

# Imputation of Missing Data Items under Linear Restrictions

Caren Tempelman

Statistics Netherlands



## Problem definition

Economic data consist of many linear constraints, such as:

$$c_1X_1 + c_2X_2 + \cdots + c_kX_k = X_{k+1}$$

If several  $X_i$ 's are missing, how do we impute these items immediately satisfying the linear constraint and preserving the distribution of the data?



## Approach

The linear constraint can be transformed as follows in order to restrict the domain of the variables to the simplex:

$$\frac{c_1 X_1}{X_{k+1}} + \frac{c_2 X_2}{X_{k+1}} + \dots + \frac{c_k X_k}{X_{k+1}} = 1 \quad \Rightarrow$$
$$\tilde{X}_1 + \tilde{X}_2 + \dots + \tilde{X}_k = 1$$



# The Dirichlet distribution and its advantages

$$f(x_1, \dots, x_k | \alpha_1, \dots, \alpha_k) = \frac{\Gamma(\sum_{i=1}^k \alpha_i)}{\prod_{i=1}^k \Gamma(\alpha_i)} \prod_{i=1}^k x_i^{\alpha_i - 1}$$

where  $x_i \geq 0$ ,  $i = 1, \dots, k$ , and  $x_1 + \dots + x_k = 1$

with parameters  $\alpha_i > 0$ ,  $i = 1, \dots, k$ .

## Advantages

- Extremely flexible in the shapes it will accommodate
- On the simplex

