

Expectation-Maximization model for substitution of missing values characterizing greenness of organic solvents

Gabriela Łuczynska ^{1,2} Francisco Pena-Pereira ³, Marek Tobiszewski ⁴ and Jacek Namieśnik* ⁴

¹ Division of Applied Mathematics and Probability, Institute of Mathematics, Faculty of Mathematics, University of Gdansk, 8 J. Bażyńskiego St., 80-309 Gdańsk, Poland; gabluczy@student.pg.edu.pl

² Department of Nonlinear Analysis and Statistics, Faculty of Applied Mathematics, Gdańsk University of Technology (GUT), 11/12 G. Narutowicza St., 80-233 Gdańsk, Poland

³ Department of Analytical and Food Chemistry, Faculty of Chemistry, University of Vigo, Campus As Lagoas-Marcosende s/n, 36310 Vigo, Spain; fjpna@uvigo.es

⁴ Department of Analytical Chemistry, Chemical Faculty, Gdańsk University of Technology (GUT), 11/12 G. Narutowicza St., 80-233 Gdańsk, Poland; marektobiszewski@wp.pl

* Correspondence: jacek.namiesnik@pg.edu.pl or chemanal@pg.edu.pl

E-M algorithm description

To complete the data we use E-M algorithm. This algorithm consists of two steps: an Expectation step or the E-step and a Maximization step or the M-step.

We observe a data vector \mathbf{y} . Let \mathbf{Y} be the random vector corresponding to the data \mathbf{y} . Let Θ be a parameter space and let \mathbf{Y} has the probability distribution function $g(\mathbf{y}; \Psi)$ where Ψ is a vector of unknown parameters from Θ .

Let \mathbf{X} be a random vector corresponding to a complete-data vector \mathbf{x} . Denote the probability distribution function of \mathbf{X} by $g_c(\mathbf{x}; \Psi)$. Let \mathbf{A} and \mathbf{B} be two samples spaces in which we observe the data \mathbf{x} and \mathbf{y} respectively. We don't observe the complete vector \mathbf{x} in \mathbf{A} but only the incomplete vector $\mathbf{y} = \mathbf{y}(\mathbf{x})$ in \mathbf{B} . Therefore we have a many-to-one mapping from \mathbf{A} to \mathbf{B} and due to the disintegration theorem [28]

$$g(\mathbf{y}; \Psi) = \int_{A(\mathbf{y})} g_c(\mathbf{x}; \Psi) dx, \quad (1)$$

where $A(\mathbf{y})$ is a subset of \mathbf{A} determined by the equation $\mathbf{y} = \mathbf{y}(\mathbf{x})$. The likelihood function for Ψ formed from the observed data \mathbf{y} is given by

$$L(\Psi) = g(\mathbf{y}; \Psi). \quad (2)$$

An estimate $\hat{\Psi}$ of Ψ can be obtained by solving the log likelihood equation

$$\frac{\partial \log L(\Psi)}{\partial \Psi} = 0, \quad (3)$$

where \log is the natural logarithm function. The E-M algorithm approaches the problem of solving the incomplete-data log likelihood equation (3) indirectly by proceeding iteratively in terms of complete-data log likelihood function $\log L_c(\Psi)$, where

$$L_c(\Psi) = g_c(\mathbf{x}; \Psi). \quad (4)$$

Since it is unobservable, it is replaced by its conditional expected value given \mathbf{y} , using the current fit for Ψ .

Let $\Psi^{(0)}$ be some initial value for Ψ . Then on the first iteration in the E-step we calculate

$$Q(\Psi, \Psi^{(0)}) := E_{\Psi^{(0)}}(\log L_c(\Psi) | \mathcal{Y}). \quad (5)$$

After that, in the M-step we maximize $Q(\Psi, \Psi^{(0)})$ with respect to Ψ over the parameter space Θ . In other words we choose such $\Psi^{(1)}$ that

$$\forall \Psi \in \Theta \quad Q(\Psi^{(1)}, \Psi^{(0)}) \geq Q(\Psi, \Psi^{(0)}). \quad (6)$$

Note that this choice doesn't have to be unique.

Thereafter steps are performed again with value $\Psi^{(1)}$ in place of $\Psi^{(0)}$. The $(k + 1)$ th iteration of the E-M algorithm requires in the E-step the calculation of

$$Q(\Psi, \Psi^{(k)}) := E_{\Psi^{(k)}}(\log L_c(\Psi) | \mathcal{Y}). \quad (7)$$

and in the M-step the choice of $\Psi^{(k+1)} \in \Theta$ which maximizes $Q(\Psi, \Psi^{(k)})$ with respect to $\Psi \in \Theta$, i.e.

$$\forall \Psi \in \Theta \quad Q(\Psi^{(k+1)}, \Psi^{(k)}) \geq Q(\Psi, \Psi^{(k)}). \quad (8)$$

The steps are carried out until the value of

$$L(\Psi^{(k+1)}) - L(\Psi^{(k)}) \quad (9)$$

is smaller than an arbitrarily amount in case of convergence of the sequence of likelihood values $(L(\Psi^{(k)}))_k$. It has been shown in [29] that this sequence is nondecreasing. Therefore it is convergent, if it is bounded above.