Number of categorical variables :type cat: int :param cat\_idx: Indexes of the categorical variables :type cat\_idx: list of ints :param nb\_of\_cat\_per\_var: Number of categories for each categorical variable :type nb\_of\_cat\_per\_var: list of ints :param SNR: Signal to Noise Ratio :type SNR: float

### Returns

Generated dataset

### Return type

`df` (*pd.DataFrame*)

`utils.create_missingness(df, proba_non_missing)`

Returns dataframe df with missing values

### Parameters

- **df** (*pd.DataFrame*) – Dataframe to mask
- **proba\_non\_missing** (*float*) – probability that a value is not missing

### Returns

masked dataframe hence containing missing values

### Return type

`data_missing_raw` (*pd.DataFrame*)

`utils.create_rare_df(f, n, S=1, K=[4], SNR=5)`

Create dataset with rare categories, as described in section 3.3 of the paper

**Parameters**

- **f** (*float*) – frequency of the rare categories
- **n** (*int*) – sample size
- **S** (*int*) – Underlying dimension.
- **K** (*list of ints*) –  $K[s]$  = number of times the variable  $s$  ( $s$  in  $\{1, \dots, S\}$ ) is duplicated in the dataset
- **SNR** (*float*) – Signal to Noise Ratio

**Returns**

Generated dataset with rare categories

**Return type**

`df_rare` (`pd.DataFrame`)

`utils.encode_dummy_variables(df, cat_var_idx)`

Encode dummy variables, returns the updated dataframe, the list of the dummy variables' names and the list of the number of different values in each categories Format of the dummy variable names : `name_of_variable + _ + variable_value`

**Args :**

`df` (`pd.DataFrame`) `cat_var` (`pd.Index`): index of the variables to encode into dummy variables

**Returns**

encoded dummy variables `dummy_var_idx` (`pd.Index`): list of the dummy variables' names  
`nb_values_per_cat` (`list`): number of different values in each categories

**Return type**

`df_dummy` (`pd.DataFrame`)

`utils.generate_random(column, missing_indices_len)`

Generate a random number between the minimum and maximum of a given variable

**Parameters**

**column** (`pd.Series`) – Column values

**Returns**

`np.random.randint(low=column.min(), high=column.max())`

**Return type**

random value (`int`)

`utils.random_imputation(df)`

Perform random imputation with random values per column

**Parameters**

**df** (`pd.DataFrame`) – `_description_`

**Returns**

Imputed dataset

**Return type**

`df_modif` (`pd.DataFrame`)



## METRICS\_FAMD MODULE

`metrics_FAMD.compute_nrmse_weighted(dfpred, dftrue)`

Computes de Normalized Root Mean Squared Error (NRMSE)

**Parameters**

- **dfpred** (*pd.DataFrame*) – dataframe gathering for continuous variables
- **dftrue** (*pd.DataFrame*) – ground truth dataframe gathering for continuous variables

**Returns**

NRMSE value for each variable

**Return type**

nrmse (float)

`metrics_FAMD.metric_fc(df_categ, true_df_categ)`

Computes the proportion of falsely classified individuals per category.

**Parameters**

- **df\_categ** (*pd.DataFrame*) – dataframe gathering only categorical variables
- **true\_df\_categ** (*pd.DataFrame*) – ground truth dataframe gathering for categorical variables

**Returns**

rate of falsely classified individuals per category

**Return type**

fc (float)



## INDICES AND TABLES

- `genindex`
- `modindex`
- `search`



## PYTHON MODULE INDEX

### a

algorithms, 1

### m

metrics\_FAMD, 5

### u

utils, 3



## INDEX

### A

algorithms  
    module, 1

### C

compute\_metrics() (in module utils), 3  
compute\_nrmse\_weighted() (in module metrics\_FAMD), 5  
create\_dataset() (in module utils), 3  
create\_missingness() (in module utils), 3  
create\_rare\_df() (in module utils), 3

### D

data\_concat() (algorithms.FAMD method), 1  
DM() (algorithms.FAMD method), 1

### E

encode\_dummy\_variables() (in module utils), 4

### F

FAMD (class in algorithms), 1

### G

generate\_random() (in module utils), 4

### I

impute() (algorithms.IterativeFAMDImputer method), 1  
inital\_impute() (algorithms.IterativeFAMDImputer method), 2  
IterativeFAMDImputer (class in algorithms), 1

### M

metric\_fc() (in module metrics\_FAMD), 5  
metrics\_FAMD  
    module, 5  
module  
    algorithms, 1  
    metrics\_FAMD, 5  
    utils, 3

### P

ponderation\_gsbs() (algorithms.FAMD method), 1  
ponderation\_gsbs() (algorithms.IterativeFAMDImputer method), 2

### R

random\_imputation() (in module utils), 4  
run\_famd() (algorithms.FAMD method), 1

### S

step3() (algorithms.FAMD method), 1

### U

utils  
    module, 3