

# Missing Data Imputation Methods for Food Composition Data

\*Gordana Ispirova<sup>1,2</sup>, Tome Eftimov<sup>1</sup>, Barbara Koroušič Seljak<sup>1</sup>

## Objective

**Scenario:** Nutrient-intake assessment, dietary guidelines

**Problem:** Missing nutrient data

**Where:** Nutrition domain

**Target:** Food composition databases (FCDBs)

**Solution:** Missing data imputation techniques

## FCDBs?

Detailed sets of information on the nutritionally important food components.

Nutrient/Food	Apple	Pear	Peach
Na	?	1	?
K	107	116	190
Mg	5	7	?

## Method

Missing data imputation techniques:

1. Fill-in with mean
2. Fill-in with median
3. Non-Negative Matrix Factorization (NMF)
4. Multiple Imputations by Chained Equations (MICE)
5. Nonparametric Missing Value Imputation using Random Forest (MissForest)
6. K-Nearest Neighbors (KNN)

FCDBs are quite incomplete, obtaining complete datasets was not an easy task, and the datasets obtained are small!

## Data?

From the national FCDBs of 10 countries, collected by EuroFIR

Nutritional values for foods from four food groups:

- Potassium in Fruits
- Sodium in Fruits
- Sodium in Vegetables
- Protein in Meat

Nutrient <sub>a</sub>	Country <sub>1</sub>	Country <sub>2</sub>	...	Country <sub>m</sub>
Food <sub>1</sub>	Value <sub>11</sub>	Value <sub>12</sub>	...	Value <sub>1m</sub>
⋮	⋮	⋮	...	⋮
Food <sub>n</sub>	Value <sub>n1</sub>	Value <sub>n2</sub>	...	Value <sub>nm</sub>

## Results and conclusions

**Evaluation:**

1. 10%, 20% and 30% of the data is set as missing
2. Calculate the missing values
3. Compare obtained values with actual values

**Evaluation criteria:**

- Mean Absolute Percentage Error (MAPE)
- Mean Kullback-Liebler Divergence Error (MKLDE)
- Relative Absolute Error (RAE)
- Root Mean Square Error (RMSE)

## Best technique?

Significantly better results than the classical fill-in with mean or median techniques.

Overall: MICE and Miss Forest yield the best results and deserve further consideration in practice.

The graphs are generated from the evaluation of the dataset with values for potassium in Fruits.

