

Ensemble smoother with multiple data assimilation as a tool for curve fitting and parameter uncertainty characterization: Example applications to fit non-linear sorption isotherms

Vanessa A. Godoy · Gian F. Napa-García ·

J. Jaime Gómez-Hernández

Received: date / Accepted: date

Abstract The ensemble smoother with multiple data assimilation (ES-MDA) coupled to a normal-score transformation is used to fit a Langmuir isotherm curve to estimate its parameters (S_m and b) and their uncertainty. The highlights of this work are three: i) the ES-MDA can be used as a curve fitting procedure, ii) the ES-MDA provides also a full uncertainty quantification about the fitted parameters and iii) for the specific case of the Langmuir isotherm, parameter S_m is well identified with little uncertainty, while parameter b is well identified with a larger uncertainty, indicative that solute concentrations are more sensitive to S_m than to b . As a by-product, the number of samples required to characterize the joint uncertainty of Langmuir isotherm parameters is also investigated; it can be concluded that the minimum number of samples to use is six, with best results

Vanessa A. Godoy ✉
Research Institute of Water and Environmental Engineering, Universitat Politècnica de València, 46022, Valencia, Spain
E-mail: godoyalmeida@gmail.com

Gian F. Napa-García
Research Institute of Water and Environmental Engineering, Universitat Politècnica de València, 46022, Valencia, Spain

J. Jaime Gómez-Hernández
Research Institute of Water and Environmental Engineering, Universitat Politècnica de València, 46022, Valencia, Spain

obtained with eight samples, a value larger than the number recommend in the literature.

Keywords Tracer tests · Inverse modeling · Solute transport · Batch test · Bayesian methods

Acknowledgements The first author acknowledges the financial support from the Schlumberger Foundation by means of the program Faculty for the Future. The last author wishes to acknowledge the financial contribution of the Spanish Ministry of Science and Innovation through project number PID2019-109131RB-I00.

Declarations

Funding

This research has been supported by the Spanish Ministry of Science and Innovation through project number PID2019-109131RB-I00 and by the Schlumberger Foundation by means of the program Faculty for the Future.

Conflicts of interest/Competing interests

Not applicable

Availability of data and material

Not applicable

Code availability

Not applicable

1 Introduction

2 The retention of chemical constituents through the transfer of ions from the aque-
3 ous phase (sorptives) to the solid phase (sorbent) is widely known as sorption.
4 To quantify the sorption characteristics of a soil, batch experiments are usually
5 performed in which the relationship between the equilibrium concentrations of the
6 sorptive and the sorbate is measured for a variety of sorptive concentrations while
7 holding temperature constant; the outcome of this experiment results in what is
8 known as a sorption isotherm, an example of which is shown in Fig. 1. The most
9 common sorption isotherms are linear, Freundlich, and Langmuir isotherms.

10 The transport of dissolved contaminants in soils mostly depends on the sorption
11 capacity of soils and rocks (Pathak and Sharma 2018). The sorption processes are
12 important in a variety of applications in the field of geosciences, including remedi-
13 ation systems design (Luo et al. 2014), the selection of a waste site (Bouchelaghem
14 2018), the design of groundwater pump-and-treat systems or the evaluation of nat-
15 ural attenuation and salinity (Park et al. 2007; Tavakoli-Kivi et al. 2019). Most
16 of these applications involve at some point the use of numerical models, in which
17 the parameters defining the isotherm equation are required as input data (Guo
18 et al. 2019; Masood and Abd Ali 2020). An accurate determination of these pa-
19 rameters is crucial for good transport predictions, especially when these models
20 are used to perform quantitative risk analysis (Capilla et al. 1998; Franssen and
21 Gómez-Hernández 2002; Gómez-Hernández and Wen 1994; Hinz et al. 1994; Li
22 et al. 2011; Zhou et al. 2011; Fairey and Wahman 2013). Hinz et al. (1994) studied
23 the sensitivity of solute transport predictions on sorption isotherm parameters by
24 quantifying the ratio of the input concentration to the maximum sorption capac-
25 ity; the results show that the retardation of contaminants are highly influenced by
26 the Langmuir isotherm parameters.

27 Determining sorption isotherm parameters based on data from laboratory ex-
28 periments is always difficult and involves uncertainties. The chemical composition
29 of the solute, the physical structure of the sorbent, the effect of temperature,

30 measurement errors, and natural variations in soils and rocks are examples of
31 sources of uncertainty in laboratory experiments. This uncertainty, then, propa-
32 gates through the isotherm fitting process onto the isotherm parameters. It is also
33 still unclear how many experiments should be run to get a good fit of the isotherm
34 curve with small estimate uncertainties. Roy et al. (1991), in his thorough analysis
35 of batch-type adsorption procedures for estimating soil attenuation of chemicals,
36 recommends a minimum of five experiments to fit the isotherm parameters.

37 Although there exist a vast number of papers related to sorption isotherms,
38 few deal with the uncertainty associated with parameter estimation from batch ex-
39 periments. Fairey and Wahman (2013) compared weighted regression with Markov
40 chain Monte Carlo (McMC) to determine the joint uncertainty of Freundlich sorp-
41 tion isotherm fitting parameters, and found that both frequentist and Bayesian
42 analyses reduced the uncertainty in parameters compared with treating the uncer-
43 tainty independently, and that the difference between the two techniques becomes
44 more pronounced as the degree of non-linearity in the isotherm increases. In a re-
45 cent study, a hierarchical Bayesian model combined with McMC was satisfactorily
46 used to estimate parameters from multiple sources of experimental data of sorp-
47 tion and to characterize the uncertainty of Langmuir isotherm parameters (Shih
48 et al. 2020).

49 In this work, several numerical examples were performed to propose a proce-
50 dure to fit non-linear isotherm curves, and to analyze how the number of experi-
51 mental pairs used affects the estimated values and their uncertainty. The example
52 is applied to the fitting of the two parameters that define the Langmuir isotherm
53 and the two parameters that define the Freundlich isotherm, the expressions of
54 which will be introduced below. Additionally, we investigate the influence of the
55 ensemble size, the measurement-error magnitude, and the number of ES-MDA
56 iterations on the uncertainty estimation.

57 The good results obtained in this specific context allows us to make a recom-
58 mendation that the ES-MDA be used for curve-fitting in the general sense, not

59 limited to fitting isotherm curves. Traditionally this type of fitting is done using
60 least-square approaches, yielding reasonable results; however least-squares will not
61 provide an estimation of the uncertainty of the estimates, something that the ES-
62 MDA will do independently of the prior distributions adopted for the parameters
63 to be fitted. Also, the analysis of the stabilization of the statistics of the posterior
64 distributions of the fitted parameters can be used as a tool to determine how many
65 samples are needed to obtain reliable parameter estimates, as will be shown.

66 The proposed procedure is based on the ensemble smoother with multiple data
67 assimilation (ES-MDA) (Emerick and Reynolds 2013; Evensen 2018), a method
68 widely used in geosciences (Emerick 2017; Chen and Oliver 2012; Todaro et al.
69 2019; 2021; Silva et al. 2021a), and it includes a normal-score transformation to
70 deal with the possible non-Gaussianity of both prior and posterior uncertainties
71 about the parameters (Capilla et al. 1999; Zhou et al. 2011; Li et al. 2012; Xu
72 and Gómez-Hernández 2015; 2016; 2018). Recently, the ES-MDA was satisfacto-
73 rily used to fit the parameters of a B-Spline curve conditioned to well-test data
74 (Silva et al. 2021b); however, this application is a common approach to solve a
75 history matching exercise in which the parameters to be identified are not mate-
76 rial parameters, such as permeability or porosity, but the geometrical parameters
77 (widths and lengths) that define the turbidite lobes in the reservoir. The fact that
78 the authors define a B-Spline curve as a function of the parameters identified does
79 not imply that they are performing curve fitting in its traditional sense as used in
80 this paper. To the best of our knowledge, this is the first paper in which the normal-
81 score ES-MDA is used for curve fitting and estimation uncertainty quantification
82 of an isotherm curve, and that investigates how many samples are necessary for a
83 proper characterization.

84 The remainder of this paper is organized as follows. After a description of the
85 ES-MDA in Sect. 2, materials and methods are reported in Sect. 3, results are
86 presented and discussed in Sect. 4, conclusions are drawn in Sect. 5, and finally
87 additional synthetic examples are presented in the Appendix.

88 The work described in this paper was presented as a poster at the 46th Annual
89 Congress of the International Association of Hydrogeologists but it was never
90 published (Gómez-Hernández et al. 2019).

91 **2 The ensemble smoother with multiple data assimilation (ES-MDA)**

92 The ES-MDA algorithm is based on the ensemble smoother (ES) (Burgers et al.
93 1998), described by Emerick and Reynolds (2013) and Evensen (2018) as an al-
94 ternative to the ensemble Kalman filter (Xu et al. 2013; Zhou et al. 2012). The
95 ES-MDA is an iterative data assimilation method that updates parameters (in our
96 case, the isotherm parameters) making use of a set of system states (in our case,
97 equilibrium concentrations) and the deviations between the predictions resulting
98 from the current parameter values with respect to the experimental observations.
99 The relationship between parameters and observations must be known and a for-
100 ward model relating parameters and state variables must be available (in our case,
101 the forward model is simply the isotherm equation).

102 The assimilation procedure used by the ES-MDA includes an initialization
103 step, to generate N_e parameter realizations through statistical or geostatistical
104 methods, a forecast step, and an update step. In the forecast step, the forward
105 model is solved for each realization i , to obtain model predictions of the system
106 state. Then, the vector \mathbf{P} of model parameters used for the forecast is updated
107 based on the discrepancy between observations and their model predictions. The
108 updated parameter vector \mathbf{P}^u is given, for each realization, by

$$\mathbf{P}_i^u = \mathbf{P}_i + \mathbf{K} \left[\mathbf{Y}^{\text{ob}} + \varepsilon_i^{\text{ob}} - \mathbf{Y}_i \right], \quad \{i = 1, \dots, N_e\}, \quad (1)$$

109 where the subscript i refers to a specific realization, \mathbf{Y}_i is the vector of model
110 predictions for realization i , \mathbf{Y}^{ob} is the vector of state observations, $\varepsilon_i^{\text{ob}}$ is the
111 vector of observation errors for realization i (the observations errors have zero
112 mean and a covariance given by matrix \mathbf{R}) and \mathbf{K} is the Kalman gain, given by

$$\mathbf{K} = \mathbf{C}_{P,Y} (\mathbf{C}_{Y,Y} + \mathbf{R})^{-1}, \quad (2)$$

113 where $\mathbf{C}_{Y,Y}$ is the auto-covariance of the state variables and $\mathbf{C}_{P,Y}$ is the cross-
 114 covariance between all parameters and state variables, which are computed from
 115 the ensemble of realizations as

$$\mathbf{C}_{P,Y} = \frac{1}{N_e - 1} \sum_{i=1}^{N_e} (\mathbf{P}_i - \bar{\mathbf{P}}) (\mathbf{Y}_i - \bar{\mathbf{Y}})^T, \quad (3)$$

$$\mathbf{C}_{Y,Y} = \frac{1}{N_e - 1} \sum_{i=1}^{N_e} (\mathbf{Y}_i - \bar{\mathbf{Y}}) (\mathbf{Y}_i - \bar{\mathbf{Y}})^T, \quad (4)$$

116 with $\bar{\mathbf{P}}$ and $\bar{\mathbf{Y}}$ being the ensemble means of parameters and predictions, respec-
 117 tively.

118 In a problem where there are n_p parameters (in our case, n_p will be two,
 119 since there are only two parameters in the Langmuir isotherm equation) and n_o
 120 observations (in our case, n_o varies between four and sixteen), vectors \mathbf{P}_i^u and \mathbf{P}_i
 121 have sizes $n_p \times 1$, vectors \mathbf{Y}_i^{ob} , $\varepsilon_i^{\text{ob}}$, and \mathbf{Y} have sizes $n_o \times 1$, the Kalman gain \mathbf{K}
 122 and the covariance $\mathbf{C}_{P,Y}$ are matrices of size $n_p \times n_o$, and the matrices $\mathbf{C}_{Y,Y}$ and
 123 \mathbf{R} are of size $n_o \times n_o$. When the observation errors are modeled as uncorrelated,
 124 \mathbf{R} is a diagonal matrix. In the covariance matrix calculation, $\bar{\mathbf{P}}$ is a column vector
 125 of size $n_p \times 1$ with the average values of each parameter computed through the
 126 realizations, $\bar{\mathbf{P}} = \frac{1}{N_e} \sum_{i=1}^{N_e} \mathbf{P}_i$, and, similarly $\bar{\mathbf{Y}}$ is a column vector of size $n_o \times 1$
 127 with the average values of each parameter computed through the ensemble of
 128 realizations, $\bar{\mathbf{Y}} = \frac{1}{N_e} \sum_{i=1}^{N_e} \mathbf{Y}_i$.

129 2.1 Dealing with nonlinear state equations

130 The resulting updated parameters from Eq. (1) will be optimal estimates if, and
 131 only if, the state equation is linear. The ES-MDA was proposed to deal with non-
 132 linear systems by iteratively applying this process of forecasting and updating

133 using the last updated parameters to make the next forecast. This iteration implies
 134 that the same data will be assimilated multiple times; for this reason, there is
 135 a need to inflate the covariance matrix by a coefficient α_j , at each iteration j ,
 136 satisfying the following equation Evensen (2018)

$$\sum_{j=1}^{N_j} \frac{1}{\alpha_j} = 1, \quad (5)$$

137 where N_j is the total number of iterations.

According to the approach proposed by Evensen (2018), to compute α_j , first, it is necessary to select any nonzero value for α'_0 , then, the remaining α'_j are calculated as

$$\alpha'_j = \frac{\alpha'_{j-1}}{\alpha_{geo}}, \quad (6)$$

where α_{geo} is a constant that controls the extent of the changes of α'_j from one iteration to the next. Finally, the values from Eq. (6) are normalized to obtain the final coefficients as

$$\alpha_j = \alpha'_j \left(\sum_{j=1}^{N_j} \frac{1}{\alpha'_j} \right). \quad (7)$$

138 We refer to Evensen (2018) for more details on the computation of the α_j . A
 139 modification in the update step is also required to consider the α_j coefficients. The
 140 update equation for the ES-MDA results then

$$\mathbf{P}_{\text{MDA},i}^u = \mathbf{P}_i + \mathbf{C}_{P,Y} (\mathbf{C}_{Y,Y} + \alpha_j \mathbf{R})^{-1} \left[\mathbf{Y}^{\text{ob}} + \sqrt{\alpha_j} \varepsilon_i^{\text{ob}} - \mathbf{Y}_i \right]. \quad (8)$$

The updating step as presented in Eq. (8) has the limitation of being suboptimal for parameters displaying a non-Gaussian distribution. To take advantage of the fact that the ES-MDA formulation is optimal when dealing with Gaussian parameters, a normal-score transformation can be performed. The advantage of using such transformation is that it can be applied to any prior distribution. After the ES-MDA formulation is applied in Gaussian space, a back transformation recovers the physical meaning of the parameters into the original non-Gaussian

space. In this work, a normal-score transformation is used at each iteration following the work by Zhou *et al.* (2011) in their proposal of the normal-score ensemble Kalman filter (NS-EnKF). The method consists in assuming a non-Gaussian prior for the parameters \mathbf{P} , which are transformed into Gaussian parameters \mathbf{G} after applying a Gaussian anamorphosis $T(\cdot)$, also known as Nataf transformation or normal-score transform (Nataf 1962)

$$\mathbf{G}_{\text{MDA},i}^u = T(\mathbf{P}_{\text{MDA},i}^u). \quad (9)$$

The forecasting step is performed using \mathbf{P} as input to the state equation, but the updating is performed on the Gaussian parameter vector \mathbf{G} computed after the Gaussian transform of \mathbf{P}

$$\mathbf{G}_{\text{MDA},i}^u = \mathbf{G}_i + \mathbf{C}_{G,Y} (\mathbf{C}_{Y,Y} + \alpha_j \mathbf{R})^{-1} [\mathbf{Y}^{\text{ob}} + \sqrt{\alpha_j} \varepsilon_i^{\text{ob}} - \mathbf{Y}_i]; \quad (10)$$

notice that the cross-covariance between parameters and state is computed on the Gaussian transform of the parameters. Finally, the updated or posterior distribution is recovered by applying the inverse transformation $T^{-1}(\cdot)$ of the Eq. (9) as

$$\mathbf{P}_{\text{MDA},i}^u = T^{-1}(\mathbf{G}_{\text{MDA},i}^u). \quad (11)$$

¹⁴¹ The different steps of the algorithm are summarized in the Algorithm 1 insert.

Algorithm 1: Iterative data assimilation

Set: N_j = The number of iterations of ES-MDA

Set: \mathbf{Y}^{ob} = Observation data (here, concentrations at the solid phase)

Set: N_e = The number of parameter realizations

Set: α_0 = Initial inflation coefficient

Set: α_{geo} = Constant that controls the extent of the change of α_j between iterations

begin

Generate an ensemble of initial parameters \mathbf{P} (here, these parameters are drawn from their prior uniform distributions)

Calculate all α_j coefficients such that

$$\alpha_j^i = \frac{\alpha_{j-1}^i}{\alpha_{geo}} \text{ and } \alpha_j = \alpha_j^i \left(\sum_{j=1}^{N_j} \frac{1}{\alpha_j^i} \right)$$

for $j \leftarrow 1$ **to** N_j **do**

for $i \leftarrow 1$ **to** N_e **do**

 Perturb the observations: $\mathbf{Y}^{\text{ob}} + \sqrt{\alpha_j} \varepsilon_i^{\text{ob}}$

 Run forward model (here, evaluate the Langmuir sorption

 isotherm for the different liquid phase equilibrium solute

 concentrations) using \mathbf{P}_i as input parameters to obtain \mathbf{Y}_i

end for

Apply a Gaussian anamorphosis: $\mathbf{G}_{\text{MDA},i}^u = \mathbf{T}(\mathbf{P}_{\text{MDA},i}^u)$

Calculate: $\mathbf{C}_{G,Y} = \frac{1}{N_e-1} \sum_{i=1}^{N_e} (\mathbf{G}_i - \overline{\mathbf{G}}) (\mathbf{Y}_i - \overline{\mathbf{Y}})^T$

Calculate: $\mathbf{C}_{Y,Y} = \frac{1}{N_e-1} \sum_{i=1}^{N_e} (\mathbf{Y}_i - \overline{\mathbf{Y}}) (\mathbf{Y}_i - \overline{\mathbf{Y}})^T$

Update:

$$\mathbf{G}_{\text{MDA},i}^u = \mathbf{G}_i + \mathbf{C}_{G,Y} (\mathbf{C}_{Y,Y} + \alpha_j \mathbf{R})^{-1} [\mathbf{Y}^{\text{ob}} + \sqrt{\alpha_j} \varepsilon_i^{\text{ob}} - \mathbf{Y}_i]$$

Back transform: $\mathbf{P}_{\text{MDA},i}^u = \mathbf{T}^{-1}(\mathbf{G}_{\text{MDA},i}^u)$

end for

end

The Langmuir isotherm is one of the most common models used for sorption in relation with transport in porous media. It explains how a solute distributes between the solid and liquid phases once equilibrium is reached. This isotherm considers that the surface at which the solute can adsorb onto the solid phase is finite and, therefore, there is a maximum adsorbed concentration possible; its expression is

$$S = \frac{S_m b C_e}{1 + b C_e}, \quad (12)$$

144 where S is the solid-phase equilibrium concentration $[\text{M}_{\text{chemical}} \text{M}_{\text{sorbent}}^{-1}]$, C_e is
 145 the liquid-phase equilibrium concentration $[\text{M}_{\text{chemical}} \text{L}_{\text{water}}^{-3}]$, S_m represents the
 146 maximum concentration of soil-adsorbed solute $[\text{M} \text{M}^{-1}]$, and b is an adsorption
 147 constant related to binding energy $[\text{L}^3 \text{M}^{-1}]$. Figure 1 shows a typical Langmuir
 148 isotherm.

149 In our case, in the context of the ES-MDA algorithm, the forward model is the
 150 Langmuir sorption isotherm function, Eq. (12), S_m and b are the model parameters,
 151 whereas S is the system state. The solute concentrations at equilibrium, C_e for the
 152 liquid phase will be the forcing terms of the forward model and they are known.
 153 The parameter vector \mathbf{P}_i for a given realization i is

$$\mathbf{P}_i = \begin{bmatrix} S_{m,i} \\ b_i \end{bmatrix}. \quad (13)$$

154 The system state vector \mathbf{Y}_i is the set of predicted solid-phase concentrations at
 155 equilibrium corresponding to the set of liquid-phase concentrations at equilibrium
 156 for which the corresponding laboratory tests have been performed

$$\mathbf{Y}_i = \begin{bmatrix} S_{1,i} \\ S_{2,i} \\ \dots \\ S_{n_o,i} \end{bmatrix}. \quad (14)$$

157 From those laboratory tests, there will be n_o observed solid-phase equilibrium
 158 concentrations resulting from the experiments

$$\mathbf{Y}^{ob} = \begin{bmatrix} S_1^{ob} \\ S_2^{ob} \\ \dots \\ S_{n_o}^{ob} \end{bmatrix}. \quad (15)$$

159 Based on a range of previous numerical experiments, several scenarios were
 160 considered in order to analyze the impact of different parameters in the estimation
 161 process of the Langmuir coefficients. More precisely, the number of realizations of
 162 the ensemble took the values of 30, 100 and 300; the number of observation n_o
 163 took all integer values between four and sixteen, the observation error standard
 164 deviation took the values of 10^{-2} , $5 \cdot 10^{-3}$ and 10^{-4} mg·g⁻¹; and the number of
 165 iterations of the ES-MDA, N_j , took all integer values between 1 and 6. The steps
 166 followed for any given scenario are described next.

167 The first step is to generate a set of data pairs (C_e, S) consistent with a Lang-
 168 muir isotherm as they could have obtained in the laboratory. These would be the
 169 value pairs that have to be curve fitted by the isotherm function. An ideal soil for
 170 which sorption follows the Langmuir isotherm and with realistic parameter values
 171 $S_m = 0.1$ mg·g⁻¹ and $b = 100$ l·mg⁻¹ (Godoy et al. 2018) is considered. Sixteen
 172 experiments are mimicked with equally spaced values of C_e in the interval between
 173 0 and 0.1 mg·l⁻¹. The corresponding observations are computed by applying the
 174 Langmuir equation and perturbing the resulting value with an error ε^{ob} drawn

175 from a Gaussian distribution with zero mean and standard deviation $\sigma_{\varepsilon_{\text{ob}}}$. From
176 these sixteen data pairs, the necessary n_o observations will be chosen.

177 Second, an ensemble of initial values for the two parameters is generated. The
178 initial parameters are drawn from the uniform distributions $S_m \sim U[0, 3] \text{ mg}\cdot\text{g}^{-1}$
179 and $b \sim U[0, 300] \text{ l}\cdot\text{mg}^{-1}$.

180 Third, the ES-MDA with a normal-score transformation, as described previ-
181 ously, was applied for each scenario.

182 Fourth, the moments (mean, standard deviation, kurtosis, and skewness) of the
183 ES-MDA final estimates were analyzed to investigate how many isotherm samples
184 are necessary to reach stable statistics and, consequently, acceptable estimates.
185 Specifically, the optimal number of samples will be determined by visual analysis
186 of the moments of the final probability distributions of the different parameters;
187 this optimal number will be achieved when the moments stabilize.

188 The results presented next correspond only to synthetic experiments, no real
189 experimental data have been used. The principal reason of this choice is that it is
190 the only way in which a comparison between the estimates and the “true” values
191 can be made and to perform an effective evaluation of the methodology. Including
192 an additional example with experimental data will not serve to verify or increase
193 the reliability of the synthetic results, as long as the soil being analyzed does
194 display an adsorption behavior suitable to be modeled by a Langmuir isotherm.
195 The method proposed does not pretend to be a method to discriminate between
196 isotherm curves, and therefore, it is important to note that including an additional
197 case with laboratory data will not improve the validity of the approach. In any
198 case, for the sake of completeness, two additional synthetic cases have been run,
199 which are discussed in the Appendix: one with a synthetic soil with a different
200 Langmuir isotherm, and another one with a Freundlich isotherm. These two ad-
201 ditional examples prove that the method is general enough for curve fitting and
202 it could be used with a different isotherm and even in a different curve-fitting
203 context.

204 4 Results and discussion

205 Figures 2, 3 and 4 serve to illustrate how the ES-MDA works. An initial ensemble
206 of 100 realizations of parameters (Fig. 2a) is generated using the above mentioned
207 uniform distributions, the pairs (S_m, b) are distributed randomly within the do-
208 main $U[0, 3] \times U[0, 300]$. Each of the points in Fig. 2a corresponds to one of the
209 isotherm curves plotted in Fig. 3a. We can see a wide scatter of potential isotherm
210 curves all of them quite far from the “true” curve. In Fig. 3 the points correspond-
211 ing to the observations are also displayed. The discrepancies between the S values
212 for the different curves and the observed ones, the intrinsic variability of the pa-
213 rameters values as measured by its covariance, and the cross-correlation between
214 parameters and predicted values serve to compute the different elements in the
215 ES-MDA equations and to update each one of the points in Fig. 2a into a new
216 pair that gives a new isotherm curve closer to the real one. After the first update,
217 the new pairs of parameters are shown in Fig. 2b, the dashed lines correspond
218 to the values of the isotherm curve used to generate the observations. It is quite
219 evident how, after one iteration, the range of variability of the S_m parameter is
220 much narrower than the initial range and quite close to the true value, whereas
221 the range of the b parameters is still scattered over the entire initial range. The
222 reason for this fast convergence of the S_m parameter is due to the higher sensitiv-
223 ity that the shape of the isotherm curve has to the S_m parameter than to the b
224 parameter. We can see how the cloud of pairs keeps reducing its spread after each
225 assimilation iteration, and at iteration #4 the cloud of parameters has collapsed
226 onto the true values with little spread. The remaining spread is a measure of the
227 residual uncertainty in their estimation. The results obtained for iterations #5 and
228 #6 are not shown since they are almost identical to those of iteration #4. Fig. 4
229 shows the histograms of the initial ensembles and after four iterations. The initial
230 histograms correspond to the starting uniform distributions and the histograms
231 after four iterations show a spike at the true value for S_m and a histogram with
232 some spread and a little bias for the b parameter. The smaller sensitivity of the

233 isotherm curves to this parameter makes it impossible to identify it more precisely.
234 This small sensitivity translates in that the set of isotherm curves corresponding
235 to this range of b values result in almost superposing curves in Fig. 3d.

236 The previous figures, which have been obtained for a specific scenario, demon-
237 strate the power of the ES-MDA to identify the Langmuir parameters, together
238 with a measure of their uncertainty. This exercise has been repeated for several
239 other scenarios with the objective to determine how many experiments should be
240 run, that is, how many observation pairs are needed to find the parameters that
241 provide the best fitting curve. One way to analyze this aspect is to seek when the
242 estimate of the statistics of the final set of parameters stabilizes with the num-
243 ber of samples, if at all. This analysis will be performed on the results after four
244 iterations of the ES-MDA; similarly as for the scenario displayed in Figures 2, 3
245 and 4, four iterations were enough for the parameter estimates to stabilize in all
246 scenarios.

247 Figure 5 shows results for the final ensemble of updated values for parameters
248 S_m . The values of the mean, standard deviation, skewness and kurtosis are shown
249 for all combinations of number of ensemble realizations, observation error and
250 observation samples. It is evident that the best results are obtained when using
251 an ensemble of 300 realizations, but good estimates of the true value (as given
252 by the ensemble mean) with little uncertainty (as given by the ensemble standard
253 deviation) can be obtained for all scenarios as soon as six observations are used.
254 The estimates stabilizes at six observations when using 300 realizations for all
255 statistics. A stable estimate for the means and standard deviations when using
256 a smaller number of realizations requires between eight and ten samples. The
257 estimated skewness and kurtosis vary more erratically for the scenarios with 30
258 and 100 realizations, in part due to the smaller number of elements to compute
259 these statistics, and in part due to the narrowness of the final distribution which
260 make these values very sensitive to small departures from the mean. It is very
261 interesting to note that the magnitude of the observation error has little or no

262 effect in the estimates of the means and standard deviations in all scenarios as
263 long as at least six observations are used in the estimation.

264 Figure 6 shows results for the final ensemble of updated values for parameter
265 b . The values of the mean, standard deviation, skewness and kurtosis are shown
266 for all combinations of number of ensemble realizations, observation error and
267 observation samples. Contrary to Fig. 5, there is not a striking difference on the
268 curves as a function of the number of realizations. The main reason for this results
269 is the already-mentioned fact that the Langmuir curve is less sensitive to parameter
270 b than to parameter S_m ; for this reason, as soon as an estimated value is relatively
271 close to the real one, the estimated Langmuir isotherm is almost indistinguishable
272 with the true one, and there is no need to update the parameter anymore. This
273 behavior is particularly noticeable in the values of the standard deviations; they do
274 not get as close to zero as for S_m but remain with non-zero values throughout and
275 with larger values when the observation errors are larger. Skewness and kurtosis
276 estimates behave as for S_m . Due to the smaller sensitivity of the isotherm to the
277 b value, an estimate based only on four samples would be enough.

278 The ES-MDA performed remarkably well for the purpose of estimating the
279 fitting parameters of a Langmuir isotherm in a wide range of scenarios, with the
280 additional benefit of providing also an estimate of their uncertainty. This uncer-
281 tainty estimate is a confidence measurement about the estimated value and it is
282 also a measurement of the sensitivity of the fitting to the parameter. After all, at
283 the end of the exercise, there is an experimental histogram showing the full distri-
284 bution of potential values for the parameters, from which the mean or the median
285 could be selected as best estimates, but from which an analysis of the parameter
286 values which are consistent with the observations can also be performed. Such an
287 analysis of the uncertainty about the estimates cannot be performed with stan-
288 dard fitting procedures, such as least-squares, that, at most, provide an estimated
289 value and an estimation error.

290 In summary and with regard to the recommendation by Roy et al. (1991)
291 that a minimum of five observations should be used to estimate the isotherm
292 parameters, we conclude that such a number would be enough for the estimation
293 of the b coefficient, but not for the estimation of S_m . Our recommendation would
294 be to increase that minimum number to six and, preferably, to eight.

295 The successful results in the application of the ES-MDA for curve fitting for the
296 three non-linear isotherm curves analyzed in the paper, makes us postulate the use
297 of the ES-MDA for general curve fitting when characterization of the uncertainty
298 about the final estimates is important.

299 **5 Conclusions**

300 In this paper, we proposed a procedure to fit sorption isotherm curves using
301 an ensemble smoother with multiple data assimilation (ES-MDA) coupled to a
302 normal-score transform. The main advantage of this approach is not only that the
303 parameter is good, but also that a characterization of the parameter estimation
304 uncertainty is obtained. In order to evaluate the proposed procedure, we performed
305 numerical examples with a variety of scenarios to additionally investigate the influ-
306 ence of the number of experimental pairs, the ensemble size, the measurement-error
307 magnitude, and the number of ES-MDA iterations on the uncertainty estimation.
308 Our results show that, since the shape of the Langmuir isotherm is much more
309 sensitive to the S_m parameter than to the b parameter, the precision in their identi-
310 fication is not the same. After four ES-MDA iterations, the cloud of experimental
311 pairs has collapsed onto the reference value for the S_m parameter and presents
312 some spread and a little bias for the b parameter. By investigating other scenarios
313 with several combinations of number of ensemble realizations, observation error
314 and observation samples we find that, for the S_m parameter, the best results are
315 obtained when using an ensemble of 300 realizations, and the use of at least six
316 observations can be enough to produce relatively good estimates of the true value
317 regardless of the scenario. An interesting finding is that when at least six obser-

318 vations are used, the magnitude of the observation error has almost no effect in
319 the estimates of the mean and standard deviation of S_m . For the b parameter,
320 there is not a clear difference on the curves as a function of the scenarios because
321 the Langmuir curve is little sensitive to this parameter. These results demonstrate
322 the power of the ES-MDA to identify the Langmuir parameters together with an
323 estimate of their uncertainty. We conclude that the actual recommendation that a
324 minimum of five observations should be used to estimate the Langmuir parameters
325 would be enough for the estimation of the b coefficient, but not for the estimation
326 of S_m . In order to correctly estimate Langmuir parameters together with their
327 uncertainty a minimum of six and, preferably, eight samples should be used. The
328 results of the application in the two additional cases included in the Appendix are
329 similar and reinforce our belief that the ES-MDA could be applied for standard
330 curve fitting when an uncertainty characterization about parameter estimates is
331 needed.

332 References

- 333 Bouchelaghem, F. (2018). Multi-scale study of pollutant transport and uptake in compacted
334 bentonite. *Mathematical Geosciences*, 50:495–523.
- 335 Burgers, G., Van Leeuwen, P. J., and Evensen, G. (1998). Analysis scheme in the ensemble
336 Kalman filter. *Monthly Weather Review*, 126:1719–1724.
- 337 Capilla, J. E., Gómez-Hernández, J. J., and Sahuquillo, A. (1998). Stochastic simulation
338 of transmissivity fields conditional to both transmissivity and piezometric head data—3.
339 application to the culebra formation at the waste isolation pilot plan (wipp), new mexico,
340 usa. *Journal of Hydrology*, 207:254–269.
- 341 Capilla, J. E., Rodrigo, J., and Gómez-Hernández, J. J. (1999). Simulation of non-gaussian
342 transmissivity fields honoring piezometric data and integrating soft and secondary infor-
343 mation. *Math. Geology*, 31:907–927.
- 344 Chen, Y. and Oliver, D. S. (2012). Ensemble randomized maximum likelihood method as an
345 iterative ensemble smoother. *Mathematical Geosciences*, 44:1–26.
- 346 Emerick, A. A. (2017). Investigation on principal component analysis parameterizations for
347 history matching channelized facies models with ensemble-based data assimilation. *Math-*
348 *ematical Geosciences*, 49:85–120.

349 Emerick, A. A. and Reynolds, A. C. (2013). Ensemble smoother with multiple data assimila-
350 tion. *Computers & Geosciences*, 55:3–15.

351 Evensen, G. (2018). Analysis of iterative ensemble smoothers for solving inverse problems.
352 *Computational Geosciences*, 22:885–908.

353 Fairey, J. L. and Wahman, D. G. (2013). Bayesian and Frequentist Methods for Estimatin-
354 ing Joint Uncertainty of Freundlich Adsorption Isotherm Fitting Parameters. *Journal of*
355 *Environmental Engineering*, 139:307–311.

356 Franssen, H. H. and Gómez-Hernández, J. (2002). 3d inverse modelling of groundwater flow at
357 a fractured site using a stochastic continuum model with multiple statistical populations.
358 *Stochastic Environmental Research and Risk Assessment*, 16:155–174.

359 Godoy, V. A., Zuquette, L. V., and Gómez-Hernández, J. J. (2018). Scale effect on hydraulic
360 conductivity and solute transport: Small and large-scale laboratory experiments and field
361 experiments. *Engineering Geology*, 243:196–205.

362 Gómez-Hernández, J. J., Napa-García, G. F., and Godoy, V. A. (2019). How to account for
363 uncertainty in the estimation of adsorption isotherm parameters. In Gómez-Hernández,
364 J. J. and Navarro, B. A., editors, *Groundwater Management and Governance: Coping*
365 *with Uncertainty, proceedings of IAHR2019, the 46th Annual Congress of the International*
366 *Association of Hydrogeologists*, page 342. Spanish Chapter of the International Association
367 of Mathematical Geosciences, Spanish Chapter of the International Association of
368 Mathematical Geosciences.

369 Gómez-Hernández, J. J. and Wen, X.-H. (1994). Probabilistic assessment of travel times in
370 groundwater modeling. *J. of Stochastic Hydrology and Hydraulics*, 8(1):19–56.

371 Guo, Z., Fogg, G. E., Brusseau, M. L., LaBolle, E. M., and Lopez, J. (2019). Modeling
372 groundwater contaminant transport in the presence of large heterogeneity: a case study
373 comparing MT3D and RWHEM. *Hydrogeology Journal*, 27:1363–1371.

374 Hinz, C., Gaston, L., and Selim, H. (1994). Effect of sorption isotherm type on predictions of
375 solute mobility in soil. *Water Resources Research*, 30:3013–3021.

376 Li, L., Zhou, H., and Gómez-Hernández, J. J. (2011). A comparative study of three-dimensional
377 hydraulic conductivity upscaling at the macro-dispersion experiment (made) site, columbus
378 air force base, mississippi (usa). *Journal of Hydrology*, 404:278–293.

379 Li, L., Zhou, H., Hendricks Franssen, H.-J., and Gómez-Hernández, J. J. (2012). Modeling
380 transient groundwater flow by coupling ensemble kalman filtering and upscaling. *Water*
381 *Resources Research*, 48:W01537.

382 Luo, Q., Wu, J., Yang, Y., Qian, J., and Wu, J. (2014). Optimal design of groundwater
383 remediation system using a probabilistic multi-objective fast harmony search algorithm
384 under uncertainty. *Journal of Hydrology*, 519:3305–3315.

385 Masood, Z. B. and Abd Ali, Z. T. (2020). Numerical modeling of two-dimensional simula-
386 tion of groundwater protection from lead using different sorbents in permeable barriers.
387 *Environmental Engineering Research*, 25:605–613.

388 Nataf, A. (1962). Determination des distribution don t les marges sont donnees. *Comptes*
389 *Rendus de l Academie des Sciences*, 225:42–43.

390 Park, D. K., Ko, N. Y., and Lee, K. K. (2007). Optimal groundwater remediation design
391 considering effects of natural attenuation processes: Pumping strategy with enhanced-
392 natural-attenuation. *Geosciences Journal*, 11:377–385.

393 Pathak, P. and Sharma, S. (2018). Sorption isotherms, kinetics, and thermodynamics of con-
394 taminants in indian soils. *Journal of Environmental Engineering*, 144:04018109.

395 Roy, W., Krapac, I., Chou, S., and Griffin, R. (1991). Batch-Type Adsorption Procedures for
396 Estimating Soil Attenuation of Chemicals.

397 Shih, C., Park, J., Sholl, D. S., Realf, M. J., Yajima, T., and Kawajiri, Y. (2020). Hierar-
398 chical Bayesian estimation for adsorption isotherm parameter determination. *Chemical*
399 *Engineering Science*, 214:115435.

400 Silva, T. M., Pesco, S., Barreto Jr, A., and Onur, M. (2021a). A new procedure for generating
401 data covariance inflation factors for ensemble smoother with multiple data assimilation.
402 *Computers & Geosciences*, page 104722.

403 Silva, T. M., Villalobos, R. S., Cardona, Y. A., Barreto, A., and Pesco, S. (2021b). Well-
404 testing based turbidite lobes modeling using the ensemble smoother with multiple data
405 assimilation. *Computational Geosciences*, 25:1139–1157.

406 Tavakoli-Kivi, S., Bailey, R. T., and Gates, T. K. (2019). A salinity reactive transport and
407 equilibrium chemistry model for regional-scale agricultural groundwater systems. *Journal*
408 *of Hydrology*, 572:274–293.

409 Todaro, V., D’Oria, M., Tanda, M. G., and Gómez-Hernández, J. J. (2019). Ensemble smoother
410 with multiple data assimilation for reverse flow routing. *Computers & Geosciences*, 131:32–
411 40.

412 Todaro, V., D’Oria, M., Tanda, M. G., and Gómez-Hernández, J. J. (2021). Ensemble smoother
413 with multiple data assimilation to simultaneously estimate the source location and the
414 release history of a contaminant spill in an aquifer. *Journal of Hydrology*, page 126215.

415 Xu, T. and Gómez-Hernández, J. J. (2015). Inverse sequential simulation: A new approach
416 for the characterization of hydraulic conductivities demonstrated on a non-Gaussian field.
417 *Water Resources Research*, 51:2227–2242.

418 Xu, T. and Gómez-Hernández, J. J. (2016). Characterization of non-Gaussian conductivities
419 and porosities with hydraulic heads, solute concentrations, and water temperatures. *Water*
420 *Resources Research*, 52:6111–6136.

- 421 Xu, T. and Gómez-Hernández, J. J. (2018). Simultaneous identification of a contaminant
422 source and hydraulic conductivity via the restart normal-score ensemble Kalman filter.
423 *Advances in Water Resources*, 112:106–123.
- 424 Xu, T., Jaime Gómez-Hernández, J., Zhou, H., and Li, L. (2013). The power of transient
425 piezometric head data in inverse modeling: an application of the localized normal-score enkf
426 with covariance inflation in a heterogenous bimodal hydraulic conductivity field. *Advances
427 in Water Resources*, 54:100–118.
- 428 Zhou, H., Gómez-Hernández, J. J., Hendricks Franssen, H.-J., and Li, L. (2011). An approach to
429 handling non-Gaussianity of parameters and state variables in ensemble Kalman filtering.
430 *Advances in Water Resources*, 34(7):844–864.
- 431 Zhou, H., Li, L., Franssen, H.-J. H., and Gómez-Hernández, J. J. (2012). Pattern recogni-
432 tion in a bimodal aquifer using the normal-score ensemble kalman filter. *Mathematical
433 Geosciences*, 44:169–185.

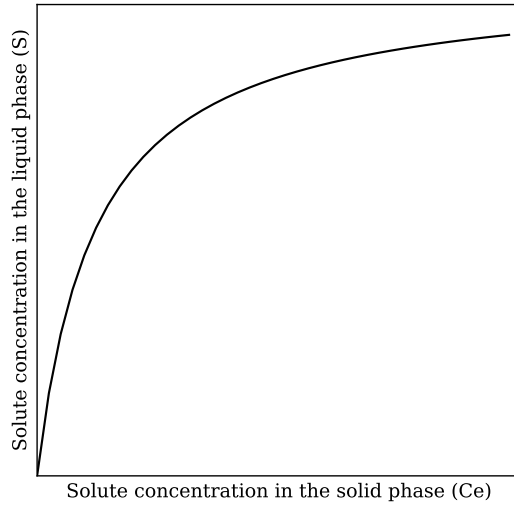


Fig. 1: Typical Langmuir isotherm.

434 A Additional Synthetic Examples

In the following, two additional examples are included to support the claims in the main body of the text; the first one is for a synthetic soil with adsorption characteristics given by a Langmuir isotherm with parameters $S_m = 7.2 \text{ mg}\cdot\text{g}^{-1}$ and $b = 0.174 \text{ l}\cdot\text{mg}^{-1}$ (?), and the second one for a soil characterized by a Freundlich isotherm

$$S = K_f C_e^{1/n} \quad (16)$$

435 with parameters $K_f = 1.5$ and $n^{-1} = 0.39$; S is given in $\text{mg}\cdot\text{g}^{-1}$ and C_e in $\text{mg}\cdot\text{l}^{-1}$ (?). Figure
 436 A.1 shows the three isotherms considered in the paper.

437 The procedure to fit the curves is the same as the one used in the body of the text.
 438 The initial sets of realizations are drawn from the following bivariate uncorrelated uniform
 439 distributions: $(S_m, b) \in U[0, 230] \times U[0, 0.8]$, and $(K_f, n^{-1}) \in U[0, 30] \times U[0.001, 0.99]$.

440 The evolution with the number of samples of the best estimate as given by the mean of the
 441 ensemble of updated parameters for the two cases can be seen in Figure1 A.2. The conclusions
 442 that can be drawn from the analysis of these figures are the same as from the analysis of

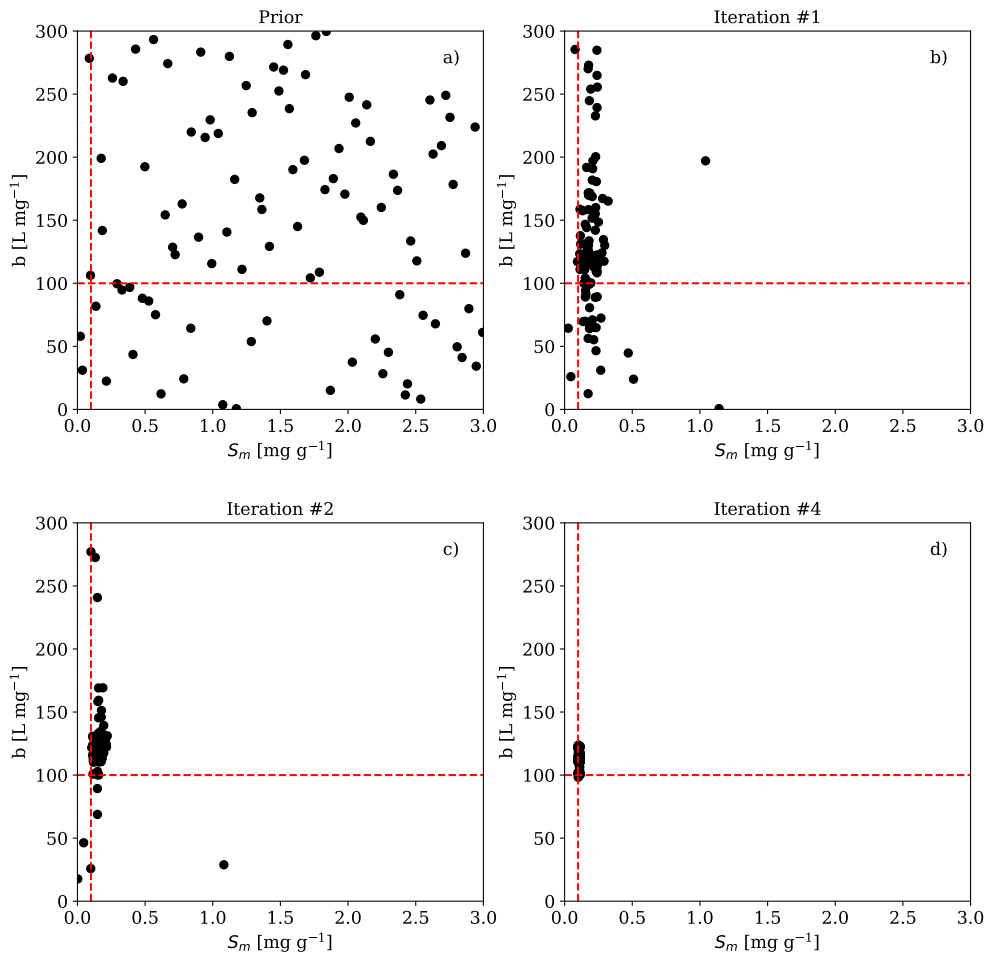


Fig. 2: Operation of the ES-MDA for: a) prior pairs of S_m and b ; parameters values obtained after ES-MDA #1, #2 and #4 iterations and 100 realizations.

443 the example in the main body. The estimated values are affected by the magnitude of the
 444 measurement errors, the larger the measurement errors, the larger the bias of the estimated
 445 value (as given by the mean of the ensemble results). When the error standard deviation
 446 is set to 1%, the estimates are quite close to the true value of the synthetic soil. For the
 447 Langmuir isotherm, the fluctuations of the mean m and mean b stabilize about six samples
 448 with a better stabilization the larger the number of realizations of the ensemble have been
 449 used. For the Freundlich isotherm, the estimation needs at least 13 samples and either 100 or
 450 300 realizations to retrieve good estimates when the error standard deviation is above 1%; for
 451 the smaller error, seven samples are necessary before the mean estimate stabilizes close to the

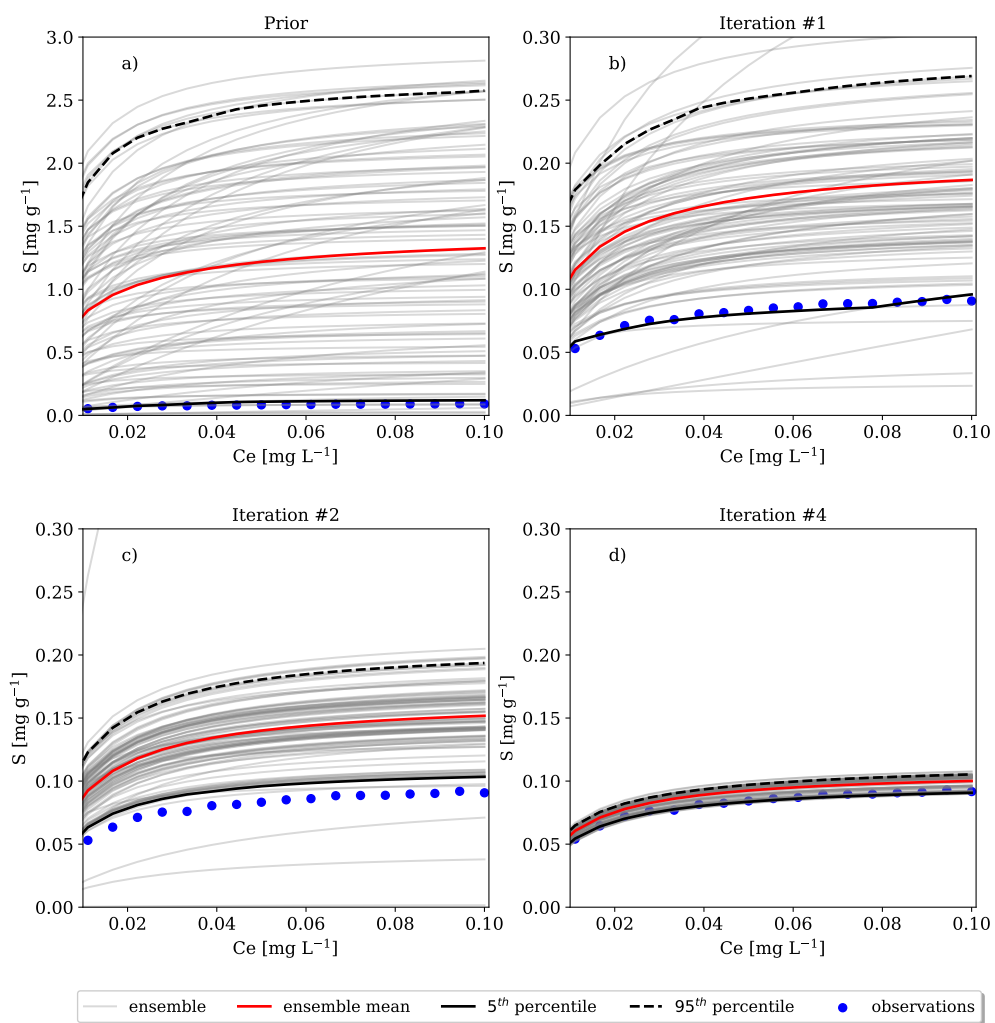


Fig. 3: Variability in the computation of the Langmuir isotherm using: **a** prior; and **b-d** updated parameters for 17 observations at iterations #1, #2 and #4, 100 realizations, and ε^{ob} of $10^{-3} \text{ mg g}^{-1}$.

452 reference value. It could be concluded that the Freundlich estimate may need more samples to
 453 ensure a good estimation of its parameters.

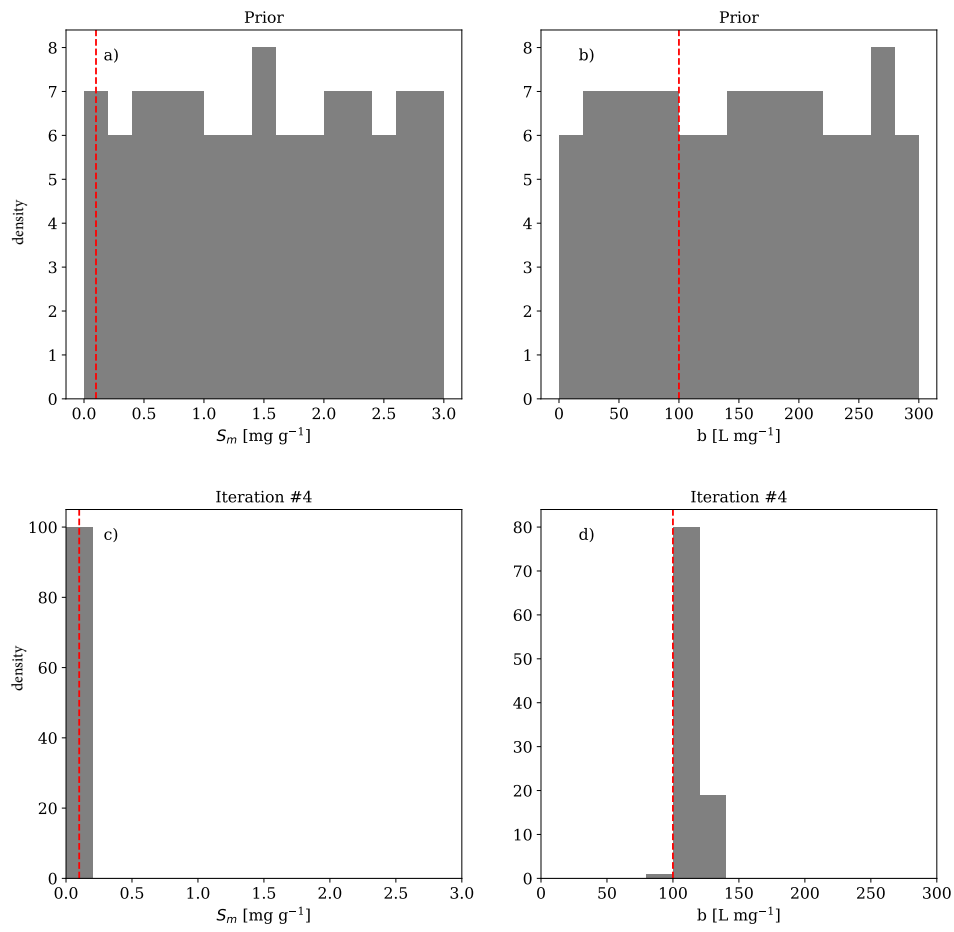


Fig. 4: The prior (a-b) and updated (c-d) histograms of the Langmuir isotherm parameters for 17 observations, 100 realizations, and ε^{ob} of $10^{-3} \text{ mg g}^{-1}$.

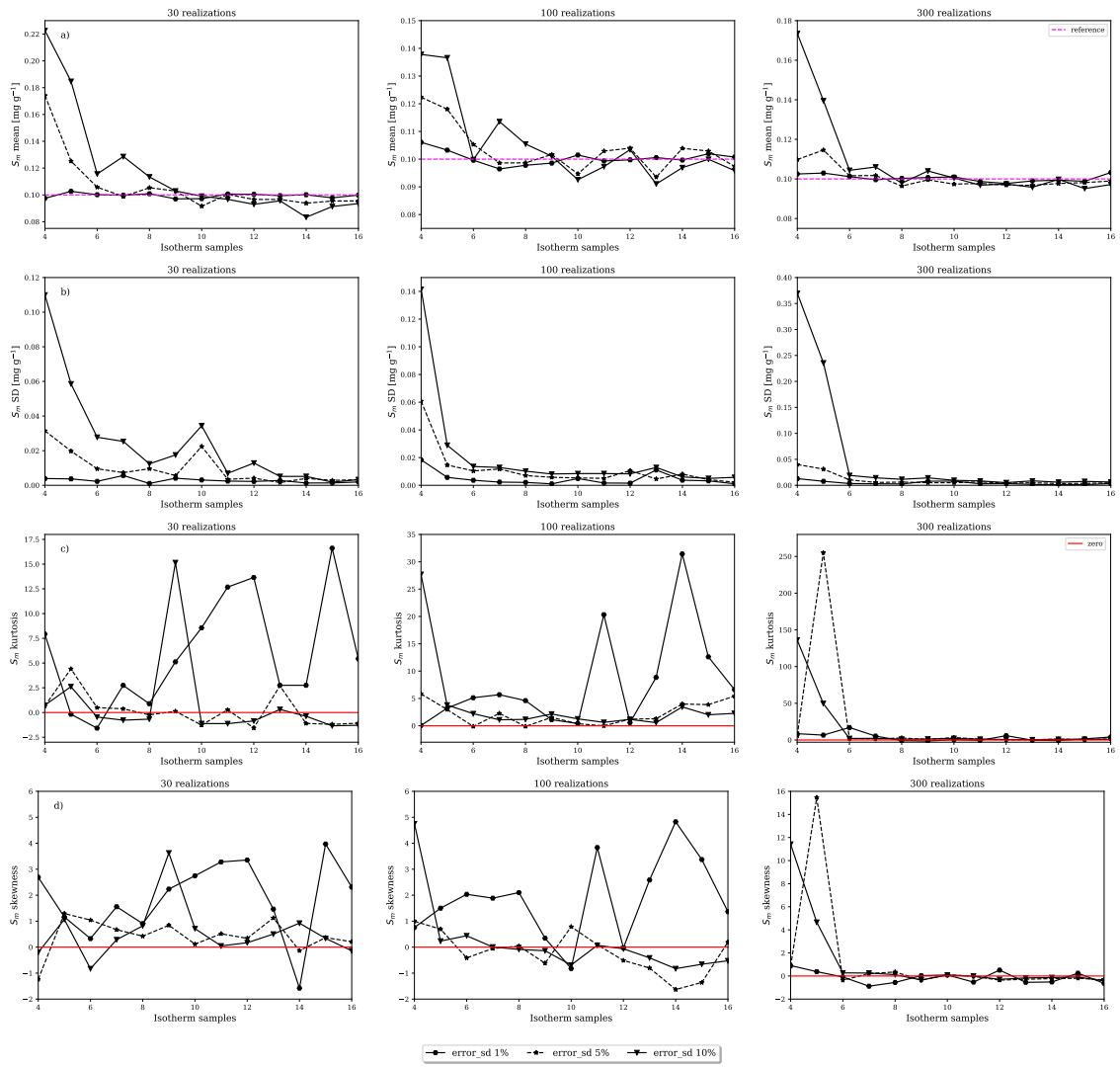


Fig. 5: From top to bottom, variation of the: a) mean, b) standard deviation, c) kurtosis, and d) skewness of S_m with the number of isotherm samples.

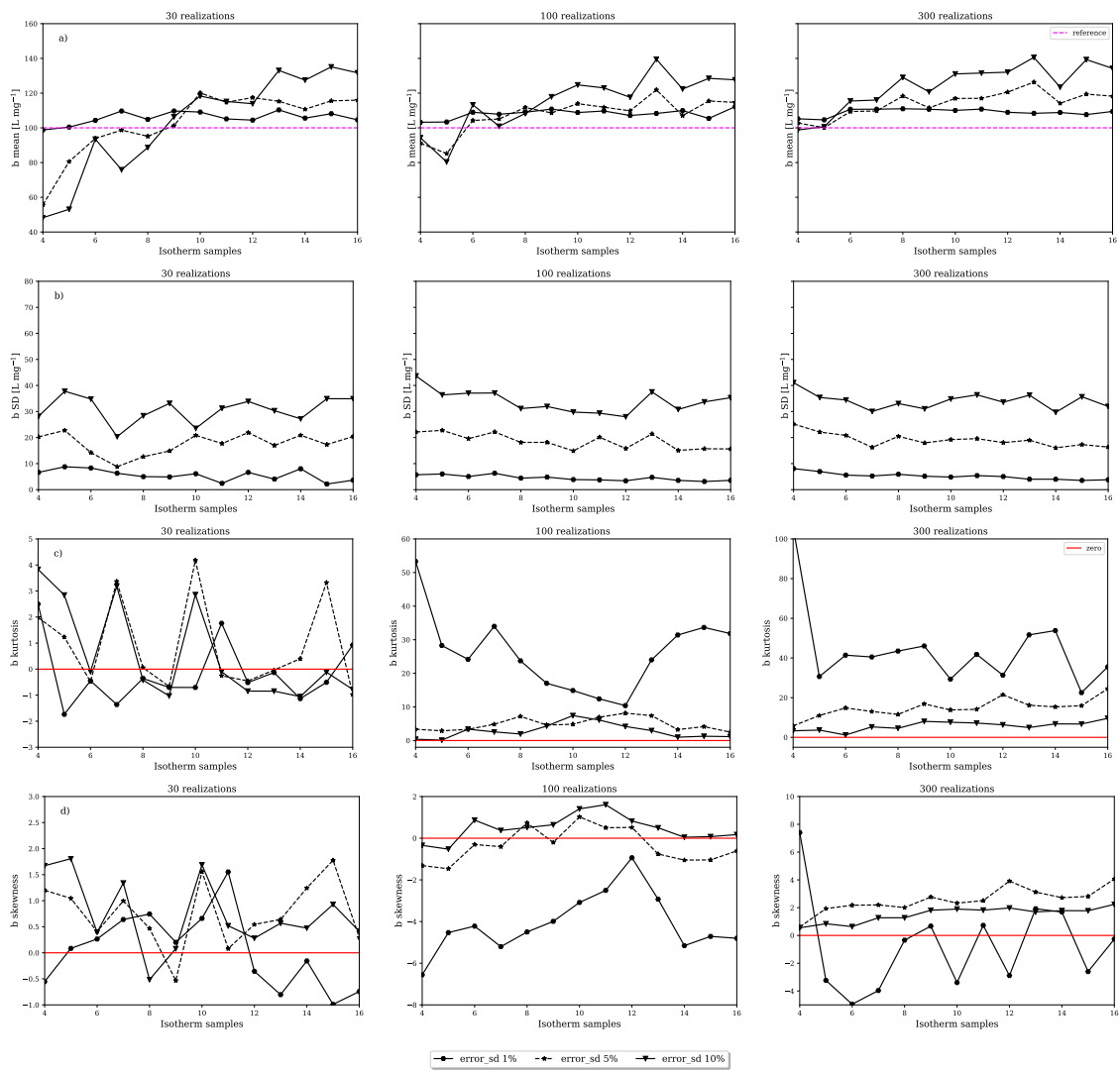


Fig. 6: From top to bottom, variation of the: a) mean, b) standard deviation, c) kurtosis, and d) skewness of b with the number of isotherm samples.

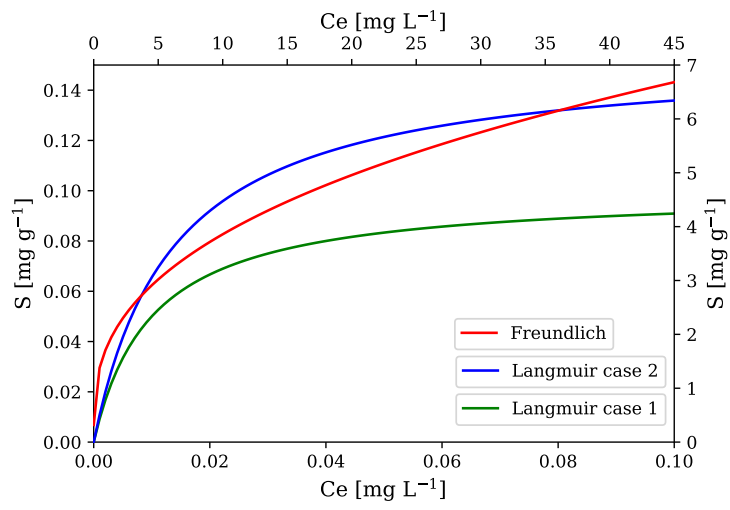


Fig. 7: The three isotherms analyzed in the paper. The Freundlich isotherm uses the left vertical axis, while the Langmuir isotherms use the right vertical axis.

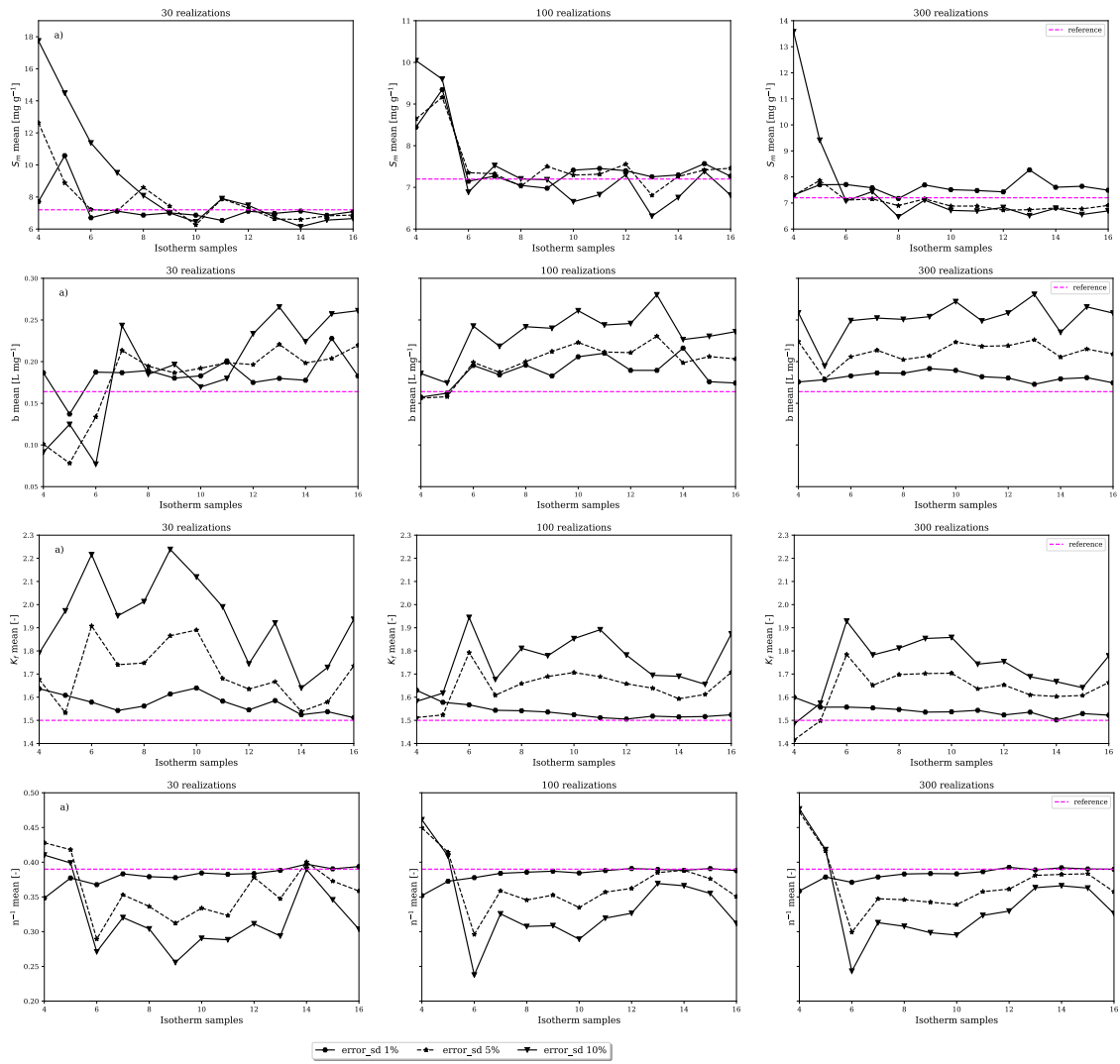


Fig. 8: From top to bottom, variation of the mean of parameter: S_m , b , K_f , and n^{-1} with the number of isotherm samples.