

Bayesian nonparametric estimation of Milky Way parameters using matrix-variate data, in a new Gaussian Process based method

Dalia Chakrabarty*

*Department of Statistics
University of Warwick
Coventry CV4 7AL, U.K.
d.chakrabarty@warwick.ac.uk*

*Department of Mathematics
University of Leicester
Leicester LE1 7RH, U.K.
dc252@le.ac.uk*

and

Munmun Biswas[†], Sourabh Bhattacharya[‡]

*Indian Statistical Institute
203, B. T. Road
Kolkata 700108, India
munmun.biswas08@gmail.com; sourabh@isical.ac.in*

Abstract: In this paper we develop an inverse Bayesian approach to find the value of the unknown model parameter vector that supports the real (or test) data, where the data comprises measurements of a matrix-variate variable. The method is illustrated via the estimation of the unknown Milky Way feature parameter vector, using available test and simulated (training) stellar velocity data matrices. The data is represented as an unknown function of the model parameters, where this high-dimensional function is modelled using a high-dimensional Gaussian Process (\mathcal{GP}). The model for this function is trained using available training data and inverted by Bayesian means, to estimate the sought value of the model parameter vector at which the test data is realised. We achieve a closed-form expression for the posterior of the unknown parameter vector and the parameters of the invoked \mathcal{GP} , given test and training data. We perform model fitting by comparing the observed data with predictions made at different summaries of the posterior probability of the model parameter vector. As a supplement, we undertake a leave-one-out cross validation of our method.

* Associate Research fellow at Department of Statistics, University of Warwick and Lecturer of Statistics at Department of Mathematics, University of Leicester

[†] PhD student in Statistics and Mathematics Unit, Indian Statistical Institute

[‡] Assistant Professor in Bayesian and Interdisciplinary Research Unit, Indian Statistical Institute

Keywords and phrases: Supervised learning, inverse problems, Gaussian Process, matrix-variate normal, transformation-based MCMC.

Received June 2014.

1. Introduction

Curiosity about the nature of the parameter space of the Milky Way that we earthlings live in, is only natural. In this paper, we discuss the learning of the parameters characterising those Milky Way features that bear influence upon the motion of individual stars that lie in the neighbourhood of the Sun. Astrophysical modelling indicates that in the solar neighbourhood, effects of different features of the Milky Way are relevant (Minchev et al., 2009; Chakrabarty, 2007; Antoja et al., 2009). Such features include an elongated bar-like structure made of stars (the stellar bar) that rotates, pivoted at the centre of the Galaxy. In addition, the spiral arms of the Galaxy are also relevant. Thus, the motions of stars in the solar neighbourhood are affected by the parameters that define these Galactic features. Included in these feature parameters are the locations of the observer of such motions—we from Earth observe such motions, so that the stellar velocities are recorded to attain the observed values, given where in the Galaxy we are measuring these velocities from. On astronomical scales, the Earth’s location in the Milky Way is equivalent to the location of the Sun inside the Galaxy. Our location in the two-dimensional (by assumption) Galactic disk, is given by the angular separation of the Sun from a chosen line (an identified axis of the aforementioned stellar bar) and the distance from the Sun to the centre of the Galaxy. These two location parameters are the components of the two-dimensional location \mathbf{S} of the observer. As motivated above, parameters of the bar, spiral pattern and other Milky Way features, can also affect the motions of stars that are observed. (See section **S-1** of the supplementary material (Chakrabarty et al., 2015) for details). Given that these galactic feature parameters affect the solar neighbourhood, if motions of a sample of stars in this neighbourhood are measured, such data will harbour information about these feature parameters. Then, the inversion of such measured motions will in principle, allow for the learning of the unknown feature parameters. This approach has been adopted in the modelling of our galaxy, to result in the estimation of the angular separation of the Sun from a chosen axis of the bar, and the distance of the Sun from the Galactic centre (Minchev et al., 2010; Fux, 2001; Dehnen, 2000; Simone et al., 2004). The other relevant feature parameters are typically held constant in such modelling.

The above inverse problem is then an example application of the method of science that is typified by attempts at learning the unknown model parameter vector given observed data, where the causal relationship between the observable and the model parameter vector \mathbf{S} , is not necessarily known. This unknown relationship or function, can itself be learnt using available “training data”. Once this function is learnt, it can in principle be inverted to predict the unknown value of \mathbf{S} at which the measured data—i.e. “test data”—is realised. Such test

data is contrasted with “training data”, which is data generated at known or chosen values of \mathbf{S} (for example, via simulations or obtained as archival data).

The learning of a high-dimensional function from available training data, using standard nonparametric methods (such as spline fitting or wavelet based learning) is expected to be unsatisfactory since modelling high-dimensional functions using splines/wavelets may fail to adequately take into account the correlation structure between the component functions. Also, the complexity of the computational task of learning the unknown function from the data—and in particular of inverting it—only increases with dimensionality. Furthermore, the additional worry in the classical approach is that parameter uncertainty is ignored, though the same can be addressed in a Bayesian framework. An added advantage of the Bayesian approach is that priors on the unknown parameters can bring in extra information into the model, allowing for a training data set of comparatively smaller size (than that required in the classical approach), to be adequate.

Solving for the value of \mathbf{S} that supports the real or test data requires operating the inverse of the learnt function on the test data. The existence and uniqueness of such solution can be questioned given that the problem may not even be well-posed in a Hadamard sense (Kabanikhin, 2008; B.Hofmann, 2011; Tarantola, 2005). The problem may even be ill-conditioned since errors in the measurement may exist. Such worries about ill-posedness and ill-conditioning are mitigated in the Bayesian framework (Carreira-Perpiñán, 2001; Stuart, 2013). In this approach, the solution entails computation of the posterior probability of the unknown \mathbf{S} (at which the test data is realised), given all data. Given the inherent inadequacies of learning using splines/wavelets discussed above, we opt to model the unknown functional relationship between data and model parameter \mathbf{S} with a high-dimensional \mathcal{GP} . Similarly, in our application of interest, the unknown functional relation between the high-dimensional observations on stellar motions and the unknown observer location vector \mathbf{S} is modelled as a high-dimensional \mathcal{GP} . In this exercise, Galactic feature parameters other than the observer location are maintained as constants.

Chakrabarty (2007) constructed four different base-astronomical models of the solar neighbourhood, each at a chosen value of the ratio of the rate of rotation of the spiral pattern (Ω_s) to that of the bar (Ω_b). Non-linear dynamical evolution of each of these four base-astronomical models were carried out by Chakrabarty (2007), resulting in four independent data sets, each consisting of n blocks of j number of k -dimensional stellar velocity vectors, where each block is generated at a chosen value of \mathbf{S} (aka, a “design point”). At each possible chosen location \mathbf{s} of the Sun, the dynamical evolution of a given base-astronomical model of the Galaxy generates a block representing the k -dimensional velocity of each of j stars, where these stars are chosen as neighbours of the Sun. Thus, there are n design points and each training data set consists of n number of $j \times k$ -matrices, with a matrix generated at the corresponding design point. There are four such training data sets generated, by performing the evolution of each base-astronomical model. In addition, there is a measured, stellar velocity data matrix—of dimensionality $j \times k$ again—available, but this time, we do not know

what is the value of \mathbf{S} at which this measured/test data has been realised. It is this unknown value of \mathbf{S} that we seek to Bayesianly learn, given the test data and one training data set at a time.

It may be asked that if a stellar velocity matrix can be generated at a chosen \mathbf{s} , via the evolution of a base-astronomical model, does this not amount to stating that the causal relationship between the observable (velocity matrix) and model parameter (\mathbf{S}) is already known? Indeed this knowledge must be embedded within the evolutionary scheme implemented on any base-astronomical model. Thus, the forward evolution of a base-astronomical model is possible (via Newton's equations of motion), in order to generate a velocity matrix at a chosen \mathbf{s} . However the inversion of this evolution—aimed at recovering the sought \mathbf{s} at which the measured velocity matrix is generated—is not possible in general, owing to non-linear dynamical effects, or chaos, that impede reversibility in evolution; see Sengupta (2003), Section 6.6 of Chakrabarty (2007), Section 7 of Fux (2001). The strength of such chaos is different in the different base-astronomical models, caused by the different values of Ω_s/Ω_b , (discussed below in Section 3). This difficulty of inversion triggers the need to learn the inverse of the function that expresses the observable as a function of \mathbf{S} , independently from each of the four available training data sets. This learnt inverse function is then to be operated upon the measured (test) data to predict the value of \mathbf{S} in the Milky Way, in each of the four cases that represent four possible astronomical models of the Milky Way. We of course, predict this value of \mathbf{S} Bayesianly, by using a high-dimensional \mathcal{GP} to model the velocity data. We then achieve a closed-form posterior probability density of the sought \mathbf{s} and relevant parameters of this \mathcal{GP} , given the test and training data. Marginal posterior distribution of the components of the sought \mathbf{s} vector are inferred using MCMC, for each base-astronomical model (i.e. each training data set) used. Our focus in this work is to make inference on all values of \mathbf{S} at which the test data is realised, in each of the four astronomical models of the Galaxy—selection of the base-astronomical model is beyond the scope of this paper (see Section 4).

In the astronomical literature, Milky Way feature parameters in the solar neighbourhood have been explored via simulation based studies (Englmaier and Gerhard, 1999; Fux, 1997) while similar estimation is performed using other (astronomical) model-based studies (Aumer and Binney, 2009; Perryman, 2012; Golubov, 2012). Chakrabarty (2007) attempted estimation of the sought Galactic parameters via a test of hypothesis exercise: a non-parametric frequentist test was designed to test for the null that the observed stellar velocity data matrix is sampled from the estimated density of a synthetic velocity data matrix generated at the corresponding chosen value of the Milky Way feature parameter vector \mathbf{S} . The p -value of the used test statistic was recorded for each choice of \mathbf{s} . The choices of \mathbf{s} at which the highest p -values were obtained, were considered better supported by the observed data. Hence the empirical distribution of these p -values in the space of \mathbf{S} , was used to provide interval estimates of the Milky Way feature vector. However, this method required computational effort and is highly data intensive since the best match is sought over a very large collection of training data points. This shortcoming had compelled Chakrabarty (2007) to

resort to an unsatisfactory coarse gridding of the space of \mathbf{S} . This problem gets acute enough for the method to be rendered useless when the dimensionality of the vector \mathbf{S} that we hope to learn, increases. Moreover, the method of quantification of uncertainty of the estimate of the location is also unsatisfactory, dependent crucially on the binning details, which in turn is bounded by cost and memory considerations.

In the method we develop here, we demonstrate the effectiveness of our Gaussian Process based method with much smaller data sets than were used in the past. The other major advantage of this presented method is that it readily allows for the expansion of dimensionality of the model parameter vector and is capable of taking measurement errors into account.

The rest of the paper is structured as follows. In Section 2, we present the details of the modelling strategy that we adopt. The treatment of measurement errors within the modelling is discussed in Section 2.6. In Section 3 we discuss the application via which the new method is illustrated while details of the inference are discussed in Section 3.1. Section 4 contains results obtained from using available real and training data. We compare the obtained results with the estimates available in the astronomical literature in Section 4.1. Section 5 presents results of model fitting by comparing test data with predictions made at different summaries of the posterior of the model parameter vector \mathbf{S} . The paper is rounded up with Section 6.

2. Model

In this section we discuss the generic methodology that we use to learn the unknown location vector of the observer in the Milky Way disk, given the matrix-variate test and training stellar velocity data. Once the method is motivated, we implement it in the following section, to perform the learning relevant to the application at hand.

If a matrix-variate observable is expressed as an unknown matrix-variate function of the model parameter \mathbf{S} , and this unknown causal relationship between observable and \mathbf{S} is modelled by a matrix-variate Gaussian Process (\mathcal{GP}), it would imply that one realisation from such a matrix-variate \mathcal{GP} would be a set of the observed matrices that will be jointly distributed as 3-tensor normal, parametrised by a mean matrix and 3 covariance matrices (Hoff, 2011). While applications of the same are being developed (Wang & Chakrabarty), here we undertake an alternative and equivalent modelling strategy. We vectorise our intrinsically matrix-variate data sets to achieve a close-form expression for the joint posterior probability of the unknown parameters that we are interested in learning from the data. This leads to the functional relationship between the data and model parameter vector being rendered vector-valued, modelled by a vector-variate \mathcal{GP} , a set of realisations from which is jointly matrix normal, parametrised by matrix-variate parameters that we intend to learn from the data, along with the unknown \mathbf{s} at which the measured data is realised.

Let j number of measurements of a k -dimensional variable be available; this

vector variable is referred to below as the “observable”. Thus the measurements of this observable constitute a $j \times k$ -dimensional matrix. We refer to the measured data as test data and seek the unknown value $\mathbf{s}^{(new)}$ of model parameter \mathbf{S} at which it is realised. Let data be generated at n known values of \mathbf{S} : $\mathbf{s}_1^*, \dots, \mathbf{s}_n^*$. Then $\{\mathbf{s}_1^*, \dots, \mathbf{s}_n^*\}$ is the design set and \mathbf{s}_i^* is the i -th design vector at which the i -th synthetic data matrix is generated, $i = 1, 2, \dots, n$. Then these n synthetic data matrices comprise a training data set. Here a data matrix is $j \times k$ -dimensional. As motivated in the introductory section, we express the relation between the observable \mathbf{V} and unknown model parameter vector \mathbf{S} as $\mathbf{V} = \boldsymbol{\xi}(\mathbf{S})$, where $\boldsymbol{\xi}(\cdot)$ is an unknown function. We train the model for $\boldsymbol{\xi}(\cdot)$ using the training data and invert the function using Bayesian means to estimate the unknown $\mathbf{s}^{(new)}$ at which the test data is realised.

As discussed above, we vectorise the intrinsically $j \times k$ -dimensional matrix-variate data sets as jk -dimensional vectors. In this treatment, as a measurement is rendered vector-valued, $\boldsymbol{\xi}(\cdot)$ is vector-valued and $\boldsymbol{\xi}(\cdot)$ can be modelled by a vector-variate \mathcal{GP} so that realisations from this \mathcal{GP} are jointly matrix normal. Thus, we consider the j number of measurements of the k -dimensional observable, as a jk -dimensional observed vector $\mathbf{v}^{(test)}$. This test data is realised at the unknown value $\mathbf{s}^{(new)}$ of \mathbf{S} . Again, a $j \times k$ -dimensional synthetic data matrix is treated as a jk -dimensional synthetic data vector \mathbf{v}_i , $i = 1, 2, \dots, n$, along the lines of the observed data. Then all the n synthetic data vectors together

comprise the training data $\mathcal{D}_s = (\mathbf{v}_1 : \mathbf{v}_2 : \dots : \mathbf{v}_n)^T$ where \mathbf{v}_i is generated at the chosen value \mathbf{s}_i^* of \mathbf{S} , $i = 1, \dots, n$. Given our treatment of \mathbf{v}_i as a jk -dimensional vector, the training data set \mathcal{D}_s is a matrix with n rows and jk columns.

Thus in this treatment, we have n jk -dimensional synthetic data vectors (inputs), each generated at a chosen value of the model parameter vector (target), i.e. we have the n observations $(\mathbf{v}_1, \mathbf{s}_1^*), \dots, (\mathbf{v}_n, \mathbf{s}_n^*)$, and the aim is to predict the unknown model parameter vector $\mathbf{s}^{(new)}$ at which the input is the test data, i.e. the data vector $\mathbf{v}^{(test)}$. In this paradigm of supervised learning akin to the discussion in Neal (1998), a predictive distribution of $\mathbf{s}^{(new)}$ is sought, conditioned on the test data $\mathbf{v}^{(test)}$ and the training data $\mathcal{D}_s = (\mathbf{v}_1 : \mathbf{v}_2 : \dots : \mathbf{v}_n)^T$.

We begin the discussion on the model by elaborating on the detailed structure of the used \mathcal{GP} . In this section we ignore measurement errors and present our model of these n vector-variate functions. Later in Section 2.6, we delineate the method used to take measurement uncertainties on board.

As the data are vectorised as jk -dimensional vectors, $\boldsymbol{\xi}(\cdot)$ is also rendered a jk -variate vector function whose ℓ -th component function is $\xi_\ell(\cdot)$. Then we can write $\mathbf{v}_i = \boldsymbol{\xi}(\mathbf{s}_i) := (\xi_1(\mathbf{s}_i), \dots, \xi_{jk}(\mathbf{s}_i))^T$, $\forall i = 1, \dots, n$. We model the jk -dimensional function $\boldsymbol{\xi}(\cdot)$ with a jk -dimensional \mathcal{GP} , so that one realisation $\{\boldsymbol{\xi}(\mathbf{s}_1), \boldsymbol{\xi}(\mathbf{s}_2), \dots, \boldsymbol{\xi}(\mathbf{s}_n)\}$, from this \mathcal{GP} , is jointly matrix normal, with adequate parametrisation. We represent this as

$$\{\boldsymbol{\xi}(\mathbf{s}_1), \boldsymbol{\xi}(\mathbf{s}_2), \dots, \boldsymbol{\xi}(\mathbf{s}_n)\} \sim \mathcal{MN}_{n,jk}(\boldsymbol{\mu}, \mathbf{A}, \boldsymbol{\Omega}), \quad (2.1)$$

where the mean matrix of this matrix normal distribution is the $n \times jk$ -dimensional matrix $\boldsymbol{\mu}$, the left covariance matrix is the $n \times n$ -dimensional \mathbf{A} and the

right covariance matrix is the $jk \times jk$ -dimensional matrix $\mathbf{\Omega}$. These individual matrix-variate parameters of this distribution stem from the parametrisation of the high-dimensional \mathcal{GP} that is used to model $\boldsymbol{\xi}(\cdot)$; we discuss such parametrisation below. Before proceeding to that, we note that Equation 2.1 is the same as saying that the likelihood is matrix normal.

2.1. Parameters of the matrix-normal distribution

Assuming $\boldsymbol{\xi}(\cdot)$ to be continuous, the applicability of a stationary covariance function is expected to suffice. We choose to implement the popularly used square exponential covariance function (Rasmussen and Williams, 2006; Schölkopf and Smola, 2002; Santner et al., 2003). This covariance function is easy to implement and renders the sampled functions smooth and infinitely differentiable. Also, we relax the choice of a zero mean function though that is another popular choice. Instead we choose to define the mean function in a way that is equivalent to the suggestion that the data is viewed as centred around a linear model with the residuals characterised by a vector-variate \mathcal{GP} (A. O'Hagan, 1978; Cressie, 1993). We then integrate over all such possible global intercepts to arrive at a result that is more general than if the mean is fixed at zero. An advantage of the non-zero mean function is that in the limit of the smoothness parameters (characterising the smoothness of the functions sampled from this \mathcal{GP}) approaching large values, the random function reduces to a linear regression model. This appears plausible, as distinguished from the result that in this limit of very large smoothness, the random function will concur with the errors, as in models with a zero mean function.

The non-zero mean function $\boldsymbol{\mu}(\cdot)$ of the \mathcal{GP} is represented as factored into a matrix \mathbf{H} that bears information about its shape and another (\mathbf{B}) that tells us about its amplitude, or the extent to which this chosen mean function deviates from being zero. Thus, $\boldsymbol{\mu}(\cdot) := \mathbf{HB}$, where

$$\begin{aligned} \mathbf{H}^T &:= [\mathbf{h}^{(m \times 1)}(\mathbf{s}_1), \dots, \mathbf{h}^{(m \times 1)}(\mathbf{s}_n)], \quad \text{with} \\ m &:= d + 1 \\ \mathbf{h}^{(m \times 1)}(\mathbf{s}_i) &= (1, s_i^{(1)}, s_i^{(2)}, \dots, s_i^{(d)})^T \end{aligned} \quad (2.2)$$

where $\mathbf{s}_i = (s_i^{(1)}, s_i^{(2)}, \dots, s_i^{(d)})^T$ for $i = 1, \dots, n$ and we have recalled the suggestion that such a non-zero mean function be expressed in terms of a few basis functions (Rasmussen and Williams, 2006), (prompting us to choose to fix this functional form such that $\mathbf{h}(\mathbf{s}) := (1, \mathbf{s})^T$ for all values of \mathbf{S}). A similar construct was used by Blight and Ott (1975) who performed a \mathcal{GP} -based polynomial regression analysis. Thus, in our treatment, $\mathbf{h}(\cdot)$ is a $(d + 1)$ -dimensional vector. The coefficient matrix \mathbf{B} is

$$\mathbf{B} = (\boldsymbol{\beta}_{11}, \dots, \boldsymbol{\beta}_{j1}, \dots, \boldsymbol{\beta}_{1k}, \dots, \boldsymbol{\beta}_{jk}) \quad (2.3)$$

where for $p = 1, \dots, j$, $p' = 1, \dots, k$, $\boldsymbol{\beta}_{pp'}$ is an m -dimensional column vector. As we choose to set $m = d + 1$, \mathbf{B} is a matrix with $d + 1$ rows and jk columns.

The covariance function of the \mathcal{GP} is again represented as factored into a matrix $\mathbf{\Omega}$ that tells us about the amplitude of the covariance and another \mathbf{A} that bears information about its shape. The amplitude matrix $\mathbf{\Omega}$ is $jk \times jk$ -dimensional and is defined as

$$\mathbf{\Omega} = \mathbf{\Sigma} \otimes \mathbf{C} \quad (2.4)$$

where $\mathbf{\Sigma}$ is the $k \times k$ matrix telling us the amplitude of the covariance amongst the j different observations, for each of the k components of the data vector, at a fixed value of \mathbf{S} . On the other hand, \mathbf{C} is the $j \times j$ matrix giving the amplitude of covariance amongst the k different components of the vector-valued observable, at each of the j observations, at a given value of \mathbf{S} . Thus in our application, an element of $\mathbf{\Sigma}$ is the matrix is the amplitude of the covariance of a given component of the velocity vectors of the different stars that are observed. This matrix can then tell us about how a given component of the velocity vectors of the different stars in the observed sample, correlate with each other. On the other hand, the matrix \mathbf{C} informs us about the amplitude of covariance amongst the different components of the velocity vectors of a given star in the sample.

We realise that under the assumption of Gaussian errors in the measurements, the error variance matrix will be added to $\mathbf{\Omega}$. We discuss this in detail later in Section 2.6.

The shape of the covariance function is borne by the matrix \mathbf{A} which is $n \times n$ -dimensional. Given our choice of square exponential covariance function, it is defined as

$$\begin{aligned} \mathbf{A}^{(n \times n)} &:= [a(\cdot, \cdot)], \quad \text{where} \\ a(\mathbf{s}, \mathbf{s}') &\equiv \exp\{-(\mathbf{s} - \mathbf{s}')^T \mathbf{Q} (\mathbf{s} - \mathbf{s}')\}, \end{aligned} \quad (2.5)$$

for any 2 values \mathbf{s} and \mathbf{s}' of \mathbf{S} . Here, $\mathbf{Q}^{(d \times d)}$ represents the inverse of the scale length that underlies correlation between functions at any two values of the function variable. In other words, \mathbf{Q} is the matrix of the smoothness parameters. Thus, \mathbf{Q} is a matrix that bears information about the smoothness of the sampled functions; it is a diagonal matrix consisting of d non-negative smoothness parameters denoted by b_1, \dots, b_d . In other words, we assume the same smoothness for each component function of $\boldsymbol{\xi}(\cdot)$. This smoothness is determined by the parameters b_1, \dots, b_d . We will learn these smoothness parameters in our work from the data. Of course, though we say that the smoothness is learnt in the data, the underlying effect of the choice of the square exponential covariance function on the smoothness of the sampled functions is acknowledged. Indeed, as Snelson (2007) states, one concern about the square exponential function is that it renders the functions sampled from it as artificially smooth. An alternative covariance function, such as the Matern class of covariances (Matern, 1986; Gneiting et al., 2010; Snelson, 2007), could give rise to sampled functions that are much rougher than those obtained using the square exponential covariance function, for the same values of the hyper-parameters of amplitude and scale that characterise these covariance functions (see Chapter 1 of Snelson's thesis).

Let $\omega_{r\ell}$ denote the (r, ℓ) -th element of $\mathbf{\Omega}$, $c_{r\ell}$ the (r, ℓ) -th element of \mathbf{C} and let $\sigma_{r\ell}$ denote the (r, ℓ) -th element of $\mathbf{\Sigma}$. Let the ℓ -th component function of $\boldsymbol{\xi}(\cdot)$ be $\xi_\ell(\cdot)$ with $\ell = m_1 k + m_2$, where $\ell = 1, \dots, jk$ and $m_2 = 1, 2, \dots, k$, $m_1 = 0, 1, \dots, j-1$. Then the correlation between the components of $\boldsymbol{\xi}(\cdot)$ yields the following correlation structures:

$$\text{corr}(\xi_{m_1 k + m_2}(\mathbf{s}_i), \xi_{m'_1 k + m_2}(\mathbf{s}_i)) = \frac{\sigma_{m_1 m'_1}}{\sqrt{\sigma_{m_1 m_1} \sigma_{m'_1 m'_1}}} \quad \forall m_2, i \text{ and } m_1 \neq m'_1 \quad (2.6)$$

$$\text{corr}(\xi_{m_1 k + m_2}(\mathbf{s}_i), \xi_{m_1 k + m'_2}(\mathbf{s}_i)) = \frac{c_{m_2 m'_2}}{\sqrt{c_{m_2 m_2} c_{m'_2 m'_2}}} \quad \forall m_1, i \text{ and } m_2 \neq m'_2 \quad (2.7)$$

$$\text{corr}(\xi_{m_1 k + m_2}(\mathbf{s}_i), \xi_{m'_1 k + m'_2}(\mathbf{s}_i)) = \frac{c_{m_2 m'_2} \sigma_{m_1 m'_1}}{\sqrt{c_{m_2 m_2} \sigma_{m_1 m_1} c_{m'_2 m'_2} \sigma_{m'_1 m'_1}}} \quad \forall i, m_1 \neq m'_1, m_2 \neq m'_2 \quad (2.8)$$

$$\text{corr}(\xi_\ell(\mathbf{s}_1), \xi_\ell(\mathbf{s}_2)) = a(\mathbf{s}_1, \mathbf{s}_2) \quad \forall \ell \text{ and } \mathbf{s}_1 \neq \mathbf{s}_2 \quad (2.9)$$

The 1st of the above 4 equations shows the correlation between the component functions for the same component of the vector-valued observable at 2 (of the j) different measurements, taken at a given value of the \mathbf{S} . For a given measurement, the correlation between 2 different components of (the k components of) the observable is given by the 2nd equation above. For a given value of \mathbf{S} , if we seek the correlation between the component functions for 2 different measurements of 2 different components of the observables, this is provided in the 3rd equation. The correlation between component functions for 2 different values of \mathbf{S} is given in the last of the above 4 equation. Then these 4 correlations give the full correlation structure amongst components of $\boldsymbol{\xi}(\cdot)$.

2.2. Likelihood

The training data is the $n \times jk$ -dimensional matrix $\mathcal{D}_s = (\mathbf{v}_1 \vdots \mathbf{v}_2 \vdots \dots \vdots \mathbf{v}_n)^T$ where \mathbf{v}_i is the jk -dimensional synthetic motion vector generated at design vector \mathbf{s}_i^* , $i = 1, 2, \dots, n$. To express the likelihood, we recall that the distribution of the training data $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n\}$, i.e. the joint distribution of $\{\boldsymbol{\xi}(\mathbf{s}_1^*), \boldsymbol{\xi}(\mathbf{s}_2^*), \dots, \boldsymbol{\xi}(\mathbf{s}_n^*)\}$ is matrix normal (Equation 2.1). In order to achieve this likelihood, we rewrite the \mathbf{S} -dependent parameters of this matrix normal distribution at the values of \mathbf{S} at which the training data \mathcal{D}_s is realised, i.e. in terms of the design vectors. Thus, we define

- the $n \times jk$ -dimensional mean function $\mathbf{H}_D \mathbf{B}$, where the linear form of the mean structure is contained in $\mathbf{H}_D^{(n \times m)} := [\mathbf{h}^{(m \times 1)}(\mathbf{s}_1^*), \dots, \mathbf{h}^{(m \times 1)}(\mathbf{s}_n^*)]$ (and the coefficient matrix \mathbf{B} is defined in Equation 2.3).
- the square exponential factor in the covariance matrix $\mathbf{A}_D^{(n \times n)} := [\exp\{-(\mathbf{s}^* - \mathbf{s}'^*)^T \mathbf{Q}(\mathbf{s}^* - \mathbf{s}'^*)\}]$ (see Equation 2.5).

Then it follows from the matrix normal distribution of Equation 2.1—with mean function defined in Equation 2.2 and Equation 2.3, and covariance matrix defined using Equation 2.5 and Equation 2.4—that \mathcal{D}_s is distributed as matrix normal with mean matrix $\mathbf{H}_D \mathbf{B}$, left covariance matrix \mathbf{A}_D and right covariance matrix $\mathbf{\Omega}$, i.e.

$$[\mathcal{D}_s \mid \mathbf{B}, \mathbf{C}, \mathbf{\Sigma}, \mathbf{Q}] \sim \mathcal{MN}_{n,jk}(\mathbf{H}_D \mathbf{B}, \mathbf{A}_D, \mathbf{\Omega}) \quad (2.10)$$

Thus, using known ideas about the matrix normal distribution - see Dawid (1981), Carvalho and West (2007) - we write

$$\begin{aligned} [\mathcal{D}_s \mid \mathbf{B}, \mathbf{C}, \mathbf{\Sigma}, \mathbf{Q}] = & \frac{1}{(2\pi)^{\frac{njk}{2}} |\mathbf{A}_D|^{\frac{jk}{2}} |\mathbf{\Omega}|^{\frac{n}{2}}} \\ & \times \exp \left\{ -\frac{1}{2} \text{tr} [\mathbf{\Omega}^{-1} (\mathcal{D}_s - \mathbf{H}_D \mathbf{B})^T \mathbf{A}_D^{-1} (\mathcal{D}_s - \mathbf{H}_D \mathbf{B})] \right\} \end{aligned} \quad (2.11)$$

The interpretation of the above is that the r -th row of $[\mathcal{D}_s \mid \mathbf{B}, \mathbf{\Sigma}, \mathbf{C}, \mathbf{Q}]$ is multivariate normal with mean corresponding to row of the mean matrix $\mathbf{H}_D \mathbf{B}$ and with covariance matrix $\mathbf{\Omega}$. Rows r and ℓ of $[\mathcal{D}_s \mid \mathbf{B}, \mathbf{\Sigma}, \mathbf{C}, \mathbf{Q}]$ has covariance matrix $a(\mathbf{s}_r, \mathbf{s}_\ell) \mathbf{\Omega}$. Similarly, the ℓ -th column of it is distributed as multivariate normal with mean being the ℓ -th column of $\mathbf{H}_D \mathbf{B}$ and with covariance matrix $\omega_{\ell,\ell} \mathbf{A}_D$, where $\omega_{r,\ell}$ denotes the (r, ℓ) -th element of $\mathbf{\Omega}$. The covariance between columns r and ℓ is given by the matrix $\omega_{r,\ell} \mathbf{A}_D$.

2.3. Estimating $\mathbf{s}^{(new)}$

In order to predict the unknown model parameter vector $\mathbf{s}^{(new)}$ when the input is the measured real data vector $\mathbf{v}^{(test)}$, we would need the posterior predictive distribution of $\mathbf{s}^{(new)}$, given $\mathbf{v}^{(test)}$ and the training data \mathcal{D}_s . This posterior predictive is usually computed by integrating over all the matrix-variate \mathcal{GP} parameters realised at the chosen design vectors $\mathbf{s}_1^*, \dots, \mathbf{s}_n^*$.

While it is possible to analytically integrate over \mathbf{B} and \mathbf{C} , $\mathbf{\Sigma}$ and \mathbf{Q} cannot be analytically integrated out. In fact, we find it useful to learn the d smoothing parameters i.e. the d diagonal elements of \mathbf{Q} , given the data. Thus, one useful advantage of our method is that the smoothness of the process does not need to be imposed by hand, but can be learnt from the data, if desired.

Given that we are then learning of $\mathbf{s}^{(new)}$, $\mathbf{\Sigma}$ and \mathbf{Q} , we rephrase our motivation as seeking to compute the joint posterior probability of $\mathbf{s}^{(new)}$, \mathbf{Q} and $\mathbf{\Sigma}$, conditional on the real data and the training data, for a choice of the design set. In fact, we achieve a closed form expression of this joint posterior of $\mathbf{s}^{(new)}$, \mathbf{Q} and $\mathbf{\Sigma}$, by integrating over the other hyper-parameters, namely, the amplitude of the mean function (\mathbf{B}) and the matrix \mathbf{C} that bears information about covariance between different components of the data vector for each of the j observations, at a fixed value of \mathbf{S} . From this closed form expression, the marginal posterior probability densities of \mathbf{Q} , $\mathbf{\Sigma}$ and any of the d components

of the $\mathbf{s}^{(new)}$ vector can be obtained, using the transformation based MCMC sampling method (Dutta and Bhattacharya, 2013) that we adopt.

Thus, for a given choice $\mathbf{s}_1^*, \dots, \mathbf{s}_n^*$ of the design vectors, the posterior distribution $[\mathbf{s}^{(new)}, \Sigma, \mathbf{Q} | \mathbf{v}^{(test)}, \mathcal{D}_s]$ is sought, by marginalising $[\mathbf{s}^{(test)}, \Sigma, \mathbf{Q}, \mathbf{B}, \mathbf{C} | \mathbf{v}^{(test)}, \mathcal{D}_s]$ over the process matrices \mathbf{B} and \mathbf{C} .

2.4. Priors used

We use uniform prior on \mathbf{B} and a simple non-informative prior on \mathbf{C} , namely, $\pi(\mathbf{C}) \propto |\mathbf{C}|^{-(j+1)/2}$. As for the priors on the other parameters, we assume uniform prior on \mathbf{Q} and use the non-informative prior $\pi(\Sigma) \propto |\Sigma|^{-(k+1)/2}$. The prior information available in the literature will be considered to select the prior on $\mathbf{s}^{(new)}$; below we use uniform priors on all components of the $\mathbf{s}^{(new)}$ vector (see Section 4 for greater details in regard to the application that we discuss later).

2.5. Posterior of $\mathbf{s}^{(new)}$ given training and test data

Since our interest lies in estimating $\mathbf{s}^{(new)}$, given the real (test) data and the simulated (training) data, as well as in learning the smoothness parameter matrix \mathbf{Q} and the matrix Σ that bears the covariance amongst the j observables, we compute the joint posterior probability density $[\mathbf{s}^{(new)}, \mathbf{Q}, \Sigma | \mathbf{v}^{(test)}, \mathcal{D}_s]$. As expressed above, we achieve this by writing $[\mathbf{s}^{(new)}, \mathbf{B}, \mathbf{C}, \mathbf{Q}, \Sigma | \mathbf{v}^{(test)}, \mathcal{D}_s]$ and marginalise over \mathbf{B} and \mathbf{C} .

To construct an expression for this posterior distribution, we first collate the training and test data to construct the augmented data set $\mathcal{D}_{aug}^T = (\mathbf{v}_1^T; \dots; \mathbf{v}_n^T; (\mathbf{v}^{(test)})^T)$. Then the set of values of the model parameter vector \mathbf{S} that supports \mathcal{D}_{aug} is $\{\mathbf{s}_1^*, \dots, \mathbf{s}_n^*, \mathbf{s}^{(new)}\}$ of which only $\mathbf{s}^{(new)}$ is unknown.

We next write the \mathbf{S} -dependent matrix-variate parameters at those values of \mathbf{S} at which the augmented data set is realised. Thus we define

- $\mathbf{H}_{\mathcal{D}_{aug}}^{((n+1) \times m)} := [\mathbf{h}^{(m \times 1)}(\mathbf{s}_1^*), \dots, \mathbf{h}^{(m \times 1)}(\mathbf{s}_n^*), \mathbf{h}^{(m \times 1)}(\mathbf{s}^{(new)})]$, where our choice of the functional form of $\mathbf{h}(\cdot)$ has been given in Section 2 and we also set $m = d + 1$,
- $\mathbf{A}_{\mathcal{D}_{aug}}^{((n+1) \times (n+1))} := [\exp\{-(\mathbf{s}'_i - \mathbf{s}'_{i'})^T \mathbf{Q} (\mathbf{s}'_i - \mathbf{s}'_{i'})\}]$ where \mathbf{s}'_i and $\mathbf{s}'_{i'}$ are members of the set $\{\mathbf{s}_1^*, \dots, \mathbf{s}_n^*, \mathbf{s}^{(new)}\}$,
- $\mathbf{M}_{aug} := \mathbf{A}_{\mathcal{D}_{aug}}^{-1} - \mathbf{A}_{\mathcal{D}_{aug}}^{-1} \mathbf{H}_{\mathcal{D}_{aug}} [\mathbf{H}_{\mathcal{D}_{aug}}^T \mathbf{A}_{\mathcal{D}_{aug}}^{-1} \mathbf{H}_{\mathcal{D}_{aug}}]^{-1} \mathbf{H}_{\mathcal{D}_{aug}}^T \mathbf{A}_{\mathcal{D}_{aug}}^{-1}$.
- $(\mathcal{D}_{aug}^T \mathbf{M}_{aug} \mathcal{D}_{aug})^{(jk \times jk)} := [\mathbf{M}_{tu}^*; t, u = 1, \dots, k]$, where \mathbf{M}_{tu}^* is a matrix with j rows and j columns. Given Σ , we define $m = d + 1$ and ψ_{tu}^{-1} as the (t, u) -th element of Σ^{-1} , so that $(n + 1 - m)k \hat{\mathbf{C}}_{GLS, aug} := \sum_{t=1}^k \sum_{u=1}^k \psi_{tu}^{-1} \mathbf{M}_{tu}^*$, where $(n + 1 - m)k \hat{\mathbf{C}}_{GLS, aug}$ is used in the closed-form expression for $[\mathbf{s}^{(new)}, \mathbf{Q}, \Sigma | \mathbf{v}^{(test)}, \mathcal{D}_s]$ that we seek.

The priors used on \mathbf{B} , \mathbf{C} , \mathbf{Q} , Σ and $\mathbf{s}^{(new)}$ are listed in Section 2.4. Using these, and recalling Equation 2.11, we get the joint posterior probability density

of all unknown parameters given all data, i.e.

$$[\mathbf{s}^{(new)}, \mathbf{Q}, \mathbf{B}, \mathbf{\Sigma}, \mathbf{C} \mid \mathbf{v}^{(test)}, \mathcal{D}_s] \\ \propto [\mathcal{D}_{aug} \mid \mathbf{B}, \mathbf{\Sigma}, \mathbf{C}, \mathbf{Q}, \mathbf{s}^{(new)}][\mathbf{B}, \mathbf{\Sigma}, \mathbf{C}, \mathbf{Q}, \mathbf{s}^{(new)}],$$

which we then marginalise over \mathbf{B} and \mathbf{C} to get the joint posterior $[\mathbf{s}^{(new)}, \mathbf{Q}, \mathbf{\Sigma} \mid \mathbf{v}^{(test)}, \mathcal{D}_s]$, as

$$[\mathbf{s}^{(new)}, \mathbf{Q}, \mathbf{\Sigma} \mid \mathbf{v}^{(test)}, \mathcal{D}_s] \\ = \int \int [\mathbf{s}^{(new)}, \mathbf{Q}, \mathbf{B}, \mathbf{\Sigma}, \mathbf{C} \mid \mathbf{v}^{(test)}, \mathcal{D}_s] d\mathbf{B} d\mathbf{C} \\ \propto |\mathbf{A}_{\mathcal{D}_{aug}}|^{-\frac{jk}{2}} \{ \mathbf{H}_{\mathcal{D}_{aug}} \}^T \{ \mathbf{A}_{\mathcal{D}_{aug}} \}^{-1} \{ \mathbf{H}_{\mathcal{D}_{aug}} \}^{-\frac{jk}{2}} \times |\mathbf{\Sigma}|^{-\frac{j(n+1-m)+k+1}{2}} \\ \times |(n+1-m)k\hat{\mathbf{C}}_{GLS,aug}|^{-\frac{(n+1-m)k}{2}} \quad (2.12)$$

Thus, we obtain a closed-form expression of the joint posterior of $\mathbf{s}^{(new)}, \mathbf{Q}, \mathbf{\Sigma}$, given training and test data, for a given choice of the design matrix (Equation 2.12), up to a normalising constant. The \mathcal{GP} prior is strengthened by the n number of samples taken from it at the training stage. We sample from the achieved posterior using MCMC techniques to achieve the marginal posterior probabilities of \mathbf{Q} , $\mathbf{\Sigma}$ or any component of $\mathbf{s}^{(new)}$, given all data. We conduct posterior inference using the TMCMC methodology (Dutta and Bhattacharya, 2013) that works by constructing proposals that are deterministic bijective transformations of a random vector drawn from a chosen distribution.

2.6. Errors in measurement

In our application, the errors in the measurements are small and will be ignored for the rest of the analysis. In general, when errors in the measurements that comprise the training data and the test data are not negligible, we assume Gaussian measurement errors ε_t , in \mathbf{v}_t , with $t = 1, 2, \dots$, such that $\varepsilon_t \sim \mathcal{N}_{jk}(\mathbf{0}, \varsigma)$, where $\varsigma = \mathbf{\Sigma}_1 \otimes \mathbf{\Sigma}_2$; $\mathbf{\Sigma}_1, \mathbf{\Sigma}_2$ being positive definite matrices. If both $\mathbf{\Sigma}_1$ and $\mathbf{\Sigma}_2$ are chosen to be diagonal matrices, then ς is a diagonal matrix; assuming same diagonal elements would simplify ς to be of the form $\varphi \times \mathbf{I}$, where \mathbf{I} is the $jk \times jk$ -th order identity matrix. This error variance matrix ς must be added to $\mathbf{\Omega}$ before proceeding to the subsequent calculations. TMCMC can be then be used to update ς .

3. Case study

Using the methodology discussed above we attempt an estimate of the unknown Milky Way feature parameter vector $\mathbf{S} \in \mathbb{R}^d$ using the available stellar velocity data. In our application, the dimensionality of \mathbf{S} is 2 as we estimate the coordinates of the radial location r_\odot of the Sun with respect to the Galactic centre and the angular separation ϕ_\odot of the Sun-Galactic centre line from a pre-set

line in the Milky Way disk (see Figure 1 in supplementary section **S-1**). Then for the Sun, $R = r_\odot$ and $\Phi = \phi_\odot$ where the variable R gives radial distance from the Galactic centre of any point on the disk of the Milky Way and the variable Φ gives the angular separation of this point from this chosen pre-set line. The reason for restricting our application to the case of $d=2$ is the existence of simulated stellar velocity data (aka training data) generated by scanning over chosen guesses for r_\odot and ϕ_\odot , with all other feature parameters held constant. If simulated data distinguished by choices of other Milky Way feature parameters become available, then the implementation of such data as training data will be possible, allowing then for the learning of Milky way parameters in addition to r_\odot and ϕ_\odot . In this method, computational costs are the only concern in extending to cases of $d > 2$; extending to a higher dimensional \mathbf{S} only linearly scales computational costs (Section 6).

Also, the stellar velocity vector is 2-dimensional, i.e. $k=2$ in this application. Then the measured data in this application is a $j \times 2$ -dimensional matrix. In our Bayesian approach, a much smaller j ($=50$) allows for inference on the unknown value $\mathbf{s}^{(new)}$ of the Milky Way feature parameter vector, than $j \sim 3000$ that is demanded by the aforementioned calibration approach used by Chakrabarty (2007).

In our application, the available data include the measured or test data and 4 sets of synthetic (or training) data sets obtained via dynamical simulations of each of 4 distinct base-astronomical models of our galaxy, advanced by Chakrabarty (2007). As the analysis is performed with each training data set at a time, we do not include reference to the corresponding base model in the used notation. The simulated data presented in Chakrabarty (2007) that we use here, is generated at 216 distinct values of \mathbf{S} , i.e. $n=216$. Thus, our design set comprises the 216 chosen values of \mathbf{S} : $\mathbf{s}_1^*, \dots, \mathbf{s}_{216}^*$. For each of the 4 base astrophysical models, at each chosen \mathbf{s}_i^* , 50 2-dimensional stellar velocity vectors are generated from dynamical simulations of that astrophysical model (of the Milky Way), performed at that value of \mathbf{S} . These 50 2-dimensional velocity vectors are treated in our work as a $50 \times 2 = 100$ -dimensional motion vector \mathbf{v}_i ; $i = 1, \dots, 216$. Then at the 216 design vectors, $\mathbf{s}_1^*, \dots, \mathbf{s}_{216}^*$, 216 motion vectors are generated: $\mathbf{v}_1, \dots, \mathbf{v}_{216}$. Then the training data in our work comprises all such motion vectors and is represented as $\mathcal{D}_s^{(216 \times 100)}$. The real or test data is treated in our work as the 100-dimensional motion vector $\mathbf{v}^{(test)}$.

As said above, there are 4 distinct training data sets available from using the 4 base astronomical models of the Milky Way, as considered by Chakrabarty (2007). The choice of the base astrophysical model is distinguished by the ratio of the rates of rotation of the spiral to the bar, Ω_s/Ω_b . That this ratio is relevant to stellar motions in the Galaxy is due to the fact that Ω_s/Ω_b can crucially control the degree of chaos in the Galactic model¹. Thus, the 4 base

¹For example, it is well known in chaos theory that when Ω_s/Ω_b is such that one of the radii at which the bar and the stellar disk resonate, concurs with a radius at which the spiral and the stellar disk resonate, global chaos is set up in the system (G. Walker and J. Ford, 1969). Chakrabarty and Sideris (2008) have corroborated that the degree of chaos is maximal in the astrophysical Galactic model marked by such a ratio ($\Omega_s/\Omega_b=22/55$). They report that

astrophysical models are differently chaotic. This results in 4 distinct simulated velocity data sets $\mathcal{D}_s^{(1)}, \mathcal{D}_s^{(2)}, \mathcal{D}_s^{(3)}, \mathcal{D}_s^{(4)}$ that bear the effects of such varying degrees of chaos, each generated at the chosen design set $\{\mathbf{s}_1^*, \dots, \mathbf{s}_n^*\}$. Details of the dynamical simulations performed on the 4 astrophysical models are given in the supplementary section **S-2**.

3.1. Details of our implementation of TMCMC

As indicated above, we use the Transformation-based MCMC (TMCMC) advanced by Dutta and Bhattacharya (2013) to conduct posterior inference. In TMCMC, high-dimensional parameter spaces are explored by constructing bijective deterministic transformations of a low-dimensional random vector. The random vector of which a proposal density is a transformation of, can be chosen to be of dimensionality between 1 and the dimensionality of the parameters under the target posterior. The acceptance ratio in TMCMC does not depend upon the distribution of the chosen random vector. In our application we use TMCMC to update the entire block $(\mathbf{s}^{(new)}, \mathbf{Q}, \mathbf{\Sigma})$ at the same time using additive transformations of a one-dimensional random variable $\epsilon \sim \mathcal{N}(0, 1)I_{\{\epsilon > 0\}}$. In the t -th iteration, the state of the unknown parameters is $(\mathbf{s}^{(new, t)}, \mathbf{Q}^{(t)}, \mathbf{\Sigma}^{(t)}) := \boldsymbol{\varphi}^{(t)}$. We update $\boldsymbol{\varphi}^{(t)}$ by setting, with probabilities π_j and $(1 - \pi_j)$, $\varphi_j^{(t+1)} = \varphi_j^{(t)} \pm c_j \epsilon$ (forward transformation) and $\varphi_j^{(t+1)} = \varphi_j^{(t)} - c_j \epsilon$ (backward transformation), respectively, where, for $j = 1, \dots, d$, π_j are appropriately chosen probabilities and c_j are appropriately chosen scaling factors. Assume that for $j_1 \in \mathcal{U}$, $\varphi_{j_1}^{(t)}$ gets the positive transformation, while for $j_2 \in \mathcal{U}^c$, $\varphi_{j_2}^{(t)}$ gets the backward transformation. Here $\mathcal{U} \cup \mathcal{U}^c = \{1, \dots, d^*\}$, where $d^* = 2d + \frac{k(k+1)}{2}$. The proposal $\boldsymbol{\varphi}^{(t+1)}$ is accepted with acceptance probability given in Supplementary Section **S-3**. Once the proposal mechanism and the initial values are decided, we discard the first 100,000 iterations of our final TMCMC run as burn-in and stored the next 1,000,000 iterations for inference. For each model it took approximately 6 hours on a laptop to generate 1,100,000 TMCMC iterations.

4. Results using real data

The training data that we use was obtained by Chakrabarty (2007), by choosing the solar radial location from the interval $[1.7, 2.3]$ in model units. This explains the motivation for selecting the bounds on r_\odot to be the edges of this interval. Here, values of distances are expressed in the units implemented in the base astrophysical models of the Milky Way. However, to make sense of the results we have obtained, these model units will need to be scaled to provide values

in models marked by slightly lower ($\Omega_s/\Omega_b=18/55$) or higher ($\Omega_s/\Omega_b = 25/55$) values of this ratio, chaos is still substantial. In the Galactic model that precludes the spiral however, chaos was quantified to be minimal. It is these 4 states of chaos - driven by the 4 values of Ω_s/Ω_b - that mark the 4 astrophysical models as distinct.

in real astronomical units of distances inside galaxies, such as the “kiloparsec” (abbreviated as “kpc”). A distance of 1 in model unit scales to $\frac{\mathcal{R}}{\hat{r}_\odot}$ kpc, where \mathcal{R} is the solar radius obtained in independent astronomical studies (Binney and Merrifield, 1998, $\mathcal{R}=8\text{kpc}$) and \hat{r}_\odot is the estimate of the solar radius in our work. The ulterior aim in estimating the solar radius is in estimating the rotational frequency Ω_b of the bar, where $\Omega_b = \frac{v_0}{\hat{r}_\odot}$, with $v_0=220\text{kms}^{-1}$ and $\mathcal{R}=8\text{kpc}$. Then, we get $\Omega_b = \frac{220}{\hat{r}_\odot} \text{kms}^{-1}/\text{kpc}$. See Section S-1 of the attached supplementary material to see a schematic representation of the central bar in the Galaxy and Section S-2 for details of the scaling between the model units and real astronomical units.

Our other estimate is of the angular separation between the long axis of the bar and the line that joins the Sun to the Galactic centre. It is suggested in past astronomical modelling work to be an acute angle (Chakrabarty, 2007; Englmaier and Gerhard, 1999; Fux, 2001). Indeed, the training data used here was generated in simulations performed by Chakrabarty (2007), in which ϕ_\odot is chosen from the interval $[0, 90^\circ]$. This motivates the consideration of the interval of $[0, 90^\circ]$ for the angular location of the Sun.

Given the bounds on r_\odot and ϕ_\odot presented above, in our TMCMC algorithm, we reject those moves that suggest r_\odot and ϕ_\odot values that fall outside these presented intervals.

The 4 astrophysical models of the Galaxy that were used to generate the 4 training data sets, are marked by the same choice of the value of Ω_b and the background Galactic model parameters, while they are distinguished by the varying choices of the ratio $\Omega_s : \Omega_b$, where the Galactic spiral pattern rotates with rate Ω_s . In fact, the astrophysical model *bar_6* is the only one that does not include the influence of the spiral pattern while the other three astrophysical models include the influence of both the bar and the spiral. For the astrophysical models *sp3bar3_18*, *sp3bar3* and *sp3bar3_25*, $\Omega_s : \Omega_b$ is respectively set to $18\Omega_b/55$, $22\Omega_b/55$, $25\Omega_b/55$. The physical effect of this choice is to induce varying levels of chaoticity in the 4 astrophysical models. Thus, Chakrabarty and Sideris (2008) confirmed that of the 4 models, *bar_6* manifests very low chaoticity while *sp3bar3* manifests maximal chaos, though both *sp3bar3_18*, *sp3bar3_25* are comparably chaotic.

Ancillary real data needs to be brought in to judge the relative fit amongst the astrophysical base models. In fact, Chakrabarty (2007) brought in extra information to perform model selection. Such information was about the observed variance of the components of stellar velocities and this was used to rule out the model *bar_6* as physically viable, though the other three models were all acceptable from the point of view of such ancillary observations that are available. This led to the inference that $\Omega_s \in [18\Omega_b/55, 25\Omega_b/55]$.

It is to be noted that if there was 1 data set and we were trying to fit 4 different models to that same data, then it is very much possible that for this 1 data set, the average of 4 models could have been achieved. However, here we are dealing with 4 base models, each of which is giving rise to a distinct training data set, in

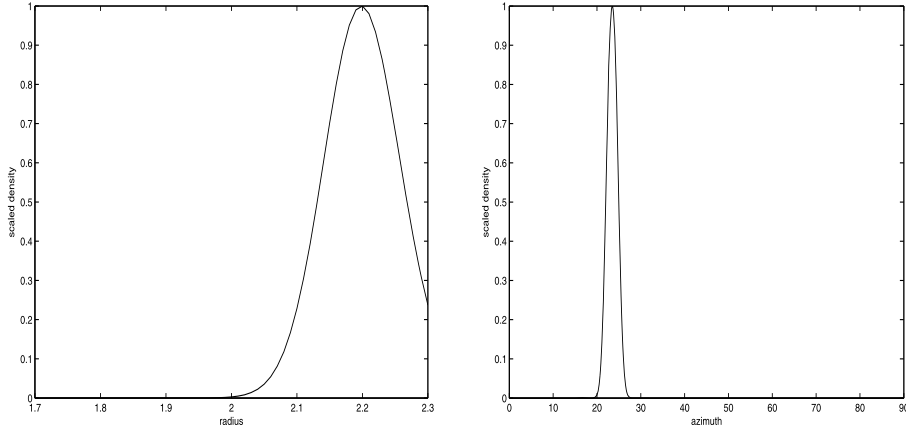


FIG 1. Posteriors of r_{\odot} in model units of r_{CR} and ϕ_{\odot} (in degrees) for the model *bar_6*.

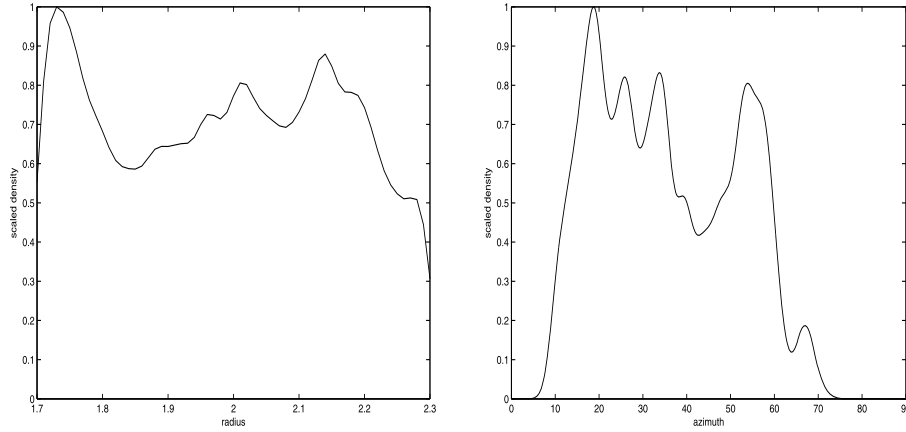


FIG 2. Posteriors of r_{\odot} in units of r_{CR} and ϕ_{\odot} (in degrees) for the model *sp3bar3*.

fact under mutually contradicting physics. Therefore, such model averaging is not relevant for this work. Cross-validation of these 4 models is indeed possible and we present this in Section **S-5** of the attached Supplementary Materials.

The marginal posterior densities of $(r_{\odot}, \phi_{\odot})$ corresponding to the 4 base astrophysical models of the Milky Way, are shown in Figures 1, 2, 3 and 4. It merits mention that the multi-modality manifest in the marginal posterior distributions in 3 of the 4 base models is not an artifact of inadequate convergence but is a direct fallout of the marked amount of chaoticity in all 3 base models except in the model *bar_6*, (Chakrabarty and Sideris, 2008). In Section **S-6**, we discuss the connection between chaos and consistency of multiple observer locations with available stellar velocity data.

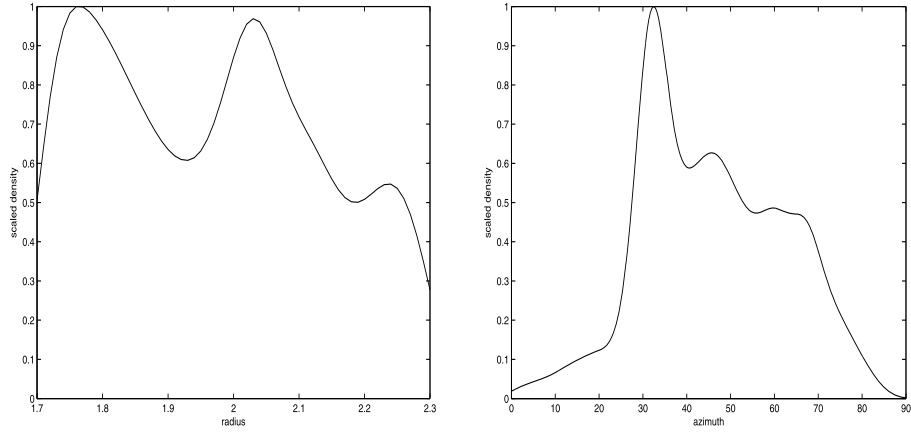


FIG 3. Posteriors of r_{\odot} in model units of r_{CR} and ϕ_{\odot} (in degrees) for the model *sp3bar3_18*.

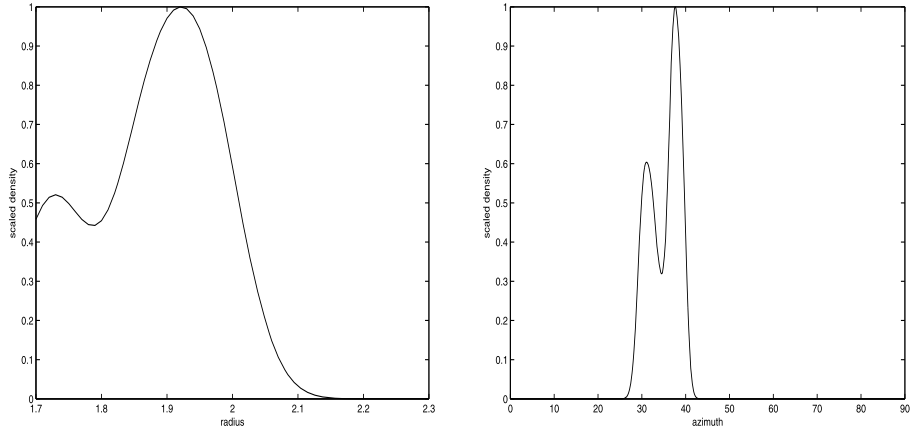


FIG 4. Posteriors of r_{\odot} in units of r_{CR} and ϕ_{\odot} (in degrees) for the model *sp3bar3_25*.

Table 1 presents the posterior mode, the 95% highest posterior density (HPD) credible region of r_{\odot} and ϕ_{\odot} respectively, associated with the four base models. Here r_{\odot} is expressed in the model units of length, i.e. in units of r_{CR} . ϕ_{\odot} is expressed in degrees. The HPDs are computed using the methodology discussed in Carlin and Louis (1996). Disjoint HPD regions, characterise the highly multimodal posterior distributions of the unknown location. Using the 95% HPDs of the estimate \hat{r}_{\odot} expressed in model units, and using the independently known astronomical measurement of the solar radial location as 8kpc, the bar rotational frequency Ω_b is computed (see third enumerated point discussed above) in Table 1.

Summaries of the posteriors (mean, variance and 95% credible interval) of the smoothness parameters b_1, b_2 and Σ are presented in Tables 2, 3. Notable

TABLE 1

Summary of the posterior distributions of the radial component r_\odot and azimuthal component ϕ_\odot of the unknown observer location vector for the 4 base astrophysical models and the unknown bar rotational frequency Ω_b computed using the 95% HPDs on the learnt radial location r_\odot in these models

| Model | r_\odot (in units of r_{CR}) | | Ω_b (in $\text{kms}^{-1}/\text{kpc}$) | | ϕ_\odot | |
|-------------------|-----------------------------------|----------------------------------|---|--|--------------|----------------|
| | Mode | 95% HPD | 95% HPD | | Mode | 95% HPD |
| <i>bar6</i> | 2.20 | [2.04, 2.30] | [56.1, 63.25] | | 23.50 | [21.20, 25.80] |
| <i>sp3bar3</i> | 1.73 | [1.70, 2.26] \cup [2.27, 2.28] | [46.75, 62.15] \cup [62.45, 62.7] | | 18.8 | [9.6, 61.5] |
| <i>sp3bar3_18</i> | 1.76 | [1.70, 2.29] | [46.75, 62.98] | | 32.5 | [17.60, 79.90] |
| <i>sp3bar3_25</i> | 1.95 | [1.70, 2.15] | [46.75, 59.12] | | 37.6 | [28.80, 40.40] |

TABLE 2

Summary of the posterior distributions of the smoothness parameters b_1, b_2 for the 4 models

| Model | b_1 | | | b_2 | | |
|-------------------|-----------|------------------------|----------------------|----------|-----------------------|----------------------|
| | Mean | Var | 95% CI | Mean | Var | 95% CI |
| <i>bar_6</i> | 0.9598155 | 3.15×10^{-9} | [0.959703, 0.959879] | 1.005078 | 2.85×10^{-9} | [1.004985, 1.005142] |
| <i>sp3bar3</i> | 0.8739616 | 6.72×10^{-7} | [0.872347, 0.875052] | 1.003729 | 8.98×10^{-7} | [1.002500, 1.005500] |
| <i>sp3bar3_18</i> | 0.9410686 | 1.46×10^{-5} | [0.938852, 0.955264] | 0.999010 | 4.08×10^{-6} | [0.997219, 1.004945] |
| <i>sp3bar3_25</i> | 0.7597931 | 5.64×10^{-10} | [0.759743, 0.759833] | 0.992174 | 2.89×10^{-9} | [0.992067, 0.992246] |

TABLE 3

Summary of the posterior distribution of the diagonal and one non-diagonal element of Σ , from the 4 base astrophysical models

| Model | σ_{11} | | σ_{22} | | σ_{12} | |
|-------------------|---|--|---|--|--|--|
| | 95% CI | | 95% CI | | 95% CI | |
| <i>bar_6</i> | [5.40×10^{-5} , 4.0×10^{-4}] | | [6.20×10^{-5} , 4.76×10^{-4}] | | [0, 1.30×10^{-5}] | |
| <i>sp3bar3</i> | [3.66×10^{-3} , 1.03×10^{-2}] | | [6.53×10^{-3} , 1.83×10^{-2}] | | [-6.40×10^{-5} , 2.68×10^{-4}] | |
| <i>sp3bar3_18</i> | [1.45×10^{-3} , 1.68×10^{-1}] | | [1.29×10^{-3} , 1.50×10^{-1}] | | [-1.19×10^{-4} , 2.16×10^{-3}] | |
| <i>sp3bar3_25</i> | [1.21×10^{-4} , 5.69×10^{-4}] | | [1.13×10^{-4} , 5.21×10^{-4}] | | [-1.00×10^{-6} , 1.50×10^{-5}] | |

in all these tables are the small posterior variances of the quantities in question; this is indicative of the fact that the data sets we used, in spite of the relatively smaller size compared to the astronomically large data sets used in the previous approaches in the literature, are very much informative, given our vector-variate \mathcal{GP} -based Bayesian approach. Owing to our Gaussian Process approach, the posterior of Σ should be close to the null matrix *a posteriori* if the choice of the design set and the number of design points are adequate. Quite encouragingly, Table 3, shows that indeed Σ is close to the null matrix *a posteriori*, for all the four models, signifying that the unknown velocity function has been learned well in all the cases.

4.1. Comparison with results in astrophysical literature

The estimates of the angular separation of the long axis of the bar from the Sun-Galactic centre line and the rotation rate of the bar compare favourably with results obtained by Chakrabarty (2007), Englmaier and Gerhard (1999), Debattista et al. (2002), Benjamin et al. (2005), Antoja et al. (2011). A salient feature of our implementation is the vastly smaller data set that we needed to invoke than any of the methods reported in the astronomical literature, in order

to achieve the learning of the two-dimensional vector \mathbf{S} - in fact while in the calibration approach of Chakrabarty (2007), the required sample size is of the order of 3,500, in our work, this number is 50. Thus, data sufficiency issues, when a concern, are well tackled by our method.

Upon the analyses of the viable astrophysical models of the Galaxy, Chakrabarty (2007) reported the result that $r_{\odot} \in [1.9375, 2.21]$ in model units while $\phi_{\odot} \in [0^{\circ}, 30^{\circ}]$, where these ranges correspond to the presented uncertainties on the estimates, which were however, rather unsatisfactorily achieved (see Section 2). The values of the components of \mathbf{S} , learnt in our work, overlap well with these results. As mentioned above, the models *sp3bar3_18*, *sp3bar3* and *sp3bar3_25* are distinguished by distinct values of the ratios of the rotational rates of the spiral pattern Ω_s to that of the bar (Ω_b) in the Galaxy. Then the derived estimate for Ω_b (Table 1) suggests values of Ω_s of the Milky Way spiral.

Another point that merits mentions is that the estimates of r_{\odot} and ϕ_{\odot} presented by Chakrabarty (2007) exclude the model *sp3bar3* which could not be used to yield estimates given the highly scattered nature of the corresponding p -value distribution. Likewise, in our work, the same model manifests maximal multi-modality amongst the others, but importantly, our approach allows for the representation of the full posterior density using which, the computation of the 95% HPDs is performed.

That the new method is able to work with smaller velocity data sets, is an important benefit, particularly in extending the application to galaxies other than our own, in which small numbers of individual stars are going to be tracked in the very near future for their velocities, under observational programmes such as PANStarrs (Johnston et al., 2009) and GAIA (Lindegren et al., 2007; Kucinskas et al., 2005, <http://www.rssd.esa.int/index.php?project=GAIA&page=index>); the sample sizes of measured stellar velocity vectors in these programmes will be much smaller in external galaxies than what has been possible in our own. At the same time, our method is advanced as a template for the analysis of the stellar velocity data that is available for the Milky Way, with the aim of learning a high-dimensional Galactic parameter vector; by extending the scope of the dynamical simulations of the Galaxy, performed on different astrophysical models of the Milky Way, the Milky Way models will be better constrained. The mission GAIA - a mission of the European Space Agency - is set to provide large sets of stellar velocity data all over the Milky Way. Our method, in conjunction with astrophysical models, can allow for fast learning of local and global model parameters of the Galaxy.

5. Model fitting

In this section we compare the test data with predictions for the observable that we make at a summary $\tilde{\mathbf{s}}$ of the posterior of the model parameter vector \mathbf{S} . To achieve this, we first need to provide a suitable estimator of the function $\xi(\cdot)$ that defines the relationship between the observable and the model parameter \mathbf{S} . We attempt to write the conditional distribution of $\xi(\tilde{\mathbf{s}})$ given the augmented data

\mathcal{D}_a that comprises training data \mathcal{D}_s , augmented by test data $v^{(test)}$. Here we consider the test data $v^{(test)}$ realised at $\mathbf{S} = \tilde{\mathbf{s}}$, where we use different candidates for $\tilde{\mathbf{s}}$. In particular, we choose $\tilde{\mathbf{s}}$ to be (1) the median $\mathbf{s}^{(median)}$ of the posterior of \mathbf{S} given \mathcal{D}_a , (2) the mode $\mathbf{s}^{(mode)}$ of this posterior, (3) or $\mathbf{s}^{(u)}$, $u=1,2,3,4$ —the end points of the disjoint 95% HPD region of the posterior of \mathbf{S} (see Table 1).

Since $\{\boldsymbol{\xi}(\mathbf{s}_1), \dots, \boldsymbol{\xi}(\mathbf{s}_n), \boldsymbol{\xi}(\tilde{\mathbf{s}})\}$ is jointly matrix-normal, $[\boldsymbol{\xi}(\tilde{\mathbf{s}})|\boldsymbol{\xi}(\mathbf{s}_1), \dots, \boldsymbol{\xi}(\mathbf{s}_n)] \equiv [\boldsymbol{\xi}(\tilde{\mathbf{s}})|\mathcal{D}_s]$, is jk -variate normal. The mean function of this multivariate normal, at different $\tilde{\mathbf{s}}$, is then compared to the test data. Thus, the estimate of the function that we seek is $\mathbb{E}[\boldsymbol{\xi}(\mathbf{S})|\mathcal{D}_s, \mathbf{S}, \mathbf{Q}]$, given the dependence of $\boldsymbol{\xi}(\cdot)$ on the smoothness parameters (elements of \mathbf{Q}) that we anticipate.

However, we only know the conditional of $\boldsymbol{\xi}(\cdot)$ on all the \mathcal{GP} parameters, including the ones that we do not learn from the data, namely \mathbf{B} and \mathbf{C} . So we need to marginalise $[\boldsymbol{\xi}(\cdot) | \boldsymbol{\Sigma}, \mathbf{B}, \mathbf{C}, \mathbf{Q}, \mathcal{D}_s]$ over \mathbf{B} and \mathbf{C} . To achieve this, we need to invoke the conditional distribution of \mathbf{B} and \mathbf{C} with respect to the other \mathcal{GP} parameters and \mathcal{D}_s . We recall the priors on the \mathcal{GP} parameters $\mathbf{B}, \boldsymbol{\Sigma}, \mathbf{C}$ (from Section 2.4) to write $\pi(\mathbf{B}, \boldsymbol{\Sigma}, \mathbf{C}) \propto |\boldsymbol{\Sigma}|^{-(k+1)/2} |\mathbf{C}|^{-(j+1)/2}$. It then follows that

$$[\mathbf{B} | \boldsymbol{\Sigma}, \mathbf{C}, \mathbf{Q}, \mathcal{D}_s] \sim \mathcal{N}_{m,jk}(\hat{\mathbf{B}}_{GLS}, (\mathbf{H}_D^T \mathbf{A}_D^{-1} \mathbf{H}_D)^{-1}, \boldsymbol{\Omega}), \quad (5.1)$$

where, we recall from Section 2.1 that we had set $m = d + 1$, with $\mathbf{S} \in \mathbb{R}^d$. Here, $\hat{\mathbf{B}}_{GLS} = (\mathbf{H}_D^T \mathbf{A}_D^{-1} \mathbf{H}_D)^{-1} (\mathbf{H}_D^T \mathbf{A}_D^{-1} \mathcal{D}_s)$. Marginalising the jk -variate normal that is the conditional $[\boldsymbol{\xi}(\cdot) | \mathbf{B}, \boldsymbol{\Sigma}, \mathbf{C}, \mathbf{Q}, \mathcal{D}_s]$ over \mathbf{B} (using Equation 5.1), it can be shown that

$$[\boldsymbol{\xi}(\cdot) | \boldsymbol{\Sigma}, \mathbf{C}, \mathbf{Q}, \mathcal{D}_s] \sim \mathcal{N}_{jk}(\boldsymbol{\mu}_2(\cdot), a_2(\cdot, \cdot) \boldsymbol{\Omega}), \quad (5.2)$$

where

$$\boldsymbol{\mu}_2(\cdot) = \hat{\mathbf{B}}_{GLS}^T \mathbf{h}(\cdot) + (\mathcal{D}_s - \mathbf{H}_D \hat{\mathbf{B}}_{GLS})^T \mathbf{A}_D^{-1} \boldsymbol{\sigma}_D(\cdot); \quad (5.3)$$

$$\begin{aligned} a_2(\mathbf{s}_1, \mathbf{s}_2) &= a_1(\mathbf{s}_1, \mathbf{s}_2) + [\mathbf{h}(\mathbf{s}_1) - \mathbf{H}_D^T \mathbf{A}_D^{-1} \mathbf{s}_D(\mathbf{s}_1)]^T (\mathbf{H}_D^T \mathbf{A}_D^{-1} \mathbf{H}_D)^{-1} \\ &\quad [\mathbf{h}(\mathbf{s}_2) - \mathbf{H}_D^T \mathbf{A}_D^{-1} \mathbf{s}_D(\mathbf{s}_2)]. \end{aligned} \quad (5.4)$$

We define $(n - m) \hat{\boldsymbol{\Omega}}_{GLS} = (\mathcal{D}_s - \mathbf{H}_D \hat{\mathbf{B}}_{GLS})^T \mathbf{A}_D^{-1} (\mathcal{D}_s - \mathbf{H}_D \hat{\mathbf{B}}_{GLS})$, i.e. $(n - m) \hat{\boldsymbol{\Omega}}_{GLS} = \mathcal{D}_s^T \mathbf{M} \mathcal{D}_s$, with $\mathbf{M} = \mathbf{A}_D^{-1} - \mathbf{A}_D^{-1} \mathbf{H}_D (\mathbf{H}_D^T \mathbf{A}_D^{-1} \mathbf{H}_D)^{-1} \mathbf{H}_D^T \mathbf{A}_D^{-1}$.

We consider the mean $\boldsymbol{\mu}_2(\cdot)$ of the conditional posterior given by (5.3) as a suitable estimator of the velocity function in our case. Note that $\boldsymbol{\mu}_2$ involves the unknown smoothness parameters; we plug-in the corresponding posterior medians 0.874254, 1.003545 for these.

It is important to mention that though the mean and variance in Equations 5.3 and Equation 5.4 were developed using \mathcal{D}_s , in our construction of the velocity function estimator $\boldsymbol{\mu}_2$, \mathcal{D}_a is implemented, where \mathcal{D}_a is obtained by augmenting \mathcal{D}_s with $\mathbf{v}^{(test)}$ that is realised at $\mathbf{S} = \tilde{\mathbf{s}}$. The underlying theory remains the same as above.

It is important to note that $\boldsymbol{\mu}_2(\mathbf{S})$, where \mathbf{S} is the unknown location, is a random variable, and even though the posterior of $\boldsymbol{\Sigma}$ is concentrated around

the null matrix, the variance of $\mu_2(\mathbf{S})$ is not $\mathbf{0}$, thanks to the fact that \mathbf{S} does not have $\mathbf{0}$ variance. Consequently, the posterior variance of $\xi(\mathbf{S})$ does not have $\mathbf{0}$ variance. To see this formally, note that

$$\begin{aligned} \text{Var} [\xi(\mathbf{S})|\mathcal{D}_a] &= \text{Var} [\mathbb{E} \{\xi(\mathbf{S})|\Sigma, \mathbf{C}, \mathbf{Q}, \mathbf{S}, \mathcal{D}_a\}] \\ &\quad + \mathbb{E} [\text{Var} \{\xi(\mathbf{S})|\Sigma, \mathbf{C}, \mathbf{Q}, \mathbf{S}, \mathcal{D}_a\}] \\ &= \text{Var} [\mu_2(\mathbf{S})|\mathcal{D}_a] + \mathbb{E} [a_2(\mathbf{S}, \mathbf{S})\Omega|\mathcal{D}_a]. \end{aligned} \quad (5.5)$$

Since the posterior $[\Sigma|\mathcal{D}_a]$ is concentrated around the $k \times k$ -dimensional null matrix, it follows that the posterior $[\Omega|\mathcal{D}_a]$ is also concentrated around the $jk \times jk$ -dimensional null matrix. Hence, in (5.5), $\mathbb{E} [a_2(\mathbf{S}, \mathbf{S})\Omega|\mathcal{D}_a] \approx \mathbf{0}^{(jk \times jk)}$. However, the first part of (5.5), $\text{Var} [\mu_2(\mathbf{S})|\mathcal{D}_a]$, is strictly (and significantly) positive, showing that the variance of the posterior of $\xi(\mathbf{S})$ is significantly positive.

The above result shows that it should not be expected that the observed test velocity data $\mathbf{v}^{(test)}$ will be predicted accurately by $\mu_2(\mathbf{s})$, for any given \mathbf{s} . This is in contrast with the usual Gaussian process emulators, where the argument of the unknown function is non-random, so that if the posterior of the function variance is concentrated around $\mathbf{0}$, then the posterior variance of the emulator would be close to $\mathbf{0}$.

In Figure 5 we illustrate, in the case of *sp3bar3* (the most chaotic model), the degree of agreement of $\mu_2(\mathbf{s})$ with $\mathbf{v}^{(test)}$ for different choices of \mathbf{s} . We compare with $\mathbf{v}^{(test)}$ the predictions $\mu_2(\mathbf{s}^{(mode)})$, $\mu_2(\tilde{\mathbf{s}})$ and $\mu_2(\mathbf{s}^{(u)})$; $u = 1, 2, 3, 4$. Here, $\mathbf{s}^{(mode)} = (1.73, 18.8^\circ)$ is the (component-wise) posterior mode and $\tilde{\mathbf{s}} = (2.2, 35^\circ)$ is a point somewhat close to the (component-wise) posterior median $\mathbf{s}^{(median)} = (1.994478, 33.59429^\circ)$ (grid-point closest to $\mathbf{s}^{(median)}$).

As observed in Figure 5 the best fit of $\mathbf{v}^{(test)}$ has been provided by $\mu_2(\tilde{\mathbf{s}})$ where $\tilde{\mathbf{s}}$ is close to the median $\mathbf{s}^{(median)}$; as the point $(\mathbf{s}^{(median)}, \mathbf{v}^{(test)})$ is in the training data constituting μ_2 , this is to be expected. The estimators $\mu_2(\mathbf{s}^{(mode)})$ and $\mu_2(\mathbf{s}^{(1)})$ perform somewhat reasonably, but the remaining estimators $\mu_2(\mathbf{s}^{(u)})$; $u = 2, 3, 4$ do not perform adequately, signifying the effect of variability of our estimator due the posterior of \mathbf{S} .

While it is the randomness of the argument \mathbf{S} of the unknown function $\xi(\cdot)$ that causes the variability of our estimator, such variability is highest in the most chaotic of the 4 base astrophysical models (*sp3bar3*), and least in the only non-chaotic base astrophysical model (*bar_6*). A similar exercise of predicting $\mathbf{v}^{(test)}$ using the training data simulated from this non-chaotic base model gives excellent fits at all the aforementioned used values of \mathbf{S} ; see Figure 6.

6. Discussions

Computational complexity scales only linearly with the dimensionality of the unknown model parameter \mathbf{S} . Thus, porting a training data comprised of n independent values of \mathbf{S} , \mathbf{s}_i , $i = 1, \dots, n$, where \mathbf{s}_i is a d -dimensional vector, $d > 2$, is not going to render the computational times infeasible. This allows for the learning of high-dimensional model parameter vectors in our method.

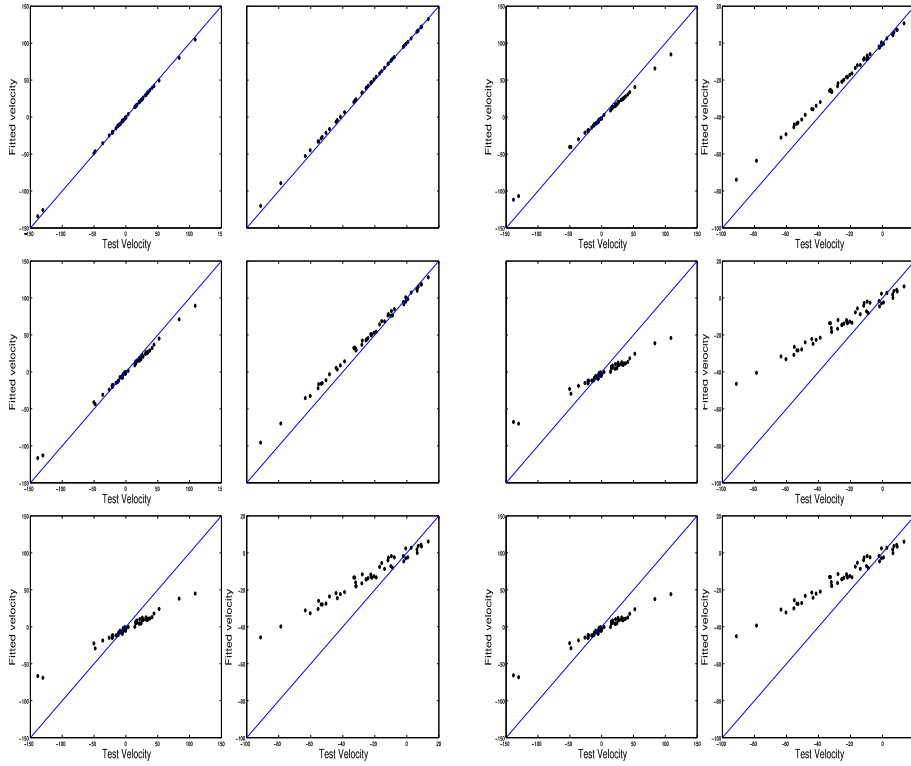


FIG 5. Prediction of $\mathbf{v}^{(test)}$ for model $sp3bar3$: plots of 2 components of $\mu_2(\mathbf{s})$ against $\mathbf{v}^{(test)}$ for $\mathbf{s} = \bar{\mathbf{s}}$ (2 left hand sided panels on the top row), $\mathbf{bs} = \mathbf{s}^{(mode)}$ (2 right panels on the top), $\mathbf{s} = \mathbf{s}^{(1)}$ (2 left panels in the middle row), $\mathbf{s} = \mathbf{s}^{(2)}$ (2 left panels in the middle row), $\mathbf{s} = \mathbf{s}^{(3)}$ (2 left panels in the lowest row), $\mathbf{s} = \mathbf{s}^{(4)}$ (2 right panels in the lowest row).

In contrast to the situation with increasing the dimensionality of the unknown model parameter, increasing the dimensionality of the measurable will but imply substantial increase in the run time, since the relevant computational complexity then scales non-linearly, as about $O(k^3)$, (in addition to the cost of k square roots), where k is the dimensionality of the observed variable. This is because of the dimensionality of the aforementioned Σ matrix is $k \times k$, and the inverse of this enters the computation of the posterior via the definition $\hat{\mathbf{C}}_{GLS,aug}$. Thus, for example, increasing the dimensions of the measurable from 2 to 4 increases the run time 8-fold, which is a large jump in the required run time. However, for most applications, we envisage the expansion of the dimensionality of the unknown model parameter, i.e. d , rather than that of the measurable, i.e. k . Thus, the method is expected to yield results within acceptable time frames, for most practical applications.

The other major benefit of our work is that it allows for organic learning of the smoothness parameters, rather than results being subject to *ad hoc* choices of the same.

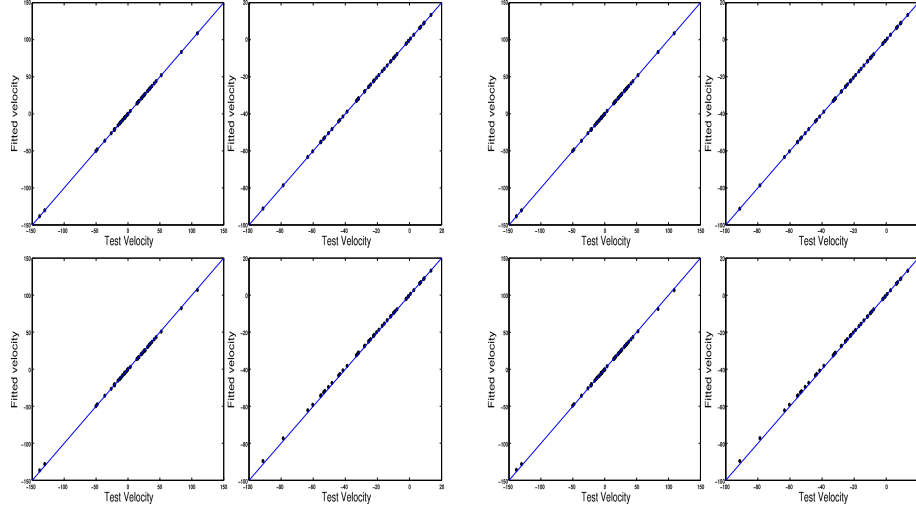


FIG 6. Prediction of $\mathbf{v}^{(test)}$ for model bar6: plots of 2 components of $\boldsymbol{\mu}_2(\mathbf{s})$ against $\mathbf{v}^{(test)}$ for $\mathbf{s} = \tilde{\mathbf{s}}$ (2 adjacent panels on the left hand side of the top row), $\mathbf{s}^{(mode)}$ (2 panels on the right of top row), $\mathbf{s}^{(1)}$ (2 panels on the left in the lower row), $\mathbf{s}^{(2)}$ (2 panels on the left in the lower row).

As more Galactic simulations spanning a greater range of model parameters become available, the rigorous learning of such Milky Way parameters using our method will become possible, given the available stellar velocity data. This will enhance the quality of our knowledge about our own galaxy. That our method allows for such learning even for under-abundant systems, is encouraging for application of a similar analysis to galaxies other than our own, in which system parameters may be learnt using the much smaller available velocity data sets, compared to the situation in our galaxy.

Supplementary Material

Supplementary section for “Bayesian Nonparametric Estimation of Milky Way Parameters Using Matrix-Variate Data, in a New Gaussian Process Based Method”

(doi: [10.1214/15-EJS1037SUPP](https://doi.org/10.1214/15-EJS1037SUPP); .pdf). Some background details on the application to the Milky Way are discussed in Section S-1 of the attached supplementary material. Section S-2 discusses the details of the dynamical simulations that lead to the training data set used in our supervised learning of the Milky Way feature parameters. In Section S-3 we present details of the TMCMC methodology that we use here. S-4 discusses the cross-validation of our model and methodology, on simulated as well the real stellar velocity data. The effect of chaos on the modality of the posterior distributions of our unknowns is discussed in Section S-5.

References

- O'HAGAN, A. (1978), "Curve Fitting and Optimal Design for Prediction," *Journal of the Royal Statistical Society B*, 40, 1–42. [MR0512140](#)
- ANTOJA, T., FIGUERAS, F., ROMERO-GÓMEZ, M., PICHARDO, B., VALENZUELA, O., and MORENO, E. (2011), "Understanding the spiral structure of the Milky Way using the local kinematic groups," *Monthly Notices of the Royal Astronomical Society*, 418, 1423–1440.
- ANTOJA, T., VALENZUELA, O., PICHARDO, B., MORENO, E., FIGUERAS, F., and FERNÁNDEZ, D. (2009), "Stellar Kinematic Constraints on Galactic Structure Models Revisited: Bar and Spiral Arm Resonances," *Astrophysical J. Letters*, 700, L78–L82.
- AUMER, M., and BINNEY, J. J. (2009), "Kinematics and history of the solar neighbourhood revisited," *Monthly Notices of the Royal Astronomical Society*, 397, 1286–1301.
- BENJAMIN, R. A., CHURCHWELL, E., BABLER, B. L., INDEBETOUW, R., MEADE, M. R., WHITNEY, B. A., WATSON, C., WOLFIRE, M. G., WOLFF, M. J., IGNACE, R., BANIA, T. M., BRACKER, S., CLEMENS, D. P., CHOMIUK, L., COHEN, M., DICKEY, J. M., JACKSON, J. M., KOBULNICKY, H. A., MERCER, E. P., MATHIS, J. S., STOLOVY, S. R., and UZPEN, B. (2005), "First GLIMPSE Results on the Stellar Structure of the Galaxy," *Astrophysical J. Letters*, 630, L149–L152.
- HOFMANN, B. (2011), Ill-posedness and regularization of inverse problems—a review on mathematical methods, in *The Inverse Problem. Symposium ad Memoriam H. v. Helmholtz, H. Lubbig (Ed)*. Akademie-Verlag, Berlin; VCH, Weinheim, pp. 45–66. [MR1399148](#)
- BINNEY, J., and MERRIFIELD, M. (1998), *Galactic Astronomy*, Princeton: Princeton University Press.
- BLIGHT, B. J. N., and OTT, L. (1975), "A Bayesian Approach to Model Inadequacy for Polynomial Regression," *Biometrika*, 62(1), 79–88. [MR0373174](#)
- CARLIN, B. P., and LOUIS, T. A. (1996), *Bayes and Empirical Bayes Methods for Data Analysis*, London: Chapman and Hall. Second Edition. [MR1427749](#)
- CARREIRA-PERPIÑÁN, M. A. (2001), Continuous latent variable models for dimensionality reduction and sequential data reconstruction, Doctoral thesis, University of Sheffield.
- CARVALHO, C. M., and WEST, M. (2007), "Dynamic Matrix-Variate Graphical Models," *Bayesian Analysis*, 2, 69–98. [MR2289924](#)
- CHAKRABARTY, D. (2007), "Phase Space around the Solar Neighbourhood," *Astronomy & Astrophysics*, 467, 145.
- CHAKRABARTY, D., BISWAS, M., and BHATTACHARYA, S. (2015), Supplementary section for "Bayesian Nonparametric Estimation of Milky Way Parameters Using Matrix-Variate Data, in a New Gaussian Process Based Method". DOI: [10.1214/15-EJS1037SUPP](#).
- CHAKRABARTY, D., and SIDERIS, I. (2008), "Chaos in Models of the Solar Neighbourhood," *Astronomy & Astrophysics*, 488, 161.

- CRESSIE, N. A. C. (1993), *Statistics for Spatial Data*, New York: Wiley. [MR1239641](#)
- DAWID, A. P. (1981), “Some Matrix-Variate Distribution Theory: Notational Considerations and a Bayesian Application,” *Biometrika*, 68, 265–274. [MR0614963](#)
- DEBATTISTA, V. P., GERHARD, O., and SEVENSTER, M. N. (2002), “The pattern speed of the OH/IR stars in the Milky Way,” *Monthly Notice of Royal Astronomical Society*, 334, 355.
- DEHNEN, W. (2000), “The Effect of the Outer Lindblad Resonance of the Galactic Bar on the Local Stellar Velocity Distribution,” *Astronomical Journal*, 11, 800.
- DUTTA, S., and BHATTACHARYA, S. (2013), “Markov Chain Monte Carlo Based on Deterministic Transformations,”. Accepted in Statistical Methods; available at [arxiv:1106.5850v3](https://arxiv.org/abs/1106.5850v3) with supplementary section in arxiv.org/pdf/1306.6684.
- ENGLMAIER, P., and GERHARD, O. (1999), “Gas dynamics and large-scale morphology of the Milky Way galaxy,” *Monthly Notices of the Royal Astronomical Society*, 304, 512–534.
- FUX, R. (1997), “3D self-consistent N-body barred models of the Milky Way. I. Stellar dynamics,” *Astronomy & Astrophysics*, 327, 983–1003.
- FUX, R. (2001), “Order and chaos in the local disc stellar kinematics induced by the Galactic bar,” *Astronomy & Astrophysics*, 373, 511–535.
- WALKER, G. and FORD, J. (1969), “Amplitude Instability and Ergodic Behavior for Conservative Nonlinear Oscillator Systems,” *Physical Review*, 188, 416–432. [MR0260978](#)
- GOLUBOV, O. (2012), Modelling the Milky Way Disk, Doctoral thesis, University of Heidelberg.
- HOFF, P. D. (2011), “Hierarchical multilinear models for multiway data,” *Computational Statistics & Data Analysis*, 55, 530–543. [MR2736574](#)
- JOHNSTON, K., BULLOCK, J. S., and STRAUSS, M. (2009), The Milky Way and Local Volume as Rosetta Stones in Galaxy Formation, in *astro2010: The Astronomy and Astrophysics Decadal Survey*, Vol. 2010, p. 142.
- KABANIKHIN, S. I. (2008), “Definitions and examples of inverse and ill-posed problems,” *J. Inv. Ill-Posed Problems*, 16, 317–357. [MR2426856](#)
- KUCINSKAS, A., LINDEGREN, L., and VANSEVICIUS, V. (2005), Beyond the Galaxy with Gaia: Evolutionary Histories of Galaxies in the Local Group, in *The Three-Dimensional Universe with Gaia*, eds. C. Turon, K. S. O’Flaherty, and M. A. C. Perryman, Vol. 576 of *ESA Special Publication*.
- LINDEGREN, L., BABUSIAUX, C., BAILER-JONES, C., BASTIAN, U., BROWN, A., CROPPER, M., HG, E., JORDI, C., KATZ, D., VAN LEEUWEN, F., LURI, X., MIGNARD, F., DE BRUIJNE, J., and PRUSTI, T. (2007), The Gaia mission: science, organization and present status, in *A Giant Step: from Milli- to Micro-arcsecond Astrometry, Proceedings IAU Symposium No. 248*, eds. W. Jin, I. Platais, and M. Perryman, pp. 217–223.
- MATERN, B. (1986), *Spatial Variation (2nd ed.)*, Springer: Springer-Verlag.

[MR0867886](#)

- MINCHEV, I., BOILY, C., SIEBERT, A., and BIENAYME, O. (2010), “Low-velocity streams in the solar neighbourhood caused by the Galactic bar,” *Monthly Notices of the Royal Astronomical Society*, 407, 2122–2130.
- MINCHEV, I., QUILLEN, A. C., WILLIAMS, M., FREEMAN, K. C., NORDHAUS, J., SIEBERT, A., and BIENAYMÉ, O. (2009), “Is the Milky Way ringing? The hunt for high-velocity streams,” *Monthly Notices of the Royal Astronomical Society*, 396, L56–L60.
- NEAL, R. M. (1998), Regression and classification using Gaussian process priors (with discussion), in *Bayesian Statistics 6*, ed. J. M. B. et. al, Oxford University Press, pp. 475–501. [MR1723510](#)
- PERRYMAN, M. (2012), *Astronomical Applications of Astrometry: Ten Years of Exploitation of the Hipparcos Satellite Data*, Cambridge: Cambridge University Press.
- RASMUSSEN, C. E., and WILLIAMS, C. K. I. (2006), *Gaussian Processes for Machine Learning*, MIT: The MIT Press. [MR2514435](#)
- SANTNER, T. J., WILLIAMS, B. J., and NOTZ, W. I. (2003), *The design and analysis of computer experiments*, Springer Series in Statistics, New York, Inc.: Springer-Verlag. [MR2160708](#)
- SCHÖLKOPF, B., and SMOLA, A. J. (2002), *Learning with Kernels*, MIT: MIT Press.
- SENGUPTA, A. (2003), “Toward a Theory of Chaos,” *International Journal of Bifurcation and Chaos*, 13, 3147–3233. [MR2031142](#)
- SIMONE, R. D., WU, X., and TREMAINE, S. (2004), “The stellar velocity distribution in the solar neighbourhood,” *Monthly Notices of the Royal Astronomical Society*, 350, 627.
- SNELSON, E. L. (2007), Flexible and efficient Gaussian process models for machine learning, Doctoral thesis, University of London. [MR3211196](#)
- STUART, A. (2013), “Bayesian Approach to Inverse Problems,”. provide an introduction to the forthcoming book *Bayesian Inverse Problems in Differential Equations* by M. Dashti, M. Hairer and A.M. Stuart; available at [arXiv:math/1302.6989](#).
- TARANTOLA, A. (2005), *Inverse Problem Theory and Methods for Model Parameter Estimation*, Philadelphia: SIAM. [MR2130010](#)
- GNEITING, T., KLEIBER, W. and SCHLATHER, M. (2010), “Matern Cross-Covariance Functions for Multivariate Random Fields,” *Journal of the American Statistical Association*, 105(491), 1167–1177. [MR2752612](#)