

Maximum Likelihood from Evidential Data: An Extension of the EM Algorithm

Thierry Dencœux

Abstract. We consider the problem of statistical parameter estimation when the data are uncertain and described by belief functions. An extension of the Expectation-Maximization (EM) algorithm, called the Evidential EM (E^2M) algorithm, is described and shown to maximize a generalized likelihood function. This general procedure provides a simple mechanism for estimating the parameters in statistical models when observed data are uncertain. The method is illustrated using the problem of univariate normal mean and variance estimation from uncertain data.

Keywords: Belief functions, Dempster-Shafer theory, Statistical inference, Uncertain data.

1 Introduction

In statistics, observations of random quantities are usually assumed to be either precise or imprecise, i.e., set-valued. The latter situation occurs, e.g., in the case of censored data, where an observation is only known to belong to a set, usually an interval. The Expectation-Maximization (EM) algorithm [4, 8] has proved to be a powerful mechanism for performing maximum likelihood parameter estimation from such incomplete data.

There are situations, however, where the observations are not only imprecise, but also *uncertain*, i.e., partially reliable [1]. Consider, e.g., a classification problem in which objects in a population belong to one and only one group. Let \mathcal{X} be the finite set of groups, and X be the group of an object randomly drawn from the population. In some applications, realizations x of X are not known with certainty. Rather, an expert provides a subjective

Thierry Dencœux

Heudiasyc, Université de Technologie de Compiègne, CNRS, Compiègne, France
e-mail: tdencœux@hds.utc.fr

assessment of x (a process known as *labeling*). This assessment may take the form of a subset $A \subseteq \mathcal{X}$, a probability distribution p on \mathcal{X} or, more generally, a mass function m on \mathcal{X} , i.e., a function $m: 2^{\mathcal{X}} \rightarrow [0, 1]$. It must be stressed that, in this example, the data generation process can be decomposed into two components: a random component, which generates a realization x from X , and a non random component, which produces a mass function m that models the expert's partial knowledge of x .

If this process is repeated n times independently, the data takes the form of n mass functions m_1, \dots, m_n , considered as a partial specification of an unknown realization x_1, \dots, x_n of an i.i.d. random sample X_1, \dots, X_n . We will refer to such data as *evidential data*. If a parametric model is postulated for X , how can the method of maximum likelihood be extended to handle such data? This is the problem considered in this paper. A generalization of the likelihood function will be proposed, and an extension of the EM algorithm, called the evidential EM (E^2M) algorithm, will be introduced for its maximization.

We may note that, in the special case where each mass functions m_i is consonant, the data can be equivalently represented as n possibility distribution $\tilde{x}_1, \dots, \tilde{x}_n$, which constitutes a *fuzzy random sample*. The problem of statistical inference from fuzzy data, which has received a lot of attention in the past few years [5, 6], is thus a special case of the problem considered here.

Early attempts to adapt the EM algorithm to evidential data, in the special case of mixture models with evidential class labels, were presented in [10, 7]. A rigorous solution to this problem, which is a special case of the general method presented in this paper, was introduced in [2].

The rest of the paper is organized as follows. The EM algorithm will first be recalled in Section 2. The extension of the likelihood function and the E^2M algorithm will then be introduced in Sections 3 and 4, respectively. Section 5 will demonstrate the application of this algorithm to the problem of univariate normal mean and variance estimation using uncertain data.

2 The EM Algorithm

The EM algorithm is a broadly applicable mechanism for computing MLEs from incomplete data, in situations where ML estimation would be straightforward if complete data were available [4]. Formally, we assume the existence of two sample spaces \mathcal{X} and \mathcal{Y} , and a many-to-one mapping φ from \mathcal{X} to \mathcal{Y} . The observed (incomplete) data \mathbf{y} are a realization from \mathcal{Y} , while the corresponding \mathbf{x} in \mathcal{X} is not observed and is only known to lie in the set

$$\mathcal{X}(\mathbf{y}) = \varphi^{-1}(\mathbf{y}) = \{\mathbf{x} \in \mathcal{X} | \varphi(\mathbf{x}) = \mathbf{y}\}.$$

Vector \mathbf{x} is referred to as the *complete data* vector. It is a realization from a random vector \mathbf{X} with p.d.f. $g_c(\mathbf{x}; \Psi)$, where $\Psi = (\Psi_1, \dots, \Psi_d)'$ is a vector of unknown parameters with parameter space Ω . The observed data likelihood $L(\Psi)$ is related to $g_c(\mathbf{x}; \Psi)$ by

$$L(\Psi) = \int_{\mathcal{X}(\mathbf{y})} g_c(\mathbf{x}; \Psi) d\mathbf{x}. \quad (1)$$

The EM algorithm approaches the problem of maximizing the observed-data log likelihood $\log L(\Psi)$ by proceeding iteratively with the complete-data log likelihood $\log L_c(\Psi) = \log g_c(\mathbf{x}; \Psi)$. Each iteration of the algorithm involves two steps called the expectation step (E-step) and the maximization step (M-step).

The E-step requires the calculation of

$$Q(\Psi, \Psi^{(q)}) = \mathbb{E}_{\Psi^{(q)}} [\log L_c(\Psi) | \mathbf{y}],$$

where $\Psi^{(q)}$ denotes the current fit of Ψ at iteration q , and $\mathbb{E}_{\Psi^{(q)}}$ denotes expectation using the parameter vector $\Psi^{(q)}$.

The M-step then consists in maximizing $Q(\Psi, \Psi^{(q)})$ with respect to Ψ over the parameter space Ω . The E- and M-steps are iterated until the difference $L(\Psi^{(q+1)}) - L(\Psi^{(q)})$ becomes smaller than some arbitrarily small amount.

3 Generalized Likelihood Function

Let us now consider the more complex situation where the relationship between the observed and complete spaces is uncertain, so that observed data \mathbf{y} can no longer be associated with certainty to a unique subset of \mathcal{X} . This situation will be formalized as follows.

Let us assume the existence of a set Θ of interpretations, one and only one of which holds, and a probability measure \Pr on Θ . If \mathbf{y} has been observed and $\theta \in \Theta$ is the true interpretation, then the complete data \mathbf{x} is known to belong to $\mathcal{X}(\mathbf{y}, \theta) \subseteq \mathcal{X}$. Having observed \mathbf{y} , the probability measure \Pr is carried to $2^{\mathcal{X}}$ by the mapping $\theta \rightarrow \mathcal{X}(\mathbf{y}, \theta)$, which defines a Dempster-Shafer mass function m on \mathcal{X} . For simplicity, we will assume from now on that Θ is finite: $\Theta = \{\theta_1, \dots, \theta_K\}$, in which case m is a discrete mass function with focal sets $\mathcal{X}_k = \mathcal{X}(\mathbf{y}, \theta_k)$ and masses $m_k = m(\mathcal{X}_k) = \Pr(\{\theta_k\})$ for $k = 1, \dots, K$.

With the same notations as in the previous section, the observed data likelihood may now be defined as:

$$\begin{aligned} L(\Psi) &= \sum_{k=1}^K m_k \int_{\mathcal{X}_k} g_c(\mathbf{x}; \Psi) d\mathbf{x} = \int_{\mathcal{X}} g_c(\mathbf{x}; \Psi) \left(\sum_{k=1}^K m_k 1_{\mathcal{X}_k}(\mathbf{x}) \right) d\mathbf{x} \\ &= \int_{\mathcal{X}} g_c(\mathbf{x}; \Psi) pl(\mathbf{x}) d\mathbf{x} = \mathbb{E}_{\Psi} [pl(\mathbf{X})], \quad (2) \end{aligned}$$

where $pl: \mathcal{X} \rightarrow [0, 1]$ is the contour function associated to m .

The generalized likelihood of Ψ is thus equal to the expectation of the plausibility contour function, with respect to the probability distribution $g_c(\mathbf{x}; \Psi)$. We can remark that, when m is consonant, the contour function can be seen as the membership function of a fuzzy subset of \mathcal{X} : $L(\Psi)$ is then the

probability of that fuzzy subset, according to Zadeh's definition of the probability of a fuzzy event [11].

In the more general setting of belief functions, $L(\Psi)$ has another interpretation that will now be explained. Let $g_c(\cdot|m; \Psi) = m \oplus g_c(\cdot; \Psi)$ denote the p.d.f. obtained by combining m with the complete data p.d.f. $g_c(\cdot; \Psi)$ using Dempster's rule [3, 9]:

$$g_c(\mathbf{x}|m; \Psi) = \frac{g_c(\mathbf{x}; \Psi)pl(\mathbf{x})}{\int_{\mathcal{X}} g_c(\mathbf{u}; \Psi)pl(\mathbf{u})d\mathbf{u}} = \frac{g_c(\mathbf{x}; \Psi)pl(\mathbf{x})}{L(\Psi)}. \quad (3)$$

The normalizing constant $L(\Psi)$ at the denominator of the above expression is equal to one minus the degree of conflict between m and $g_c(\mathbf{x}; \Psi)$. Consequently, maximizing $L(\Psi)$ amounts to *minimizing the conflict* between the observations (represented by m) and the parametric model $g_c(\cdot; \Psi)$.

The expression of the observed data likelihood (2) can often be simplified by making independence assumptions. Let us assume that the observed data $\mathbf{x} = (x_1, \dots, x_n)$ is a realization from a random vector $\mathbf{X} = (X_1, \dots, X_n)$. In many applications, we can make the following assumptions:

A1: Stochastic independence of the r.v. X_1, \dots, X_n :

$$g_c(\mathbf{u}; \Psi) = \prod_{i=1}^n g_c(u_i; \Psi), \quad \forall \mathbf{u} = (u_1, \dots, u_n) \in \mathcal{X}.$$

A2: The plausibility contour function $pl(\mathbf{x})$ can be written as

$$pl(\mathbf{u}) = \prod_{i=1}^n pl_i(u_i), \quad \forall \mathbf{u} = (u_1, \dots, u_n) \in \mathcal{X},$$

where pl_i is the contour function corresponding to the marginal mass function m_i on x_i .

It should be noted that Assumption A2 is totally unrelated to A1: it is not a property of the random variables X_1, \dots, X_n , but of the uncertain observation process. It is actually a weaker form of the *cognitive independence* assumption, as defined by Shafer [9].

Under Assumptions A1 and A2, the observed data log likelihood can be written as a sum of n terms:

$$\log L(\Psi) = \sum_{i=1}^n \log \mathbb{E}_{\Psi} [pl_i(X_i)].$$

4 The Evidential EM Algorithm

To maximize function $L(\Psi)$ defined by (2), we propose to adapt the EM algorithm as follows. Let the E-step now consist in the calculation of the expectation of $\log L_c(\Psi)$ with respect to $g_c(\cdot|m; \Psi^{(q)})$ defined by (3):

$$Q(\Psi, \Psi^{(q)}) = \mathbb{E}_{\Psi^{(q)}} [\log L_c(\Psi) | m] = \frac{\int \log(L_c(\Psi)) g_c(\mathbf{x}; \Psi^{(q)}) pl(\mathbf{x}) d\mathbf{x}}{L(\Psi^{(q)})}. \quad (4)$$

The M-step is unchanged and requires the maximization of $Q(\Psi, \Psi^{(q)})$ with respect to Ψ . The E²M algorithm alternately repeats the E- and M-steps above until the increase of observed-data likelihood becomes smaller than some threshold. The following theorem shows that E²M algorithm inherits the monotonicity property of the EM algorithm, which ensures convergence provided the sequence of incomplete-data likelihood values remains bounded from above.

Theorem 1. *Any sequence $L(\Psi^{(q)})$ for $q = 0, 1, 2, \dots$ of likelihood values obtained using the E²M algorithm is non decreasing, i.e., it verifies*

$$L(\Psi^{(q+1)}) \geq L(\Psi^{(q)}), \quad \forall q. \quad (5)$$

Proof. The proof is similar to that of Dempster et al. [4]. □

To conclude this section, we may note that the p.d.f. $g_c(\mathbf{x}|m; \Psi)$ and, consequently, the E²M algorithm depend only on the contour function $pl(\mathbf{x})$ and they are unchanged if $pl(\mathbf{x})$ is multiplied by a constant. Consequently, the results are unchanged if m is converted into a probability distribution by normalizing the contour function.

5 Normal Mean and Variance Estimation

To illustrate the above algorithm, let us assume that the complete data $\mathbf{x} = (x_1, \dots, x_n) \in \mathcal{X} = \mathbb{R}^n$ is a realization from an i.i.d. random sample from a univariate normal distribution $\mathcal{N}(\mu, \sigma^2)$. The parameter vector is thus $\Psi = (\mu, \sigma)$. The observed data has the form $\mathbf{y} = (\mathbf{y}_1, \dots, \mathbf{y}_n)$ with $\mathbf{y}_i = (w_i, \alpha_i)$, where w_i is an estimate of x_i (provided, e.g., by a sensor), and $\alpha_i \in [0, 1]$ is a degree of confidence in that estimation. For each \mathbf{y}_i , there are two interpretations θ_{i1} and θ_{i2} . Under interpretation θ_{i1} , we admit that $x_i = w_i$; under interpretation θ_{i2} , we know only that $x_i \in \mathbb{R}$. The probability for interpretation θ_{i1} to be correct is α_i , which can thus be interpreted as a degree of reliability of the piece of information \mathbf{y}_i . The induced mass function m_i on \mathbb{R} is defined by

$$m_i(\{w_i\}) = \alpha_i, \quad m_i(\mathbb{R}) = 1 - \alpha_i.$$

The corresponding contour function is defined by

$$pl_i(x) = \alpha_i \delta(x - w_i) + 1 - \alpha_i$$

for all $x \in \mathbb{R}$, where $\delta(\cdot)$ is the Dirac Delta function.

Let $g_c(\cdot; \mu, \sigma)$ denote the normal p.d.f. with mean μ and standard deviation σ . The observed data log likelihood is

$$\log L(\mu, \sigma) = \sum_{i=1}^n \log \left(\int_{-\infty}^{\infty} g_c(x; \mu, \sigma) p l_i(x) dx \right) = \sum_{i=1}^n \log (\alpha_i g_c(w_i; \mu, \sigma) + 1 - \alpha_i),$$

which is to be maximized with respect to μ and σ .

The complete data log likelihood is

$$\begin{aligned} \log L_c(\mu, \sigma) &= -\frac{n}{2} \log(2\pi) - n \log \sigma - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 = \\ &= -\frac{n}{2} \log(2\pi) - n \log \sigma - \frac{1}{2\sigma^2} \left(\sum_{i=1}^n x_i^2 - 2\mu \sum_{i=1}^n x_i + n\mu^2 \right). \end{aligned}$$

Consequently,

$$\begin{aligned} Q(\Psi, \Psi^{(q)}) &= -\frac{n}{2} \log(2\pi) - n \log \sigma \\ &\quad - \frac{1}{2\sigma^2} \left(\sum_{i=1}^n \beta_i^{(q)} - 2\mu \sum_{i=1}^n \gamma_i^{(q)} + n\mu^2 \right), \quad (6) \end{aligned}$$

where $\gamma_i^{(q)}$ and $\beta_i^{(q)}$ denote, respectively, the expectations of X and X^2 with respect to the conditional probability distribution

$$g_c(\cdot | m_i; \Psi^{(q)}) = g_c(\cdot; \mu^{(q)}, \sigma^{(q)}) \oplus m_i$$

defined by

$$g_c(x | m_i; \Psi^{(q)}) = \frac{g_c(x; \Psi^{(q)}) p l_i(x)}{\int_{-\infty}^{+\infty} g_c(u; \Psi^{(q)}) p l_i(u) du} = \frac{g_c(x; \Psi^{(q)}) [\alpha_i \delta_{w_i}(x) + (1 - \alpha_i)]}{\alpha_i g_c(w_i; \Psi^{(q)}) + 1 - \alpha_i}.$$

The following equalities thus hold:

$$\gamma_i^{(q)} = \frac{\alpha_i g_c(w_i; \Psi^{(q)}) w_i + (1 - \alpha_i) \mu^{(q)}}{\alpha_i g_c(w_i; \Psi^{(q)}) + 1 - \alpha_i} \quad (7)$$

and

$$\beta_i^{(q)} = \frac{\alpha_i g_c(w_i; \Psi^{(q)}) w_i^2 + (1 - \alpha_i) \left[\left(\mu^{(q)} \right)^2 + \left(\sigma^{(q)} \right)^2 \right]}{\alpha_i g_c(w_i; \Psi^{(q)}) + 1 - \alpha_i}. \quad (8)$$

The maximum of $Q(\Psi, \Psi^{(q)})$ defined by (6) is obtained for the following values of μ and σ :

$$\mu^{(q+1)} = \frac{1}{n} \sum_{i=1}^n \gamma_i^{(q)} \quad (9)$$

and

$$\sigma^{(q+1)} = \sqrt{\frac{1}{n} \sum_{i=1}^n \beta_i^{(q)} - (\mu^{(q+1)})^2}. \tag{10}$$

In E-step of the E²M algorithm for this problem thus consists in the calculation of $\gamma_i^{(q)}$ and $\beta_i^{(q)}$ for all i using (7) and (8), respectively. The M-step then updates the estimates of μ and σ using (9) and (10). The algorithm stops when the relative increase of the observed data likelihood becomes less than some threshold ε .

Example 1. To illustrate the application of the above algorithm to a situation where data are unreliable, we considered the following experiments. Random samples of size $n = 100$ were drawn from a standard normal distribution. For each realization x_i , a number α_i was drawn from the uniform distribution $\mathcal{U}_{[0,1]}$. With probability α_i , w_i was defined as x_i , and with probability $1 - \alpha_i$ it was set to $x_i + \varepsilon_i$, with $\varepsilon_i \sim \mathcal{N}(0, s^2)$. Parameters μ and σ were estimated using the E²M algorithm based on the data (w_i, α_i) , $i = 1, \dots, n$. The experiment was repeated $N = 100$ times and mean squared errors on μ and σ were computed. The results are shown in Figure 1. Our approach was compared with the simple strategy that consists in estimating μ and σ by the sample mean and standard deviation of the w_i for all i such that $\alpha_i \geq c$, for different choices of c . We can see that the E²M algorithm is much more robust than this simple reference method. Further experiments involving comparisons with more sophisticated alternative estimators are under way.

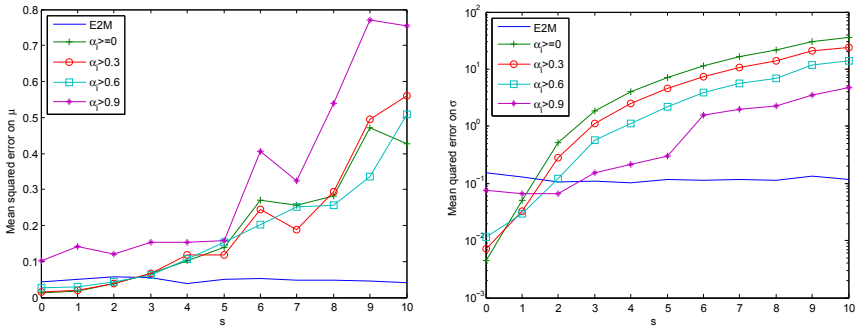


Fig. 1 Mean squared errors on μ (left) and σ (right, logarithmic y scale) as functions of the noise standard deviation s for the E²M algorithm and alternative methods (see details in text).

6 Conclusion

An iterative procedure for estimating the parameters in a statistical model using evidential data has been proposed. This procedure, which generalizes the EM algorithm, minimizes the degree of conflict between the uncertain

observations and the parametric model. It provides a general mechanism for statistical inference when the observed data are uncertain. It remains an open problem to determine the conditions under which the obtained estimator is consistent. This is the topic of on-going research.

References

1. Aggarwal, C.C., Yu, P.S.: A survey of uncertain data algorithms and applications. *IEEE Trans. Knowl. Data Eng.* 21(5), 609–623 (2009)
2. Côme, E., Oukhellou, L., Denœux, T., Aknin, P.: Learning from partially supervised data using mixture models and belief functions. *Pattern Recogn.* 42(3), 334–348 (2009)
3. Dempster, A.P.: Upper and lower probabilities induced by a multivalued mapping. *Ann. Math. Statist.* 38, 325–339 (1967)
4. Dempster, A.P., Laird, N.M., Rubin, D.B.: Maximum likelihood from incomplete data via the EM algorithm. *J. Roy. Statist. Soc. Ser. B* 39, 1–38 (1977)
5. Denœux, T., Masson, M.-H., Hébert, P.-A.: Nonparametric rank-based statistics and significance tests for fuzzy data. *Fuzzy Sets Syst.* 153, 1–28 (2005)
6. Gebhardt, J., Gil, M.A., Kruse, R.: Fuzzy set-theoretic methods in statistics. In: Slowinski, R. (ed.) *Fuzzy sets in decision analysis, operations research and statistics*, pp. 311–347. Kluwer Academic Publishers, Boston (1998)
7. Jraidi, I., Elouedi, Z.: Belief classification approach based on generalized credal EM. In: Mellouli, K. (ed.) *ECSQARU 2007. LNCS (LNAI)*, vol. 4724, pp. 524–535. Springer, Heidelberg (2007)
8. McLachlan, G.J., Krishnan, T.: *The EM Algorithm and Extensions*. Wiley, New York (1997)
9. Shafer, G.: *A mathematical theory of evidence*. Princeton University Press, Princeton (1976)
10. Vannoorenberghe, P., Smets, P.: Partially supervised learning by a credal EM approach. In: Godó, L. (ed.) *ECSQARU 2005. LNCS (LNAI)*, vol. 3571, pp. 956–967. Springer, Heidelberg (2005)
11. Zadeh, L.A.: Probability measures of fuzzy events. *J. Math. Anal. Appl.* 10, 421–427 (1968)