

Quantitative selection of sample structures in small-angle scattering using Bayesian methods

Yui Hayashi,^a Shun Katakami,^a Shigeo Kuwamoto,^b Kenji Nagata,^c Masaichiro Mizumaki^d and Masato Okada^{a*}

^aGraduate School of Frontier Sciences, University of Tokyo, Kashiwa, Chiba 277-8561, Japan, ^bJapan Synchrotron Radiation Research Institute, Sayo, Hyogo 679-5198, Japan, ^cNational Institute for Materials Science, Tsukuba, Ibaraki 305-0047, Japan, and ^dFaculty of Advanced Science and Technology, Kumamoto University, Kumamoto 860-8555, Japan. *Correspondence e-mail: okada@edu.k.u-tokyo.ac.jp

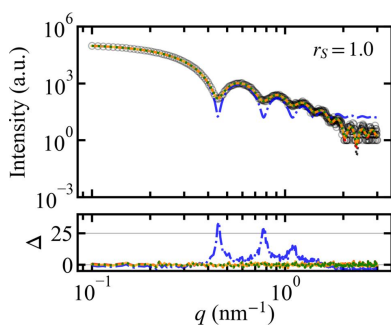
Small-angle scattering (SAS) is a key experimental technique for analyzing nanoscale structures in various materials. In SAS data analysis, selecting an appropriate mathematical model for the scattering intensity is critical, as it generates a hypothesis of the structure of the experimental sample. Traditional model selection methods either rely on qualitative approaches or are prone to overfitting. This paper introduces an analytical method that applies Bayesian model selection to SAS measurement data, enabling a quantitative evaluation of the validity of mathematical models. The performance of the method is assessed through numerical experiments using artificial data for multicomponent spherical materials, demonstrating that this proposed analysis approach yields highly accurate and interpretable results. The ability of the method to analyze a range of mixing ratios and particle size ratios for mixed components is also discussed, along with its precision in model evaluation by the degree of fitting. The proposed method effectively facilitates quantitative analysis of nanoscale sample structures in SAS, which has traditionally been challenging, and is expected to contribute significantly to advancements in a wide range of fields.

1. Introduction

In recent years, the analysis of nanoscale structures of materials has become increasingly important in advancing the development of new materials and understanding biological phenomena. Small-angle scattering (SAS) is a fundamental experimental method for analyzing such nanoscale structures. It involves irradiating substances with X-rays or neutron beams and analyzing the resulting scattering intensity data at small angles, typically 5° or less (Guinier & Fournet, 1955). SAS is versatile and applicable to a wide array of heterogeneous materials including nanoparticles, polymers, soft materials and fibers, and is utilized across many fields including materials science, chemistry and biology.

SAS measurement data are expressed in terms of scattering intensity that corresponds to a scattering vector, a physical quantity representing the scattering angle. Data analysis requires selection and parameter estimation of a mathematical model of the scattering intensity that contains information about the structure of the specimen. This selection process is critical as it involves assumptions about the structure of the specimen.

Traditionally, model selection in SAS data analysis has been performed by listing candidate models according to theoretical or empirical rules, conducting parameter fitting against the measurements, and comparing suitability using criteria such as χ -squared error, among other criteria (Breßler *et al.*,



2015; Da Vela & Svergun, 2020; Kline, 2006; Larsen *et al.*, 2018a; Pedersen, 1997; Schneidman-Duhovny *et al.*, 2010; Svergun *et al.*, 1995). Alternatively, models may be chosen on the basis of the general shape of the measurement data. However, these methods each have drawbacks: the former risks overfitting, which can lead to an overestimation of the model's degrees of freedom (Rambo & Tainer, 2013), while the latter yields only qualitative model selections. Furthermore, quantitatively evaluating the reliability of the results is challenging with traditional methods.

In this study, we propose a novel framework for SAS model selection that quantitatively assesses the validity of mathematical models that represent specimen structures in measurements. This approach uses Bayesian model selection within the framework of Bayesian inference, a method increasingly applied to analysis of various types of physical experimental data (Nagata *et al.*, 2012, 2019; Rappel *et al.*, 2020; Machida *et al.*, 2021; Moriguchi *et al.*, 2022; Nagai *et al.*, 2021; Kashiwamura *et al.*, 2022; Katakami *et al.*, 2022; Nelson & Prescott, 2019; Orioli *et al.*, 2020; Scheres, 2012; Ueda *et al.*, 2023). In the context of SAS data analysis, Bayesian inference has been applied to various cases, including ensembles of protein structures (Antonov *et al.*, 2016), regularization methods in parameter fitting (Larsen *et al.*, 2018b), indirect Fourier transforms (Hansen, 2000; Larsen & Pedersen, 2021), the estimation of particle size distributions (Asahara *et al.*, 2021) and specimen parameter estimates (Hayashi *et al.*, 2023). The method solves inverse problems by establishing the likelihood, which is the data generation model, and the prior distribution, which corresponds to the prior knowledge about the target being estimated. The posterior distribution is then calculated according to the model and parameters with the acquired data using Bayes' theorem. In our proposed method, the posterior probability of the data generation model is calculated under the measured data using the exchange Monte Carlo method (Hukushima & Nemoto, 1996), also known as parallel tempering, and then the resulting values are compared among the candidate models while concurrently obtaining Bayesian estimates of the model parameters. Moreover, since the validity of the measured data model is obtained as a posterior probability, the reliability of the results can be quantified by comparing these probabilistic values.

We conducted numerical experiments to assess the effectiveness of our proposed method. These experiments are based on synthetic data used to estimate the number of distinct components in a specimen, which was modeled as a mixture of monodisperse spheres of varying radii, scattering length densities and volume fractions. The results demonstrate the high accuracy, interpretability and stability of our method, even in the presence of measurement noise. To discuss the utility of the proposed method, we compare our approach with traditional model selection methods based on the reduced χ -squared error.

The structure of this paper is as follows. We first formalize the proposed analytical method, and then describe the model of multicomponent monodisperse spheres used in our numerical experiments. In Section 4, we detail the setup and

the results of these experiments using the proposed method to estimate the number of mixed components in the synthetic data. We then discuss the analytical capabilities of our method and the performance of the traditional method based on the degree of fitting. We conclude with implications and potential applications of our method.

2. Formulation of the proposed framework

In this section, we present a detailed formulation of our algorithm for selecting mathematical models for SAS specimens using Bayesian model selection. The pseudocode for this algorithm is provided in Algorithm 1.

Algorithm 1 : Proposed framework for quantitative selection of specimen model.

Require: The measured data, $\mathcal{D} = \{q_i, y_i\}_{i=1}^N$. The number of replicas, L . The inverse temperature, $\{\beta_l\}_{l=1}^L$ where $0 = \beta_1 < \dots < \beta_L = 1$. The burn-in, \mathcal{T}_0 . The number of samples, \mathcal{T}_1 . The step size for the Metropolis algorithm, $\{e_l\}_{l=1}^L$. Candidate models, $\{K_m\}_{m=1}^M$. The prior distribution of models, $\varphi(K)$. The prior distribution of parameters, $\varphi(\Xi|K)$.

Ensure: $\forall l \in \{1, \dots, L\}, \forall \tau \in \{\mathcal{T}_0 + 1, \dots, \mathcal{T}_1\}, \Xi_l^\tau \sim p(\Xi_l|\mathcal{D}, \beta_l), Z(K), F(K), p(K|\mathcal{D})$.

```

1: for  $m \in \{1, \dots, M\}$  do
2:   Initialize array of sampled parameters,  $\Psi = \{\}$ .
3:   for  $l \in \{1, \dots, L\}$  do
4:      $\Xi_l^0 \sim \varphi(\Xi|K_m)$ 
5:   end for
6:   for  $\tau \in \{1, \dots, \mathcal{T}_0 + \mathcal{T}_1\}$  do
7:     for  $l \in \{1, \dots, L\}$  do
8:       Propose the following state,  $\Xi' = \Xi_l^{\tau-1} + \epsilon \times \text{Uniform}(-1, 1)$ .
9:       Calculate the acceptance ratio,  $\alpha = p(\Xi'|\mathcal{D}, K_m, \beta_l) / p(\Xi_l^{\tau-1}|\mathcal{D}, K_m, \beta_l)$ .
10:      if  $\text{Uniform}(0, 1) < \alpha$  then
11:         $\Xi_l^\tau = \Xi'$ 
12:      else
13:         $\Xi_l^\tau = \Xi_l^{\tau-1}$ 
14:      end if
15:    end for
16:    for  $l \in \{1, \dots, L-1\}$  do
17:      Calculate the probability of exchanging states,  $p(\Xi_l^\tau \leftrightarrow \Xi_{l+1}^\tau)$ .
18:      if  $\text{Uniform}(0, 1) < p(\Xi_l^\tau \leftrightarrow \Xi_{l+1}^\tau)$  then
19:        Swap the  $\Xi_l^\tau$  for the  $\Xi_{l+1}^\tau$ .
20:      end if
21:    end for
22:    if  $\tau > \mathcal{T}_0$  then
23:      Append the  $\{\Xi_l^\tau\}_{l=1}^L$  to the  $\Psi$ .
24:    end if
25:  end for
26:  Calculate the marginal likelihood  $Z(K_m)$  and the Bayesian free energy  $F(K_m)$ .
27: end for
28: Calculate the likelihood of the model against the measured data  $\{p(K_m|\mathcal{D})\}_{m=1}^M$ .

```

2.1. Bayesian model selection

The process of generating experimental measurement data is generally described by a probabilistic model that includes noise components. The SAS measurement data consist of scattering intensities that correspond to the scattering vector. As the scattering intensity is a measure of the number of incident photons on the detector, the scattering intensity values are assumed to follow a Poisson distribution (Katakami *et al.*, 2022; Kirian *et al.*, 2011; Liebi *et al.*, 2015; Nagata *et al.*, 2019). Let $I_K(q, \Xi)$ be the mathematical model of scattering intensity characterized by the parameter K for sample parameters Ξ and scattering vector magnitude q . The likelihood, which is the probability of generating the measured value y , is then given by

$$p(y|q, \Xi, K) = \frac{I_K(q, \Xi)^y \exp[-I_K(q, \Xi)]}{y!}. \quad (1)$$

Assuming that the measurement data $\mathcal{D} = \{q_i, y_i\}_{i=1}^N$, consisting of N data points, are samples from an independent and identically distributed population under K and Ξ , the likelihood is expressed by

$$p(\mathcal{D}|\Xi, K) = \prod_{i=1}^N p(y_i|q_i, \Xi, K). \quad (2)$$

Here, we introduce the Poisson cost function to transform the likelihood of the measured data in equation (2) as

$$E(\Xi, K) = \frac{1}{N} \sum_{i=1}^N \left[I_K(q_i, \Xi) - y_i \log I_K(q_i, \Xi) + \sum_{j=1}^{y_i} \log j \right]. \quad (3)$$

The likelihood is thus expressed as

$$p(\mathcal{D}|\Xi, K) = \exp[-NE(\Xi, K)]. \quad (4)$$

Let $\varphi(K)$ be the prior distribution of the parameter K that characterizes the model, and $\varphi(\Xi|K)$ be the prior distribution of the model parameters Ξ . Then, from Bayes' theorem, the posterior distribution of the parameters given the measurement data can be written as

$$p(\Xi|\mathcal{D}, K) = \frac{p(\mathcal{D}|\Xi, K) \varphi(\Xi|K) \varphi(K)}{\int p(\mathcal{D}, \Xi, K) d\Xi} \quad (5)$$

$$= \frac{\exp[-NE(\Xi, K)] \varphi(\Xi|K)}{Z(K)}, \quad (6)$$

$$Z(K) = \int \exp[-NE(\Xi, K)] \varphi(\Xi|K) d\Xi. \quad (7)$$

$Z(K)$ is the marginal likelihood, which corresponds to the normalization constant of the posterior parameter distribution. The probability of model K given the data \mathcal{D} , denoted $p(K|\mathcal{D})$, is given by

$$p(K|\mathcal{D}) = \frac{\int p(\mathcal{D}, \Xi, K) d\Xi}{\sum_K \int p(\mathcal{D}, \Xi, K) d\Xi} \quad (8)$$

$$= \frac{\int \exp[-NE(\Xi, K)] \varphi(\Xi|K) \varphi(K) d\Xi}{\sum_K \int \exp[-NE(\Xi, K)] \varphi(\Xi|K) \varphi(K) d\Xi} \quad (9)$$

$$= \frac{\exp[-F(K)] \varphi(K)}{\sum_K \exp[-F(K)] \varphi(K)}, \quad (10)$$

where

$$F(K) = -\log Z(K). \quad (11)$$

$F(K)$ is referred to as the Bayesian free energy, also known as the stochastic complexity. The posterior probability of the model, $p(K|\mathcal{D})$, can be rephrased as the validity of model K for the measurement data \mathcal{D} . In other words, calculating and comparing the value of $p(K|\mathcal{D})$ for all candidate models $\{K\}$ thus enables quantitative model selection. Note that in Bayesian model selection the parameter K does not need to appear explicitly within the mathematical model of the specimen. This means that the method is also applicable to

analyses such as comparing models with different sample shapes, like cylinders and ellipsoids, or evaluating the validity of spherical versus cylindrical models for the aspect ratios of colloidal particles.

2.2. Calculation of marginal likelihood

In our Bayesian model selection method, the Bayesian free energy $F(K)$ and the probability $p(K|\mathcal{D})$ are calculated and compared for all candidate models. This computation relies on determining the marginal likelihood $Z(K)$, as expressed in equation (7). The marginal likelihood generally involves multi-dimensional integration, which can be computationally intensive and unstable. To address this challenge, our framework uses replica-exchange Monte Carlo (REMC) to calculate the marginal likelihood (Hukushima & Nemoto, 1996). This method facilitates sampling from the desired probability distribution at multiple inverse temperatures, referred to as replicas, using the Markov-chain Monte Carlo method (MCMC) to exchange states strategically between adjacent inverse temperatures at arbitrary intervals, thus avoiding local minima. To calculate the marginal likelihood using REMC, we establish a series of L inverse temperatures $\{\beta_l\}_{l=1}^L$ that satisfy the relation

$$0 = \beta_1 < \dots < \beta_L = 1. \quad (12)$$

Sampling from the joint probability distribution at each inverse temperature gives

$$p(\Xi_1, \dots, \Xi_L|\mathcal{D}, K, \beta_1, \dots, \beta_L) = \prod_{l=1}^L p(\Xi_l|\mathcal{D}, K, \beta_l), \quad (13)$$

where Ξ_l denotes the model parameter at the l th inverse temperature β_l . The posterior distribution $p(\Xi_l|\mathcal{D}, K, \beta_l)$ satisfies the following relation:

$$p(\Xi_l|\mathcal{D}, K, \beta_l) \propto \exp[-N\beta_l E(\Xi_l, K)] \varphi(\Xi_l|K). \quad (14)$$

These distributions are sampled using MCMC at each inverse temperature, as expressed in equation (14), and states at adjacent inverse temperatures are periodically exchanged with a probability that satisfies the detailed balance condition. The probability of exchanging the l th and $(l+1)$ th states, $p(\Xi_l \leftrightarrow \Xi_{l+1})$, is

$$p(\Xi_l \leftrightarrow \Xi_{l+1}) = \min \left[1, \frac{p(\Xi_{l+1}|\mathcal{D}, K, \beta_l) p(\Xi_l|\mathcal{D}, K, \beta_{l+1})}{p(\Xi_l|\mathcal{D}, K, \beta_l) p(\Xi_{l+1}|\mathcal{D}, K, \beta_{l+1})} \right] \quad (15)$$

$$= \min \left(1, \exp \left\{ N(\beta_{l+1} - \beta_l) \times [E(\Xi_{l+1}, K) - E(\Xi_l, K)] \right\} \right). \quad (16)$$

The marginal likelihood expressed in equation (7) can be efficiently determined using samples from various inverse temperatures sampled by REMC. The marginal likelihood $Z(K, \beta)$ at inverse temperature β is expressed as

$$Z(K, \beta) = \int \exp[-N\beta E(\Xi, K)] \varphi(\Xi|K) d\Xi. \quad (17)$$

In this case, the target marginal likelihood expressed in equation (7) is equivalent to $Z(K, \beta = 1)$. Using the relation in equation (12), $Z(K, \beta = 1)$ can be expressed as follows:

$$Z(K, \beta = 1) = \prod_{l=1}^{L-1} \frac{Z(K, \beta_{l+1})}{Z(K, \beta_l)} \quad (18)$$

$$= \prod_{l=1}^{L-1} \left(\exp[-N(\beta_{l+1} - \beta_l) E(\Xi, K)] \right)_{p(\Xi_l | \mathcal{D}, K, \beta_l)} \quad (19)$$

In equation (19), the symbol $\langle \cdot \rangle_{p(\Xi_l | \mathcal{D}, K, \beta_l)}$ denotes the expectation value with respect to $p(\Xi_l | \mathcal{D}, K, \beta_l)$. Computing equation (19) using sampling with REMC provides the marginal likelihood expressed in equation (7). Once the marginal likelihood $Z(K)$ is determined, we can find the Bayesian free energy expressed in equation (11) and the posterior probability of model K given the measurement data expressed in equation (10). For the numerical experiments presented here, we used the Metropolis method (Metropolis *et al.*, 1953) for MCMC sampling of the posterior distributions at each inverse temperature, as expressed in equation (14).

2.3. Estimation of model parameters

During the marginal likelihood calculation the posterior distribution of $p(\Xi_L | \mathcal{D}, K, \beta_L = 1)$ is obtained, which simply represents the Bayesian estimate of the model parameters (Hayashi *et al.*, 2023). Therefore, the parameter estimation is conducted simultaneously with performing the Bayesian model selection. Since the posterior distribution is sampled using REMC sampling, it can provide a global parameter estimate solution. The reliability of the estimation can be assessed from the statistical properties of the sampled posterior distribution.

In Bayesian estimation, the maximum *a posteriori* (MAP) solution provides a point estimate of the parameters. The MAP solution Ξ_{MAP} for the parameters of model K is expressed by this equation from equation (14):

$$\Xi_{\text{MAP}} = \underset{\Xi}{\operatorname{argmax}} \exp[-NE(\Xi, K)] \varphi(\Xi | K). \quad (20)$$

3. Formulation of a multicomponent monodisperse spheres model

In this section, we describe a model for the scattering intensity of a dilute sample comprising multicomponent monodisperse spheres (Guinier & Fournet, 1955; Hashimoto, 2022). This model serves as the basis for evaluating the performance of the proposed method.

Let \mathbf{e}_i and \mathbf{e}_s represent the unit vectors in the direction of the wavevector of the incident and scattered beams, respectively. If \mathbf{e}_i and \mathbf{e}_s form an angle 2θ , and the wavelength of the beam is λ , then the scattering vector \mathbf{q} is given by

$$\mathbf{q} = \frac{2\pi}{\lambda} (\mathbf{e}_s - \mathbf{e}_i). \quad (21)$$

In this paper, we consider isotropic scattering and focus on the scattering vector's magnitude q , defined as

$$q = |\mathbf{q}| = \frac{4\pi}{\lambda} \sin \theta. \quad (22)$$

Monodisperse spheres are spherical particles of uniform radius. The scattering intensity $I(q, \xi)$ of a specimen composed of sufficiently dilute monodisperse spheres of a single type for the scattering vector magnitude q is given by

$$I(q, \xi) = SV \left[\frac{\sin(qR) - qR \cos(qR)}{(qR)^3} \right]^2 + B, \quad (23)$$

where $V = \frac{4}{3}\pi R^3$. If the difference in scattering length density between the solute and solvent of the specimen is $\Delta\rho$ and the volume fraction is ϕ , then $S = (3\Delta\rho)^2\phi$. The parameters ξ of this model are the particle size R , the scale S and the background B .

To formulate the scattering intensity of a specimen composed of K types of monodisperse sphere, we assume a dilute system and denote the particle size of the k th component in the sample as R_k and the scale as S_k . The scattering intensity of a sample composed of K types of monodisperse sphere is then given by

$$I_K(q, \Xi) = \sum_{k=1}^K S_k V_k \left[\frac{\sin(qR_k) - qR_k \cos(qR_k)}{(qR_k)^3} \right]^2 + B, \quad (24)$$

where we assume that $V_k = \frac{4}{3}\pi R_k^3$. The model parameters Ξ for the scattering intensity $I_K(\cdot)$ are $\Xi = \{R_k, S_k\}_{k=1}^K, B$.

4. Numerical experiments

Here, we present numerical experiments to evaluate the model selection among models with K ranging from one to four components to demonstrate the capabilities of the proposed framework. We apply the framework to synthetic data generated to represent a system with two types ($K = 2$) of monodisperse sphere, as described by equation (24). Bi-component spherical specimens, depicted in Fig. 1, correspond to model scenarios where two types of particle differing in size

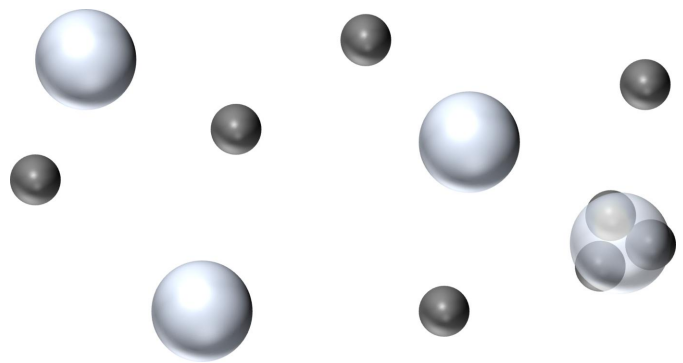


Figure 1
An illustration of a mixture of two types of spherical specimen. This shows scenarios with two components ($K = 2$), including mixtures of spherical particles of different sizes or volume fractions, and aggregates from a single particle type approximated as a large sphere.

or volume fraction are mixed, or cases in which particles of a single type aggregate into larger spherical entities.

In typical SAS experiments, the scale parameter S_k in equation (24) tends to be small. Therefore, we normalize the scale parameter S_k as

$$\bar{S}_k = S_k \times 10^8. \quad (25)$$

Accordingly, we refer to the model parameters as $\Xi = \{R_k, \bar{S}_k\}_{k=1}^K, B\}$.

The numerical experiments reported in this section were conducted with a burn-in period of 10^5 and a sample size of 10^5 for the REMC. We set the number of replicas for REMC, the values of inverse temperature and the step size of the Metropolis method taking into consideration the state exchange rate and the acceptance rate.

4.1. Generation of synthetic data

The scattering intensity in SAS experiments, which is typically recorded as count data, is subject to Poisson noise, as described by equation (1). We therefore generated synthetic data \mathcal{D} using the following procedure:

(i) Set the number of data points to $N = 400$ and define the scattering vector magnitudes at N equally spaced points within the interval $[0.1, 3]$ to obtain $\{q_i\}_{i=1}^{N=400}$ (nm^{-1}).

(ii) Assume $K = 2$ and set the true model parameters Ξ to Ξ^* .

(iii) Calculate the scattering intensity at the scattering vector magnitudes $\{q_i\}_{i=1}^N$ obtained in Step (i), using the model in equation (24) and Ξ^* . Introduce a pseudo-measurement time T to adjust the noise intensity in the data, to obtain $\{I(q_i, \Xi^*, K) T\}_{i=1}^N$.

(iv) Generate measurement values $\{y_i\}_{i=1}^N$ as Poisson-distributed random numbers with means of $\{I(q_i, \Xi^*, K) \times T\}_{i=1}^N$ to create the synthetic data set $\mathcal{D} = \{q_i, y_i\}_{i=1}^N$.

In this section, we consider cases with pseudo-measurement times of $T = 1$ and $T = 0.1$. Generally, smaller values of T indicate greater effects from measurement noise.

4.2. Setting the prior distributions

In the Bayesian model selection framework, prior knowledge concerning the parameters Ξ and the model-characterizing parameter K is set as their prior distributions.

In this numerical experiment, the prior distributions for the parameters Ξ were set as Gamma distributions based on the pseudo-measurement time T used during data generation, while the prior for K was a discrete uniform distribution over the interval $[1, 4]$.

$$\varphi(\Xi|K) = \varphi(B) \prod_{k=1}^K \varphi(R_k) \varphi(\bar{S}_k), \quad (26)$$

$$\varphi(R_k) = \text{Gamma}(R_k; \alpha = 1.2, \beta = 20) \quad (27)$$

$$= \frac{\exp(-R_k/\beta)}{\beta^\alpha \Gamma(\alpha)} R_k^{\alpha-1}, \quad (28)$$

$$\varphi(\bar{S}_k) = \begin{cases} \text{Gamma}(\bar{S}_k; 1.05, 300 \times T) & \text{if } \{\bar{S}_k\}_{k=1}^K \text{ is in} \\ & \text{descending order,} \\ 0 & \text{otherwise,} \end{cases} \quad (29)$$

$$\varphi(B) = \text{Gamma}(B; 1.05, 0.02 \times T), \quad (30)$$

$$\varphi(K) = \text{DiscreteUniform}(K; 1, 4). \quad (31)$$

Fig. 2 shows the prior distributions for each parameter, as described in equations (28), (29), (30) and (31).

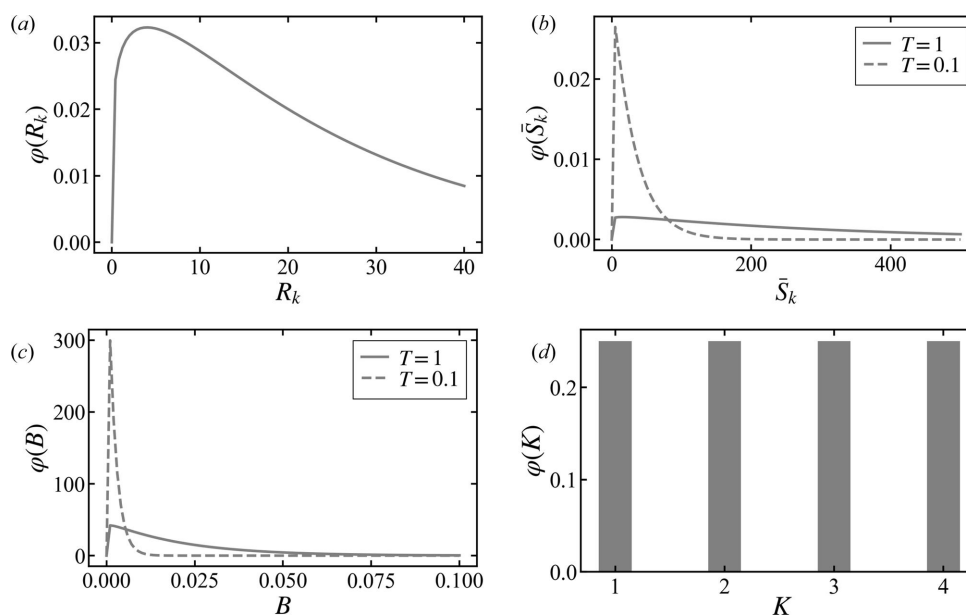


Figure 2

Plots of the prior distributions for various parameters. (a) Prior distribution of R_k , $\varphi(R_k)$. (b) Prior distribution of \bar{S}_k , $\varphi(\bar{S}_k)$. (c) Prior distribution of B , $\varphi(B)$. (d) Prior distribution of K , $\varphi(K)$.

4.3. Results for two-component monodisperse spheres based on scale ratio

The ratio of the scale parameters S_1 and S_2 for spheres 1 and 2 during data generation, denoted r_S , is defined as

$$r_S = \frac{\bar{S}_2}{\bar{S}_1}. \quad (32)$$

Next, we present the results from applying our proposed method to analyzing six types of data generated by varying the

value of r_S for pseudo-measurement times of $T = 1$ and 0.1 . Table 1 displays the parameter values used to generate the synthetic data.

Fig. 3 displays the fitting results and residual plots for synthetic data generated with the parameter values from Table 1. For each model ($K = 1, K = 2, K = 3$ and $K = 4$), 1000 samples were randomly selected from their respective posterior parameter distributions to plot these curves. Here, the residual Δ , which normalizes the difference between the model predictions and the observed data points (q, y) using the model parameters Ξ , is given by

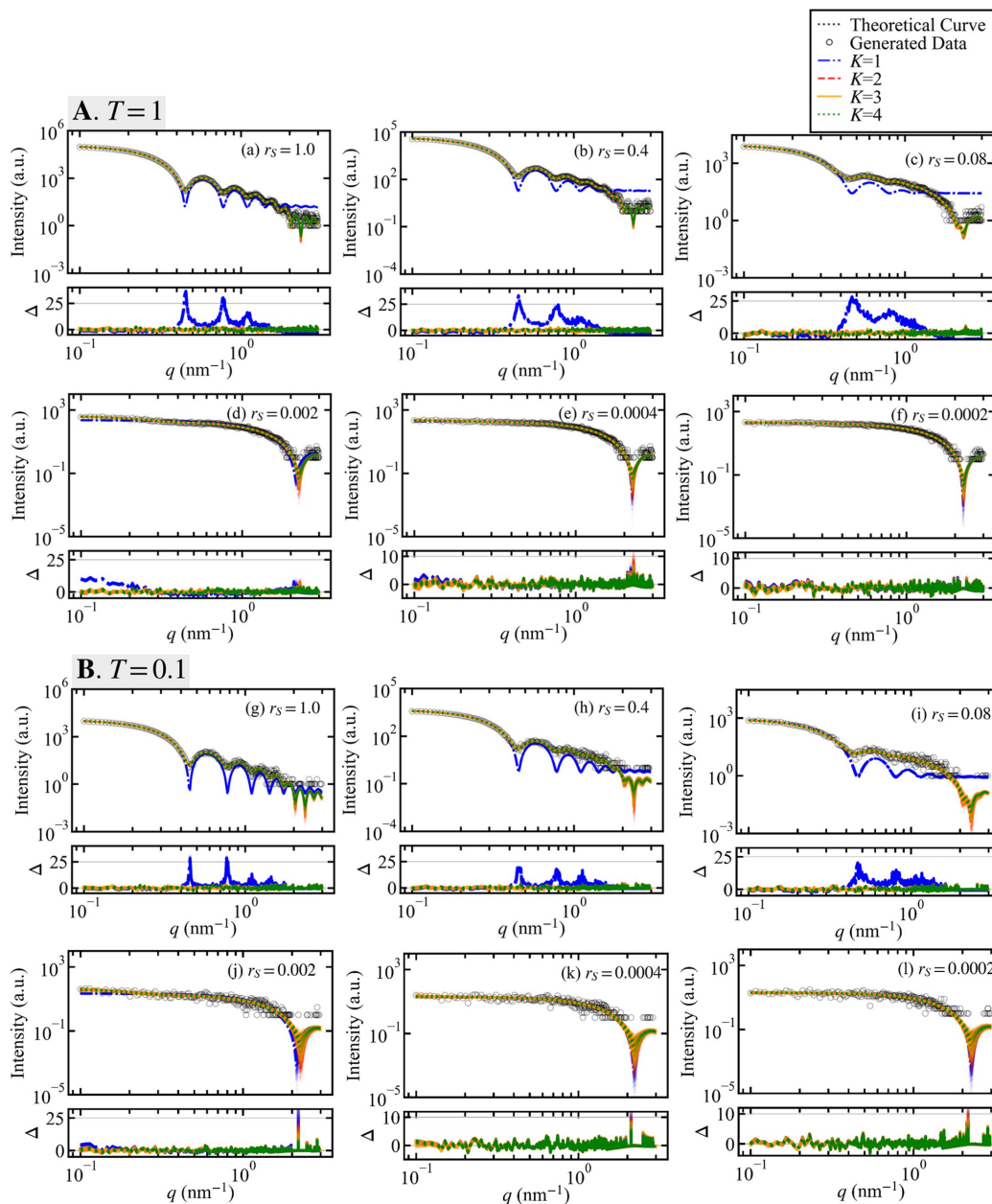


Figure 3 Fitting to synthetic data generated at various r_S values and residual plots. Panels **A** and **B** show cases for pseudo-measurement times of $T = 1$ and $T = 0.1$, respectively. In plots (a)–(f) and (g)–(l), the scale ratio r_S is displayed in descending order for $T = 1$ and $T = 0.1$, respectively. Black circles represent the generated data and the black dotted lines indicate the true scattering intensity curves. For models $K = 1, K = 2, K = 3$ and $K = 4$, the fitting curves and residual plots are represented by blue dashed–dotted lines, red dashed lines, orange solid lines and green dotted lines, respectively. Fitting curves were plotted using 1000 parameter samples that were randomly selected from the posterior probability distributions for each model. The width of the distribution of these fitting curves reflects the confidence level at each point.

Table 1

Parameter values used for data generation with varying r_S .

	Sphere 1	Sphere 2
Radius R (nm)	2	10
Scale \bar{S}	250	{250, 100, 20, 0.5, 0.1, 0.05}
Background B (cm ⁻¹)	0.01	
Pseudo-measurement time T	{1, 0.1}	

$$\Delta = \frac{I_K(q, \Xi) - y}{I_K(q, \Xi)}. \quad (33)$$

The fitting curves in Fig. 3 illustrate that the intensity and spread of these curves are indicative of confidence levels, where darker areas and narrower spreads denote higher confidence levels for the respective model.

In Fig. 3 the plots labeled (a) to (c) and (g) to (i) demonstrate, through the residual plots, that the model with $K = 1$ predominantly fails to represent the data accurately. However, we can also see that the fitting curves for models with $K = 2-4$ are almost identical in shape. The data shown in plots (d)–(f) and (j)–(l) are difficult to distinguish from the well known scattering data of a single type of monodisperse sphere ($K = 1$), making it challenging to qualitatively compare the goodness of fit among the models with $K = 1-4$.

Fig. 4 presents the Bayesian model selection results using our proposed framework. Within Fig. 4, panel **A** contains results for the case with $T = 1$ and panel **B** contains results with $T = 0.1$, each showing the probability $p(K|\mathcal{D})$ of model K based on the synthetic data \mathcal{D} for each scale ratio r_S . Here, ten data sets were created for each parameter value by varying the

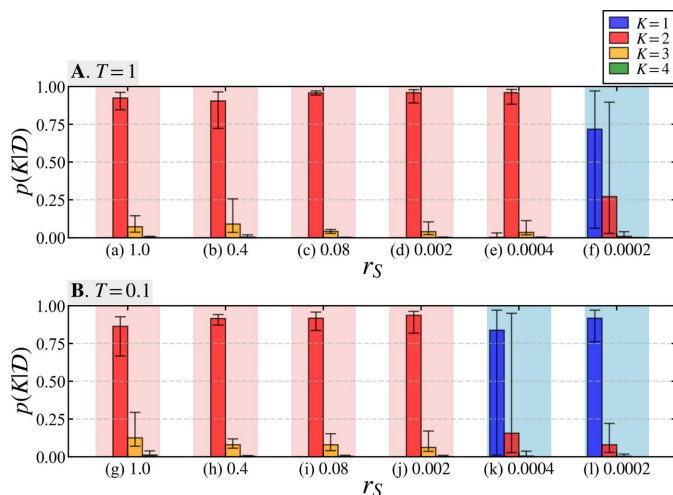


Figure 4

Results of Bayesian model selection among models $K = 1-4$ for varying r_S values. Panel **A** shows the posterior probability for each model using data generated with a pseudo-measurement time of $T = 1$, and panel **B** shows results for $T = 0.1$. In cases (a)–(f) and (g)–(l), the scale ratio r_S is displayed in descending order for $T = 1$ and $T = 0.1$, respectively. The height of each bar corresponds to the average values calculated for ten data sets generated with different random seeds, with maximum and minimum values shown as error bars. Areas highlighted in red indicate cases where, on average, the highest probability was found for the true model with $K = 2$, while blue backgrounds indicate that models other than $K = 2$ were associated with the highest probability on average.

Table 2

The number of times each model was associated with the highest probability in numerical experiments for ten data sets generated with different random seeds at each r_S value.

In cases (a)–(f) and (g)–(l), the scale ratio r_S is displayed in descending order for $T = 1$ and $T = 0.1$, respectively. The most frequently counted case for each r_S value is shown in bold.

(**A**) $T = 1$

r_S	K			
	1	2	3	4
(a) 1.0	0	10	0	0
(b) 0.4	0	10	0	0
(c) 0.08	0	10	0	0
(d) 0.002	0	10	0	0
(e) 0.0004	0	10	0	0
(f) 0.0002	8	2	0	0

(**B**) $T = 0.1$

r_S	K			
	1	2	3	4
(g) 1.0	0	10	0	0
(h) 0.4	0	10	0	0
(i) 0.08	0	10	0	0
(j) 0.002	0	10	0	0
(k) 0.0004	9	1	0	0
(l) 0.0002	10	0	0	0

random seed during data generation, and the average value of $p(K|\mathcal{D})$ is indicated by the height of the bar graph, with error bars indicating the maximum and minimum values. For the relatively large scale ratios r_S in plots (a)–(e) in panel **A**, the true model with $K = 2$ has a high probability, while the average value of $p(K|\mathcal{D})$ is highest for $K = 1$ in plot (f). In panel **B**, the true model with $K = 2$ is associated with high probability in cases (g)–(j), while in cases (k) and (l) $K = 1$ is associated with the highest probability.

Table 2 summarizes the number of times each model was found to have the highest probability in numerical experiments using the ten separate data sets shown in Fig. 4. For values of $r_S = 0.0004$ and above (Table 2, part **A**) and for $r_S = 0.002$ and above (Table 2, part **B**), the model with $K = 2$ was associated with the highest probability in all ten data sets. This demonstrates the high accuracy of the proposed method and its robustness to measurement noise. In cases (f), (k) and (l) of Fig. 4 and Table 2, the model with $K = 1$ was found to have the highest probability in nearly all of the ten data sets. These results were used to inform a discussion of the suitable analysis range of r_S using the proposed method, as addressed in the next section.

4.4. Results for two-component monodisperse spheres based on radius ratio

During synthetic data generation, the ratio of the radii R_1 and R_2 of spheres 1 and 2, denoted r_R , was defined as

$$r_R = \frac{R_1}{R_2}. \quad (34)$$

In this setup, we generated seven types of data by varying the value of r_R for pseudo-measurement times of $T = 1$ and $T = 0.1$.

Fig. 5 displays fitting curves and residual plots for the models $K = 1$, $K = 2$, $K = 3$ and $K = 4$, calculated from 1000 samples randomly selected from their respective posterior parameter distributions. These samples are derived from synthetic data generated with the parameter values given in Table 3. The residuals Δ were calculated using equation (33).

Aside from the cases of $r_R = 0.5$ in plots (d) and (k), the profiles of the data in Fig. 5 are very similar to those of a single monodisperse sphere, and the fitting curves for models $K = 1$ to $K = 4$ are nearly identical in shape. In contrast, the data for

Table 3

Parameter values used for data generation when varying r_R .

	Sphere 1	Sphere 2
Radius R (nm)	{9.9, 9.7, 9.5, 0.5, 0.5, 0.4, 0.3}	10
Scale \bar{S}	250	100
Background B (cm ⁻¹)	0.01	
Pseudo-measurement time T	{1, 0.1}	

cases (d) and (k) with $r_R = 0.5$ have a complex profile, and the model with $K = 1$ represents the data poorly.

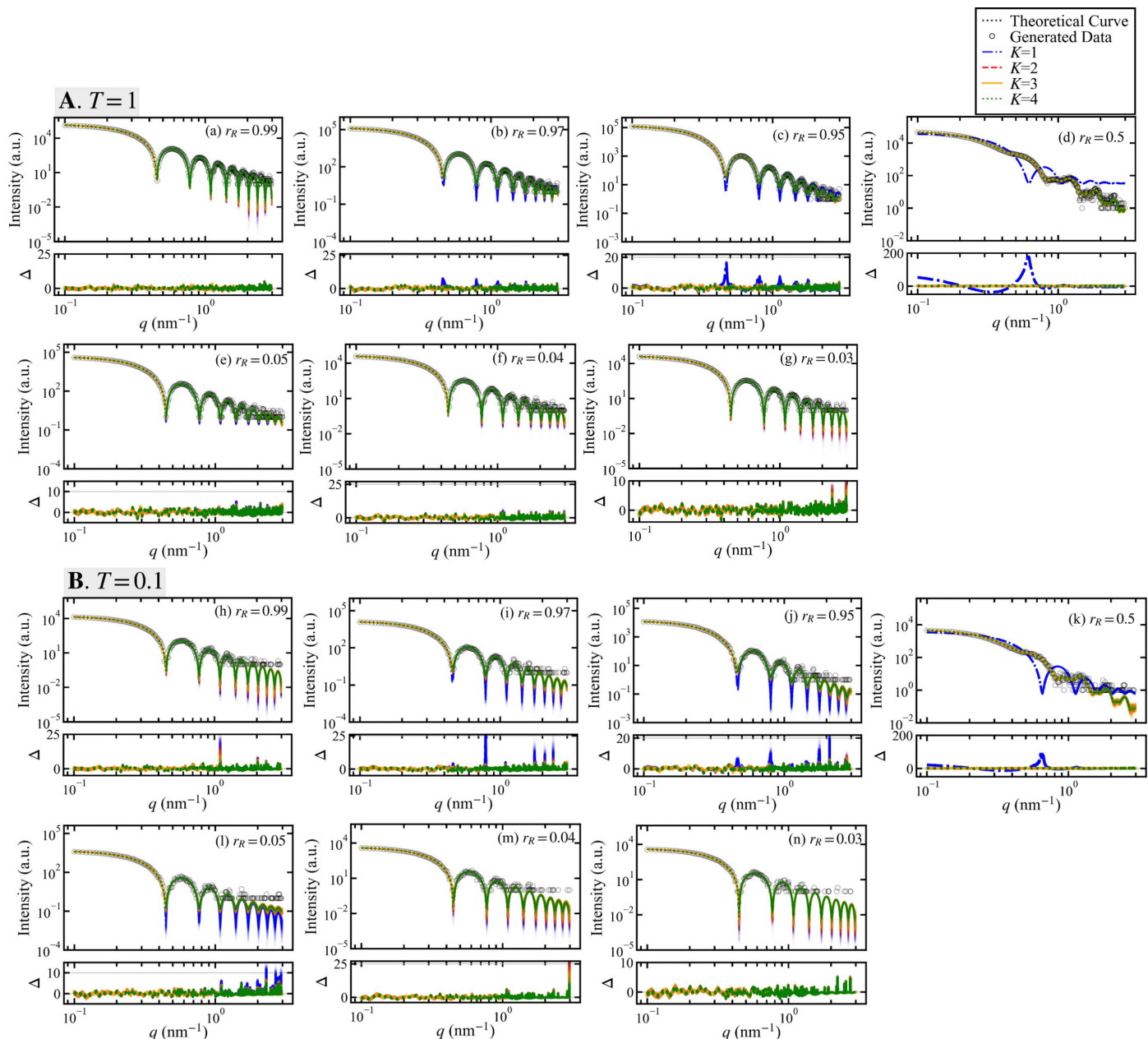


Figure 5

Fitting to synthetic data generated at various r_R values and residual plots. Panels **A** and **B** show cases for pseudo-measurement times of $T = 1$ and $T = 0.1$, respectively. In plots (a)–(g) and (h)–(n), the radius ratio r_R is displayed in descending order for $T = 1$ and $T = 0.1$, respectively. Black circles represent the generated data and the black dotted lines indicate the true scattering intensity curves. For models $K = 1$, $K = 2$, $K = 3$ and $K = 4$, the fitting curves and residual plots are represented by blue dashed–dotted lines, red dashed lines, orange solid lines and green dotted lines, respectively. Fitting curves were plotted using 1000 parameter samples that were randomly selected from the posterior probability distributions for each model. The width of the distribution of these fitting curves reflects the confidence level at each point.

Table 4

The number of times each model was most highly supported in numerical experiments for ten data sets generated by varying r_R values.

In cases (a)–(g) and (h)–(n), the radius ratio r_R is displayed in descending order for $T = 1$ and $T = 0.1$, respectively. The cases with the highest counts for each r_R value are shown in bold.

(A) $T = 1$

r_R	K			
	1	2	3	4
(a) 0.99	9	1	0	0
(b) 0.97	0	10	0	0
(c) 0.95	0	10	0	0
(d) 0.5	0	10	0	0
(e) 0.05	0	10	0	0
(f) 0.04	1	9	0	0
(g) 0.03	10	0	0	0

(B) $T = 0.1$

r_R	K			
	1	2	3	4
(h) 0.99	10	0	0	0
(i) 0.97	2	8	0	0
(j) 0.95	0	10	0	0
(k) 0.5	0	10	0	0
(l) 0.05	1	9	0	0
(m) 0.04	7	3	0	0
(n) 0.03	10	0	0	0

Fig. 6 displays the results of Bayesian model selection using synthetic data generated by varying the radius ratio r_R . Ten data sets were created for each parameter value by varying the random seed during data generation. The average value of

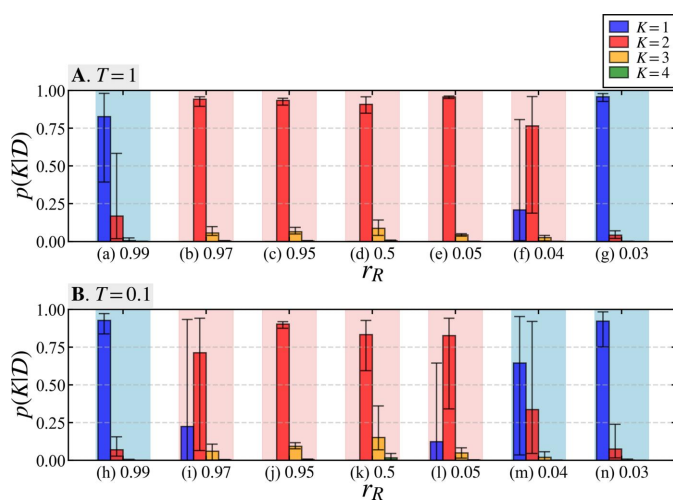


Figure 6

Results of Bayesian model selection among models $K = 1$ – 4 for varying r_R values. Panel **A** shows the posterior probability of each model using data generated with a pseudo-measurement time of $T = 1$, and panel **B** shows results for $T = 0.1$. In cases (a)–(g) and (h)–(n), the radius ratio r_R is displayed in descending order for $T = 1$ and $T = 0.1$, respectively. The height of each bar corresponds to the average values calculated for ten data sets generated with different random seeds, with the maximum and minimum values shown as error bars. Areas highlighted in red indicate cases where the true model $K = 2$ was most highly supported, while the blue backgrounds indicate that the likelihood of a model other than $K = 2$ was the highest.

$p(K|\mathcal{D})$ is indicated by the height of the bar graph, with the maximum and minimum values shown as error bars. Unlike the results for the variations in scale ratio shown in Fig. 4, the model selection procedure fails not only at a radius ratio r_R close to 0 but also at values close to 1, with $K = 1$ being the most highly supported. In the case of $r_R = 0.04$, the result for $T = 1$ in case (f) supports the true model $K = 2$, but for $T = 0.1$ in case (m) the alternative model $K = 1$ is the most highly supported. However, in cases (b)–(f) and (i)–(l), the true model $K = 2$ is associated with a high average probability (Fig. 6).

Table 4 presents the results of numerical experiments for the ten separate data sets shown in Fig. 6, summarizing the number of times each model $K = 1$ – 4 was most highly supported. Near the analytical limits of the proposed method, there are cases where the supported model changes depending on the data, as shown in Table 4, entries (a), (f), (i), (l) and (m).

5. Discussion

In Section 4, we conducted numerical experiments to determine the number of components K in two-component monodisperse sphere specimens using the proposed method through model selection applied to artificial measurement data. In this section, we discuss the analytical limits of our method under the settings of this study with respect to the scale ratio r_S and radius ratio r_R of the specimen's two components, as well as the performance of the conventional model selection method based on the reduced χ -squared error.

5.1. Limitations of the proposed method

The experiments detailed in Section 4 explored the selection of the number of components K for two-component monodisperse spheres using the proposed Bayesian method. We observed certain analytical limitations for various values of the scale ratio r_S and radius ratio r_R . In practical data analysis applications using the proposed method, it is advisable to conduct preliminary tests using synthetic data with noise intensity and anticipated parameter values similar to those of the measured data. This step can help ensure a more reliable analysis, as detailed below.

The scale parameter S is a value that is multiplied by the square of the difference in scattering length density between the solvent and the specimen, as well as the volume fraction. This can cause r_S to become extremely small when there is little difference in scattering length density between the solvent and a component of the specimen, or when there is a significant difference in the mixing ratio of the components. The results in Fig. 4 and Table 2 for a pseudo-measurement time of $T = 1$ (panel **A**) indicate that the model selection favored non-true models at a scale ratio of $r_S = 0.0002$. Similarly, for $T = 0.1$ (panel **B**), non-true models were favored at scale ratios of $r_S = 0.0004$ and $r_S = 0.0002$, indicating that these cases exceed the analytical capabilities of the proposed

method. These findings imply that, within the experimental parameters of this study, the proposed method reliably identifies the true model with a high probability for scale ratios up to $r_S = 0.0004$ at $T = 1$ and up to $r_S = 0.002$ at $T = 0.1$.

In Section 4.4, we investigated the effect of varying the radius ratio r_R . When components of different radii are mixed, it is important to consider not only simple mixtures but also instances of aggregated specimens. The findings shown in Fig. 6 and Table 4 indicate that the proposed method reaches its analytical limits as r_R approaches 1 and as it approaches 0. As r_R nears 1, the scattering profiles of the two-component system become similar to that of a single-component system, leading to the selection of the single-component model ($K = 1$). We found an analytical limit at $r_R = 0.99$ for both $T = 1$ and $T = 0.1$. The results for $r_R = 0.97$ show that at $T = 0.1$, which has a higher noise intensity than $T = 1$, the posterior probability of the single-component model ($K = 1$) increases, resulting in an unstable analysis. Conversely, as r_R approaches 0 with the results $r_R = 0.03$ at $T = 1$ and $r_R = 0.04$ and $r_R = 0.03$ at $T = 0.1$, the single-component model ($K = 1$) is associated with high probability, indicating an analytical limit. Overall, the proposed method demonstrates the ability to select the true model with high probability for radius ratios ranging from $r_R = 0.04$ to 0.99 at $T = 1$, and from $r_R = 0.05$ to 0.99 at $T = 0.1$.

5.2. Model selection based on χ -squared error

In SAS data analysis, selecting an appropriate mathematical model for the analysis is a crucial but challenging process. In this subsection, we compare the conventional model selection method based on the χ -squared error with the results of model selection using our proposed method.

Conventionally, model selection is performed by minimizing the χ -squared error through fitting with candidate models, and the model with a reduced χ -squared closest to 1 is considered to be the best representation of the data (Pedersen, 1997). The χ -squared error χ^2 is given by the following equation:

$$\chi^2 = \sum_{i=1}^N \frac{[I_K(q_i, \Xi) - y_i]^2}{I_K(q_i, \Xi)}. \quad (35)$$

The reduced χ -squared χ_r^2 is obtained by dividing the χ -squared error by the degrees of freedom, dof:

$$\chi_r^2 = \frac{\chi^2}{\text{dof}}. \quad (36)$$

The degrees of freedom dof are calculated by subtracting the number of model parameters from the number of data points N . For the model represented by equation (24), it is given by

$$\text{dof} = N - (2K + 1). \quad (37)$$

In the following, we discuss the results of selecting the model with the closest reduced χ -squared χ_r^2 to 1 for models $K = 1$ – 4 using the same data generated with different random seeds for each of the six types of r_S determined by the parameters shown in Table 1 for $T = 1$, as in Section 4.3. Since it is difficult to obtain a global optimum solution using conven-

Table 5
Model selection results based on reduced χ -squared values.

The table shows the number of times each model had the closest reduced χ -squared value to 1 for ten data sets generated with different random seeds for each r_S setting $T = 1$. Labels (a) to (f) refer to the settings in Figs. 3–4 and Table 2. The cases with the highest level of support for each data set are shown in bold.

r_S	K			
	1	2	3	4
(a) 1.0	0	2	8	0\sim
(b) 0.4	0	0	9	1
(c) 0.08	0	0	9	1
(d) 0.002	0	0	10	0
(e) 0.0004	0	4	5	1
(f) 0.0002	0	2	8	0

tional fitting methods such as the quasi-Newton method and conjugate gradient method, we evaluate the reduced χ -squared χ_r^2 using the parameters that minimize χ^2 among the parameters sampled from the posterior distribution in the experiment of Section 4.3.

First, we present the results for the data shown in Fig. 3 plot (a), generated using the model with $K = 2$. Fig. 7 shows the fitting results and reduced χ -squared for models $K = 1$ – 4 obtained by minimizing the χ -squared error for the data in Fig. 3 plot (a).

Fig. 7 shows the results obtained by the conventional method, indicating that the model with $K = 3$, even though it is not the true model, is considered most appropriate for the data due to its reduced χ -squared value being closest to 1. The fitting curves for models $K = 2, 3$ and 4 exhibit nearly identical shapes, which complicates the determination of the most suitable model based solely on their appearance.

Table 5 shows the aggregated results of calculating the reduced χ -squared for models $K = 1$ – 4 and counting the

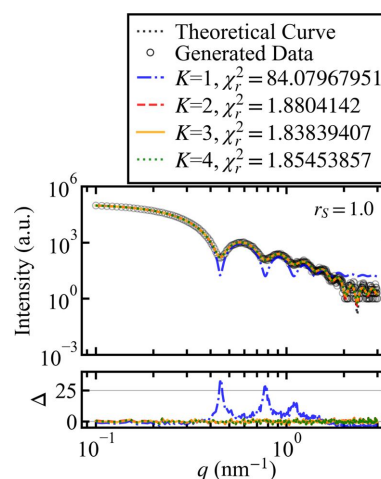


Figure 7
The fitting results and residual plots for the data shown in Fig. 3(a) were derived using parameters that minimize the χ -squared error from the posterior probability distributions for models ranging from $K = 1$ to $K = 4$. For each of these models, the fitting curves and their corresponding residual plots are represented by blue dashed–dotted lines, red dashed lines, orange solid lines and green dotted lines, respectively. The legend indicates the reduced χ -squared values for each model ($K = 1$ to $K = 4$).

number of times each model was closest to 1 for the data sets (a)–(f) generated by varying the scale ratio r_S of the model with $K = 2$ in Section 4.3.

The results shown in Table 5 indicate that the model with $K = 3$ is the most supported for all data sets (a)–(f). This is thought to be because minimizing the χ -squared error failed to address the noise in the data adequately, ultimately leading to an overestimation of the model's degrees of freedom. This implies that it is difficult to select the true model using the conventional method of comparing the χ_r^2 values among candidate models.

On the other hand, the results of the proposed method shown in Table 2 part A demonstrate that the true model $K = 2$ is supported ten out of ten times for cases (a)–(e). Within the analyzable range discussed in the previous subsection, the proposed method enables accurate model selection that takes into account the model degrees of freedom.

6. Conclusions

In this paper, we have introduced a Bayesian model selection framework for SAS data analysis that quantitatively evaluates model validity through posterior probabilities. We have conducted numerical experiments using synthetic data for a two-component system of monodisperse spheres to assess the performance of the proposed method.

We have identified the analytical limits of the proposed method, under the settings of this study, with respect to the scale and radius ratios of two-component spherical particles, and compared the performance of traditional model selection methods based on the reduced χ -squared.

The numerical experiments and subsequent discussion reveal the range of parameters that can be analyzed using the proposed method. Within that range, our method provides stable and highly accurate model selection, even for data with significant noise or in situations in which qualitative model determination is challenging. In comparison with the traditional method of selecting models based on fitting curves and data residuals, it was found that the proposed method offers greater accuracy and stability.

SAS is used to study specimens with a variety of structures other than spheres, including cylinders, core–shell structures, lamellae and more. The proposed method should be applied to other sample models to determine the feasibility of expanding the analysis beyond the case examined here to broader experimental settings. Future work could benefit from using the proposed method to conduct real data analysis, which is expected to yield new insights through our more efficient analysis approach.

Funding information

This work was supported by JST CREST (grant Nos. PMJCR1761 and JPMJCR1861) from the Japan Science and

Technology Agency (JST) and by a JSPS KAKENHI Grant-in-Aid for Scientific Research (A) (grant No. 23H00486).

References

- Antonov, L. D., Olsson, S., Boomsma, W. & Hamelryck, T. (2016). *Phys. Chem. Chem. Phys.* **18**, 5832–5838.
- Asahara, A., Morita, H., Ono, K., Yano, M., Mitsumata, C., Shoji, T. & Saito, K. (2021). AAAI Spring Symposium, 22–24 March 2021, held online. MLPS.
- Breßler, I., Kohlbrecher, J. & Thünemann, A. F. (2015). *J. Appl. Cryst.* **48**, 1587–1598.
- Da Vela, S. & Svergun, D. I. (2020). *Curr. Res. Struct. Biol.* **2**, 164–170.
- Guinier, A. & Fournet, G. (1955). *Small Angle Scattering of X-rays*. New York: Wiley.
- Hansen, S. (2000). *J. Appl. Cryst.* **33**, 1415–1421.
- Hashimoto, T. (2022). *Principles and Applications of X-ray, Light and Neutron Scattering*. Berlin: Springer Nature.
- Hayashi, Y., Katakami, S., Kuwamoto, S., Nagata, K., Mizumaki, M. & Okada, M. (2023). *J. Phys. Soc. Jpn.* **92**, 094002.
- Hukushima, K. & Nemoto, K. (1996). *J. Phys. Soc. Jpn.* **65**, 1604–1608.
- Kashiwamura, S., Katakami, S., Yamagami, R., Iwamitsu, K., Kumazoe, H., Nagata, K., Okajima, T., Akai, I. & Okada, M. (2022). *J. Phys. Soc. Jpn.* **91**, 074009.
- Katakami, S., Sakamoto, H., Nagata, K., Arima, T. H. & Okada, M. (2022). *Phys. Rev. E*, **105**, 065301.
- Kirian, R. A., Schmidt, K. E., Wang, X., Doak, R. B. & Spence, J. C. (2011). *Phys. Rev. E*, **84**, 011921.
- Kline, S. R. (2006). *J. Appl. Cryst.* **39**, 895–900.
- Larsen, A. H., Arleth, L. & Hansen, S. (2018b). *J. Appl. Cryst.* **51**, 1151–1161.
- Larsen, A. H., Dorosz, J., Thorsen, T. S., Johansen, N. T., Darwish, T., Midtgaard, S. R., Arleth, L. & Kastrup, J. S. (2018a). *IUCrJ*, **5**, 780–793.
- Larsen, A. H. & Pedersen, M. C. (2021). *J. Appl. Cryst.* **54**, 1281–1289.
- Liebi, M., Georgiadis, M., Menzel, A., Schneider, P., Kohlbrecher, J., Bunk, O. & Guizar-Sicairos, M. (2015). *Nature*, **527**, 349–352.
- Machida, A., Nagata, K., Murakami, R., Shinotsuka, H., Shouno, H., Yoshikawa, H. & Okada, M. (2021). *Sci. Technol. Adv. Mater.* **1**, 123–133.
- Metropolis, N., Rosenbluth, A. W., Rosenbluth, M. N., Teller, A. H. & Teller, E. (1953). *J. Chem. Phys.* **21**, 1087–1092.
- Moriguchi, R., Tsutsui, S., Katakami, S., Nagata, K., Mizumaki, M. & Okada, M. (2022). *J. Phys. Soc. Jpn.* **91**, 104002.
- Nagai, K., Anada, M., Nakanishi-Ohno, Y., Okada, M. & Wakabayashi, Y. (2021). *J. Appl. Cryst.* **54**, 1023.
- Nagata, K., Muraoka, R., Mototake, Y. I., Sasaki, T. & Okada, M. (2019). *J. Phys. Soc. Jpn.* **88**, 044003.
- Nagata, K., Sugita, S. & Okada, M. (2012). *Neural Netw.* **28**, 82–89.
- Nelson, A. R. J. & Prescott, S. W. (2019). *J. Appl. Cryst.* **52**, 193–200.
- Orioli, S., Larsen, A. H., Bottaro, S. & Lindorff-Larsen, K. (2020). *Prog. Mol. Biol. Transl. Sci.* **170**, 123–176.
- Pedersen, J. S. (1997). *Adv. Colloid Interface Sci.* **70**, 171–210.
- Rambo, R. P. & Tainer, J. A. (2013). *Nature*, **496**, 477–481.
- Rappel, H., Beex, L. A., Hale, J. S., Noels, L. & Bordas, S. P. A. (2020). *Arch. Comput. Methods Eng.* **27**, 361–385.
- Scheres, S. H. (2012). *J. Struct. Biol.* **180**, 519–530.
- Schneidman-Duhovny, D., Hammel, M. & Sali, A. (2010). *Nucleic Acids Res.* **38**, W540–W544.
- Svergun, D., Barberato, C. & Koch, M. H. J. (1995). *J. Appl. Cryst.* **28**, 768–773.
- Ueda, H., Katakami, S., Yoshida, S., Koyama, T., Nakai, Y., Mito, T., Mizumaki, M. & Okada, M. (2023). *J. Phys. Soc. Jpn.* **92**, 054002.