

Published in final edited form as:

Neuroimage. 2021 December 15; 245: 118715. doi:10.1101/2021.04.05.438429.

Warped Bayesian Linear Regression for Normative Modelling of Big Data

Charlotte J. Frazza^{a,b}, Richard Dinga^a, Christian F. Beckmann^{a,b,d}, Andre F. Marquand^{a,b,c}

^aDonders Centre for Cognitive Neuroimaging, Donders Institute for Brain, Cognition and Behaviour, Radboud University, Nijmegen, the Netherlands

^bDepartment of Cognitive Neuroscience, Radboud University Medical Centre, Nijmegen, the Netherlands

^cDepartment of Neuroimaging, Centre for Neuroimaging Sciences, Institute of Psychiatry, King's College London, London, UK

^dOxford Centre for Functional Magnetic Resonance Imaging of the Brain (FMRIB), University of Oxford, Oxford, UK

Abstract

Normative modelling is becoming more popular in neuroimaging due to its ability to make predictions of deviation from a normal trajectory at the level of individual participants. It allows the user to model the distribution of several neuroimaging modalities, giving an estimation for the mean and centiles of variation. With the increase in the availability of big data in neuroimaging, there is a need to scale normative modelling to big data sets. However, the scaling of normative models has come with several challenges.

So far, most normative modelling approaches used Gaussian process regression, and although suitable for smaller datasets (up to a few thousand participants) it does not scale well to the large cohorts currently available and being acquired. Furthermore, most neuroimaging modelling methods that are available assume the predictive distribution to be Gaussian in shape. However, deviations from Gaussianity can be frequently found, which may lead to incorrect inferences, particularly in the outer centiles of the distribution. In normative modelling, we use the centiles to give an estimation of the deviation of a particular participant from the 'normal' trend. Therefore, especially in normative modelling, the correct estimation of the outer centiles is of utmost importance, which is also where data are sparsest.

Here, we present a novel framework based on Bayesian Linear Regression with likelihood warping that allows us to address these problems, that is, to scale normative modelling elegantly to big data cohorts and to correctly model non-Gaussian predictive distributions. In addition, this method provides also likelihood-based statistics, which are useful for model selection.

To evaluate this framework, we use a range of neuroimaging-derived measures from the UK Biobank study, including image-derived phenotypes (IDPs) and whole-brain voxel-wise measures derived from diffusion tensor imaging. We show good computational scaling and improved accuracy of the warped BLR for certain IDPs and voxels if there was a deviation from normality of these parameters in their residuals.

The present results indicate the advantage of a warped BLR in terms of; computational scalability and the flexibility to incorporate non-linearity and non-Gaussianity of the data, giving a wider range of neuroimaging datasets that can be correctly modelled.

Keywords

Machine learning; UK Biobank; Big Data; Bayesian Linear Regression; Normative Modelling

1 Introduction

Big data has become more widely available in neuroimaging (UK Biobank, ENIGMA, ABCD study, PNC, among others) [1], [2], [3], [4]. This has ignited a renewed interest in modelling normal brain development, to estimate quantitative brain-behaviour mappings and capture deviations from such models to derive neurobiological markers of different psychiatric disorders. These neurobiological markers could move us closer towards individualized and precision medicine [5]. Until now, the neurobiological markers for psychiatric disorders have been mostly developed with studies that used classifiers trained in a case-control setting. Counter-intuitively, an increase in sample size has shown to reduce the accuracy of classifying cases from controls for psychiatric disorders [6]. One of the main reasons for this decrease in accuracy has been posed to be the heterogeneity in the participants both biologically and behaviorally, which can only fully be captured by a large data set [6]. Normative modelling is an emerging method used to understand this heterogeneity in the population. Similar to growth charts in pediatric medicine, which describe the distribution of height or weight of children according to their age and sex, normative models can be used to model the distribution of neuroimaging derived phenotypes in a population, including the mean and centiles of variation [7], according to age, gender, or other demographic or clinical variables [8]. The deviations from this normative range can be quantified statistically, for example as Z-scores, which have been linked to several psychiatric disorders [7], [9], [10], [11], [12], [13].

Although promising, there are still certain challenges in performing normative modelling on big neuroimaging data. First of all, Normative models have been mainly developed using Gaussian process regression. [14]. Gaussian process regression is flexible and accurate, but a drawback is its computational complexity, which is governed by the need to compute the exact inverse of the covariance matrix. This inversion scales poorly with an increase in data points [15]. Therefore, using these models on large datasets requires extensive computational power and is often not feasible (typically beyond a few thousand subjects). Furthermore, most normative models assume the modelled distribution is Gaussian. However, distributions diverging from Gaussianity are frequently found in specific neuroimaging modalities. These non-Gaussian signals cannot be accounted for using a standard normative model based on Gaussian process regression. We argue that modelling non-Gaussianity is important in general and is frequently overlooked by the neuroimaging community in that most regression methods used in practice –often implicitly– assume Gaussian residuals. Thus, there is a need to develop methods that can flexibly handle the computational demand and non-Gaussianity of big data sets.

In this paper, we propose a next-generation framework based on Bayesian linear regression (BLR) to address these challenges. We introduce an extension of the BLR method for accurately modelling non-Gaussian distributions using a likelihood warping technique, giving a warped BLR model. The new framework has several benefits over previously used methods: (i) A BLR model can use a linear combination of non-linear basis functions (such as B-splines) which can be considered to provide a low-rank approximation of the Gaussian process regression models [16]. However, the BLR model has considerably better computational scaling, since the complexity of the model is fixed according to a set of basis-functions. Therefore, the model can be scaled much more easily to large datasets. Furthermore, a set of model coefficients can be estimated that can easily be shared without the need to share the data (e.g. to compute a cross-covariance matrix for new data points), thus making it easier to make predictions on new datasets. (ii) The non-Gaussianity of the residuals can be modelled by the flexible warping of the Gaussian function, which gives more options to model different types of neuroimaging data that cannot be accurately modelled using a standard BLR. (iii) Efficient model selection criteria are naturally defined for the warped BLR through the marginal likelihood and can be calculated in closed form. The marginal likelihood gives a balance between model complexity and model fit. This can aid in choosing the optimal model for a specified imaging modality.

We will demonstrate this model by testing it on different types of neuroimaging data derived from the UK Biobank dataset. The UK Biobank dataset has several magnetic resonance imaging (MRI) imaging modalities, including structural and functional brain data. With over 40,000 participants' MRI data from 40 to 80 years old, this provides a rich set of different neuroimaging data and defines a benchmark for future population-based studies. In this work, we will present the warping function and recommend how to use it for several data modalities. First, we give an illustrative example using image-derived phenotypes (IDPs), which are convenient and widely used summary measures of brain function and structure [17]. Specifically, we will show a detailed example of estimating a normative model for white matter hyperintensities (WMHs). WMHs have been shown before to demonstrate quite non-Gaussian behaviour [18], and are therefore a good example where the warped BLR could be preferred over the B-spline BLR. Second, we show the scalability of this method by performing a whole-brain analysis for certain diffusion tensor imaging (DTI) measures. We use DTI measurements, as there are large associations with age and we expect certain non-linear and non-Gaussian trends in the data [19].

Finally, we want to evaluate the link between brain imaging abnormality scores and behaviour. Therefore, deviations from normal brain functioning are associated with cognitive functioning. The deviations are captured by Z-scores, which are shown to correlate with measures of intelligence in the UK Biobank dataset, such as; numerical memory, reaction time and visual memory.

In summary, the main contributions of the paper are to give: (i) a new comprehensive framework for big data normative modelling; (ii) the introduction of the novel methodological approach for modelling non-Gaussian response variables; (iii) an extensive and didactic evaluation of this framework on the UK Biobank cohort and (iv) a demonstration of the 'Predictive Clinical Neuroscience software toolkit' ([PCNtoolkit](#)) for

big data normative modelling. Ultimately, we hope this paper will give deeper insight into the application of normative models on different types of neuroimaging modalities.

2 Materials and methods

2.1 Sample

All the data used came from the UK Biobank imaging dataset [1]. Full details on the design of the study and the preprocessing steps can be found in subsequent papers [17], [20]. Briefly, the data used contains around 10,000 participants of the 2017 release and additional longitudinal data of around 5,000 subjects of the 2020 release. The participants were between 40 and 80 years of age, with around 47 % males.

In this study, two types of analyses were performed using different datasets. For the first analysis, a dataset containing IDPs was used. For consistency with existing work, the IDPs were processed using FUNPACK [21], which is an automatic normalisation, parsing and cleaning kit, developed at the Wellcome Centre for Integrative Neuroimaging. The IDPs include three main imaging modalities: structural, functional and diffusion brain imaging. Among these IDPs, there are very gross measures, such as the total amount of brain volume, to more detailed measurements, such as the connectivity between two brain regions. In total 819 neuroimaging IDPs were used for subsequent analysis, see B.1 for the list of IDPs used. Furthermore, we also tested our model on another set of IDPs; 150 FreeSurfer measures, which were preprocessed with [FreeSurfer](#) v6.1.0, using a $T2$ -weighted image where available, see B.1 for the list of the FreeSurfer measures used.

For the second analysis, a whole-brain model was built, using voxel-wise fractional anisotropy (FA) and mean diffusivity (MD) measures. The data were processed using the UKB pipelines; including the DTI fitting tool DTI-FIT and a tract-based spatial statistics (TBSS) style analysis, which gave us the skeletonised DTI files. In total, around 10,000 participants with dMRI-scans passed the quality control applied by the UK Biobank [17]. Afterwards, we added two extra exclusion criteria. First, participants were removed if their Z-score of the discrepancy between the dMRI image and the structural T1 image was higher than three, based on data-field 25731 in the UK Biobank. Second, participants were removed if their Z-score of the number of outlier slices was higher than three, which is a reflection of the movement of the participant during the scan, based on data-filed 25746-2.0 in the UK Biobank. For the covariates we used age, gender and dummy coded site variables.

2.2 Cognitive data

We used the cognitive phenotypes that were extracted from the UK biobank using FUNPACK [21] to evaluate the cognitive associations with the deviations from the normative model. These measures are derived from the 13 cognitive tests present in the UK Biobank, see the [UKB showcase](#). The tests were administered using a touchscreen questionnaire and included numerical memory, reaction time, fluid intelligence, visual memory and prospective memory. Later other tests that measured executive function, declarative memory and non-verbal reasoning were added [22]. For full details on the

different cognitive tests applied in UK Biobank see [23]. An overview of all the measures used in this study is presented in the supplementary E.6.

2.3 Normative model formulation

We use a flexible normative modelling framework to model different types of neuroimaging data. We have N subjects with brain data $\{y_n\}_{n=1}^N$, each of dimension D (e.g. the number of voxels or IDPs) and acquired from one of S different scanning sites. We use \mathbf{Y} to denote an $N \times D$ matrix containing these variables, where y_{nd} denotes the n -th subject and d -th neuroimaging variable. Since the neuroimaging variables are estimated separately here, we simplify the notation by using \mathbf{y} to denote the vector of observations from a single variable and y_n for a single observation. In general, we want to predict the distribution of the value for each voxel or brain region, the dependent variable (\mathbf{y}), from a set of covariates $\{\mathbf{x}_n\}_{n=1}^N$ (e.g. age, gender or site), the independent variables. In this paper, we adopt a straightforward approach to model nonlinear relationships, by applying a basis expansion to the independent variables. A common approach is to use polynomials, but these can be a poor choice, as they can induce global curvature [24]. Here we apply a common B-spline basis expansion (specifically, cubic splines with 5 evenly spaced knot points), although other approaches are also possible. We denote this expansion by $\phi(\mathbf{x})$, with Φ an $N \times K$ matrix containing the basis expansion for all subjects. In the applied model, y is assumed to be the result of a linear combination of the B-spline basis function transformation plus a noise term:

$$y = \mathbf{w}^T \phi(\mathbf{x}) + \epsilon_s \quad (1)$$

With \mathbf{w} the estimated vector of weights and $\epsilon_s = \mathcal{N}(0, \beta_s^{-1})$ a Gaussian noise distribution for site s , with mean zero and a noise precision term β_s (i.e. the inverse variance). All the noise precision terms from the different sites will be combined in a vector $\boldsymbol{\beta}$ and the site precision matrix $\Lambda_{\boldsymbol{\beta}}$ which has $\boldsymbol{\beta}$ along the leading diagonal and is the inverse of the site covariance matrix $\Sigma_{\boldsymbol{\beta}} = \Lambda_{\boldsymbol{\beta}}^{-1}$. Note that we allow the noise precision to vary across sites in order to accommodate inter-site variation along with site-specific intercepts (i.e. dummy coded site regressors in the design matrix). We have shown previously that this approach provides an efficient way to accommodate site effects in normative modelling [25].

Following similar derivations as given by Huertas et al. [16], we consider a BLR model, placing a Gaussian prior over our model parameters $p(\mathbf{w}|\boldsymbol{\alpha}) = \mathcal{N}(\mathbf{w}|0, \Lambda_{\boldsymbol{\alpha}}^{-1})$, with $\boldsymbol{\alpha}$ the hyper-parameters that the weights depend on. The Gaussian prior is assumed to have a mean zero and a precision matrix $\Lambda_{\boldsymbol{\alpha}}$, with the precision matrix the inverse of the covariance matrix $\Sigma_{\boldsymbol{\alpha}} = \Lambda_{\boldsymbol{\alpha}}^{-1}$. As shown in Huertas et al. [16], $\Lambda_{\boldsymbol{\alpha}}$ can be quite general, but here we use a simple isotropic precision matrix $\Lambda_{\boldsymbol{\alpha}} = \alpha \mathbf{I}$. The Gaussian prior choice allows us to compute the posterior distribution of \mathbf{w} in a closed form:

$$p(\mathbf{w} | \mathbf{y}, \Phi, \alpha, \boldsymbol{\beta}) = \frac{\text{likelihood} \times \text{prior}}{\text{marginal likelihood}} = \frac{\prod_n p(y_n | \Phi, \boldsymbol{\beta}, \mathbf{w}) p(\mathbf{w} | \boldsymbol{\alpha})}{p(\mathbf{y} | \Phi, \alpha, \boldsymbol{\beta})} \quad (2)$$

The posterior for each subject can then be found using the standard derivations of the posterior [26]:

$$\begin{aligned}
 p(\mathbf{w} \mid \mathbf{y}, \Phi, \alpha, \beta) &= \mathcal{N}(\mathbf{w} \mid \bar{\mathbf{w}}, \mathbf{A}^{-1}) \\
 \mathbf{A} &= \Phi^T \Lambda_\beta \Phi + \Lambda_\alpha \\
 \bar{\mathbf{w}} &= \mathbf{A}^{-1} \Phi^T \Lambda_\beta \mathbf{y}
 \end{aligned} \tag{3}$$

We use a Type II maximum likelihood approach (i.e. empirical Bayes), optimizing the denominator of the posterior to find the optimal hyper-parameters α and β . This gives an automatic trade-off between model fit and model complexity. The marginal likelihood is maximized by minimizing the negative log likelihood (NLL):

$$\begin{aligned}
 \text{NLL} &= -\log(p(\mathbf{y} \mid \alpha, \beta)) \\
 &= -\log\left(\int p(\mathbf{y} \mid \mathbf{w}, \beta) p(\mathbf{w} \mid \alpha) d\mathbf{w}\right) \\
 &= -\left(\frac{N}{2} \log|\Lambda_\beta| - \frac{ND}{2} \log 2\pi - \frac{N}{2} \log|\Lambda_\alpha| - \frac{N}{2} \log|\mathbf{A}|\right. \\
 &\quad \left. - \frac{1}{2} \sum_{n=1}^N (\mathbf{y} - \Phi \bar{\mathbf{w}})^T \Lambda_\beta (\mathbf{y} - \Phi \bar{\mathbf{w}}) - \bar{\mathbf{w}}^T \Lambda_\alpha \bar{\mathbf{w}}\right)
 \end{aligned} \tag{4}$$

The optimal hyper-parameters α and β are often estimated using a conjugate gradient optimisation of the NLL, where the derivatives can be computed directly. However, here we used Powell's method to fit the hyper-parameters. Powell's method is a derivative-free method, which in this case is faster, because computing the derivatives of the marginal likelihood with respect to the hyper-parameters is computationally very expensive. Finally, the predictive distribution is given by:

$$\hat{y} = \mathcal{N}(\bar{\mathbf{w}}^T \phi(\mathbf{x}), \phi(\mathbf{x})^T \mathbf{A}^{-1} \phi(\mathbf{x}) + \beta_s^{-1}) \tag{5}$$

2.3.1 Likelihood warping—In order to model non-Gaussian error distributions, we employed a ‘warped’ likelihood [27]. This involves applying a non-linear monotonic warping function ϕ_i to the input data during the model fit, with the index i indicating a different warping function (e.g. SinArcsinh, Box-Cox etc.). This is similar to the classical statistical technique of variable transformation, but has the advantage that the parameters of the transformation are optimised during model fitting, to provide the optimal mapping that ensures that model residuals have a Gaussian form. The warped functions are chosen such that they have a closed form inverse and are differentiable, which has several benefits: first, non-Gaussian data can be mapped (i.e. warped) exactly to better match Gaussian modelling assumptions or the predictions can be warped back to the original non-Gaussian space; second, it allows inference, prediction and computation of error measures all in closed form; finally, we can construct compositions of functions from the invertible monotonic warping functions that can greatly improve the expressivity of the model in transforming

non-Gaussian distributed data \mathbf{y} to a Gaussian form, \mathbf{z} , where inference is straightforward [28]. This is done by applying a compositional warping function φ to the observations \mathbf{y} :

$$\begin{aligned}\varphi(\cdot) &= \varphi_i(\varphi_{i-1}(\dots(\varphi_2(\varphi_1(\cdot))))\dots)) \\ \mathbf{z} &= \varphi(\mathbf{y}; \boldsymbol{\gamma})\end{aligned}\quad (6)$$

With $\boldsymbol{\gamma}$ denoting the hyper-parameter(s) of different warping functions. The warping transformation allows us to compute error measures in the warped space and to describe the deviations of subjects under a Gaussian error distribution in the form of pseudo Z statistics, even if the original data distribution is non-Gaussian.

The optimal hyper-parameters ($\boldsymbol{\alpha}$, $\boldsymbol{\beta}$ and $\boldsymbol{\gamma}$) are calculated by minimizing the warped NLL. The warped NLL can be found by accounting for the change of variables in the probability density function [28]:

$$p_{\mathbf{y}}(\mathbf{y}) = p_{\mathbf{z}}(\varphi(\mathbf{y})) |\nabla \varphi(\mathbf{y})|$$

With $\nabla \varphi(\cdot)$ the Jacobian of the transformation, which is diagonal and therefore we can simplify as a product of the individual terms:

$$p_{\mathbf{y}}(\mathbf{y}) = p_{\mathbf{z}}(\varphi(\mathbf{y})) \prod_{i=1}^n \frac{d\varphi(y_n)}{dy}$$

If we take the negative log of this equation the warped NLL will remain the same as equation 4, except for replacing the y by the transformed $\varphi(\mathbf{y})$ and the inclusion of the Jacobian term that takes the change of volume induced by the warping into account, thereby ensuring a valid probability measure, for details see [28]:

$$\begin{aligned}\text{Warped NLL} &= -\log(p(\mathbf{y} | \boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\gamma})) \\ &= \text{NLL} - \sum_{n=1}^N \log \frac{d\varphi(y_n)}{dy}\end{aligned}\quad (7)$$

2.3.2 Computational complexity—The optimization of the hyper-parameters is controlled by the minimization of the warped NLL. The warped NLL consists of the basic BLR NLL term and the log-derivatives of the warping φ_i functions, which are known in closed-form by construction. The complexity of the warped BLR model is then roughly the same as the classic BLR. However, the warped NLL is optimized for an extra hyper-parameter $\boldsymbol{\gamma}$, which could lead to the presence of more local minima, making the optimization process slightly slower [28].

2.3.3 Warped composition function—Different elementary functions can be used to create the warped composition function φ . For this paper, we test affine, Box-Cox and SinhArcsinh transformations and compositions of these transformations:

$$\begin{aligned}
\varphi_{Affine}(\mathbf{y}; \gamma) &= a + b\mathbf{y} \\
\varphi_{Box-Cox}(\mathbf{y}; \gamma) &= \frac{\text{sgn}(\mathbf{y})|\mathbf{y}|^\lambda - 1}{\lambda} \\
\varphi_{SinhArcsinh}(\mathbf{y}; \gamma) &= \sinh(b * \text{arcsinh}(\mathbf{y}) - a)
\end{aligned} \tag{8}$$

With γ the respective parameters of the different warping functions. For the SinArcsinh warping we also applied a reparametrization [29], as this empirically gave more stable results:

$$\begin{aligned}
\varphi_{SinhArcsinh}(\mathbf{y}; \gamma) &= \sinh(b * \text{arcsinh}(\mathbf{y}) + \epsilon * b) \\
a &= -\epsilon * b
\end{aligned}$$

2.4 Model selection

We evaluate the models using different model selection criteria. First, we calculate the explained variance (EV) of the model. It is expected that the gain in fit for the warped BLR will be highly dependent on the flexibility of the model. Therefore, the Bayesian Information Criterion (BIC) is also considered:

$$BIC = k * \log(N) + 2 * NLL$$

Which penalises for model complexity. Here N denotes the number of participants in the training set, NLL the negative log-likelihood. k is the number of free parameters. Note that we use the marginalized form of the NLL, which already takes into account the number of estimated coefficients. Therefore, the BIC only needs to be corrected for the added complexity of the degrees of freedom of the model (i.e. the parameters that are not integrated out). For the standard BLR this is two, one for the precision over the weights and one for the precision over the noise (α and β respectively). For the warped SinArcsinh BLR two extra degrees of freedom are added for the shape parameters (a and b). The BIC gives a good trade-off between the extra flexibility found in the warped BLR model and the better fit of the model. Finally, the mean standardized log-likelihood (MSLL) is used as a third model criterion. The MSLL takes into account the mean error and the estimated prediction variance.

2.5 Deviance scores and correlation to cognitive phenotypes

We want to find a statistical estimate of how much each participant deviates from the normal range. This is done by computing a Z-score for each subject n , also denoting explicitly the dependence on each voxel or IDP d :

$$z_{nd} = \frac{y_{nd} - \hat{y}_{nd}}{\sqrt{\sigma_d^2 + (\sigma_*^2)_d}} \tag{9}$$

With, \hat{y}_{nd} the predicted mean and y_{nd} the true response. Normalized by $\sigma_d^2 = (\beta_s^{-1})_d$ the estimated noise variance (i.e. reflecting variation in the data) and $(\sigma_*^2)_d = \phi(\mathbf{x})^T \mathbf{A}_d^{-1} \phi(\mathbf{x})$ the variance attributable to modelling uncertainty for the d -th voxel. For the warped statistic, we compute the Z-scores in the warped (i.e. Gaussian) space. The true response variables are warped to the Gaussian space to ensure the underlying assumption of normality is satisfied by the construction of the warping functions.

Afterwards, to ensure our model can also be applied for behavioural and clinical estimations, we look at the correlations between the Z-scores from the IDPs and the whole brain analysis, and the cognitive scores of the UK Biobank. For the IDPs, we directly correlate the Z-scores and the cognitive phenotypes through a Spearman correlation. For the whole-brain analysis, we first make a summary statistic of the Z-scores by calculating the extreme value distribution. We model the extreme value distribution by looking at the mean of the top 1% of the deviations across the whole brain [10]. The extreme value statistics give the largest deviations per subject from the normal pattern, which have shown to be strongly correlated to behaviour [10], [30]. Afterwards, we apply a principal component analysis (PCA) on the cognitive phenotypes to give a one-factor solution. This first component has been shown to be correlated to the ‘general’ factor of cognitive ability or the ‘g-factor’ [31]. Lastly, we compute the Spearman coefficient between the first principal component and the summary deviation score.

3 Results

3.1 Performance of the warped Bayesian linear regression model for IDPs

All the statistical analyses were performed in Python version 3.8, using the [PCNtoolkit](#). The BLR algorithm from the PCNtoolkit was chosen for all experiments. We considered age, binary gender and binary site ID within the covariance matrix. We used a standard BLR or we transformed the age covariate with a B-spline of order three with three knots. The Powell method was selected for the optimizer. We randomly split the dataset into 50% training and 50% test and reported all the error metrics on the test set. In the PCNtoolbox, several warpings can be chosen depending on the imaging modality one wants to model. We tested several warping functions (affine, Box-Cox and SinhArcsinh) and compositions of these warping functions. Preliminary testing showed that the SinhArcsinh warping gave the best fit compared to the alternatives evaluated. Therefore, in this paper, only the results of the SinhArcsinh warping are presented.

In figure 1, Bland-Altman plots are shown comparing the standard BLR and the B-spline BLR. The figure presents different model selection criteria: MSLL and BIC (EV can be seen in supplement figure A.8). The plots demonstrate that for most IDPs a non-linear B-spline BLR model performs better than a standard BLR. Indicating that non-linearity is a key component that should be accounted for in modelling neuroimaging data.

In figure 2, Bland-Altman plots are shown that compare the B-spline BLR and the warped BLR models for all IDPs, using the MSLL and BIC (EV can be seen in supplement figure A.8). We also plotted the difference in absolute values of the skewness and kurtosis.

In figure 3, the same plots are shown for the FreeSurfer measures. We included them separately, as they were preprocessed separately (i.e. we did not use the IDPs provided by UK Biobank and instead ran the FreeSurfer reconstructions manually). The plots show that for specific IDPs the warped BLR performs better than the B-spline BLR. When we examined these IDPs more closely, it was noted that they demonstrated distinct non-Gaussian behaviour. An example of such behaviour is given down below with the WMHs (white matter hyper-intensities). In the supplementary table C.3, we provide a summary of some of the results for different IDPs that can help inform which neuroimaging modalities are best modelled with the warped BLR. For an indication of the effect sizes of the model selection criteria for the different model settings, see supplementary tables D.4 and D.5. Note also that the MSLL and EV do not clearly reflect differences in the shape of the predictive distribution. For example, for the IDPs, there is no average difference between the warped and non-warped model (Fig. 2 panel A and supp. fig. A.8 panel B), yet the warped model consistently yields a predictive distribution –and resultant Z-score distribution– that is less (or equivalently) skewed and kurtotic (Fig. 2 panels C and D).

In figure 4 and 5, we show the results of an illustrative analysis predicting WMH load across ageing to demonstrate how the performance of the warped BLR model compares to a B-spline BLR. The figures show the B-spline BLR and warped BLR results for WMHs at one-time point and the longitudinal data of two-time points. The results demonstrate that (i) the non-linearity of the data is sufficiently captured with a B-spline transformed BLR (ii) the WMHs show a distinctly non-Gaussian variance pattern, which is better predicted by the warped BLR. Thus, indicating that if the data has a non-Gaussian distribution for the residuals a warped BLR is preferred over a B-spline BLR.

3.1.1 Correlation deviance scores WMHs and cognitive phenotypes—We also wanted to correlate the warped BLR model output of the WMHs to behavioural variables to ensure that the model can be used for behavioural predictions. We loaded all cognitive phenotypes available in UK Biobank according to the FUNPACK categorization, including: reaction time, numeric memory, prospective memory etc. (for a full list of the cognitive phenotypes used, see the supplementary table E.6). We calculated the deviance Z-scores according to formula 9. Afterwards, we calculated the Spearman correlation between the cognitive phenotypes and the Z-scores. Numeric memory (ID: 4259, ‘Digits entered correctly’) was modestly but significantly correlated with the warped Z-scores: $\rho = -0.0331$, $p = 0.0262$. In other words, if a participant’s WMH deviation from normal development increases the number of correctly remembered digits drops.

Lastly, to illustrate the value of normative models in a longitudinal context, we tested for an association between change in WMHs and change in cognitive phenotypes of the longitudinal data to see if WMH load is correlated to cognitive decline. We performed a statistical Wilcoxon rank-sum test on the participants’ cognitive phenotypes contrasting subjects that have a difference in the Z-scores > 0.5 , which corresponds to a difference in half a standard deviation, versus the participants that do not. Intuitively, this contrasts individuals who are following an expected trajectory of ageing with those who deviate from such a trajectory. Highly significant associations were found with the reaction time (ID: 404, ‘Duration to first press of snap-button in each round’) $W = 5.5641$, $p < 0.001$ and with

the Trail Making Test (ID: 6771, 'Errors before selecting correct item in alphanumeric path (trail #2)') $W = 8.3105$, $p < 0.001$. The results show an association between the change in cognition and the change in WMH deviance scores.

3.2 Scalability to a whole brain voxelwise based analysis

For the follow-up analysis, we evaluated the warped BLR approach on a whole-brain level for two DTI imaging modalities (FA and MD). The results of these two modalities were very similar and therefore we will only present the results for FA here. We separated the entire dataset into 80% training data and 20% testing data. First, we computed the time complexity per model fit (e.g. for one voxel) with varying number of subjects using the B-spline BLR model setting and compared it to the Gaussian process regression setting (Figure 6). This demonstrates the clear computational advantage of the BLR setting for the whole brain analysis.

Afterwards, we tested different model settings for the imaging modalities including a standard BLR, B-spline BLR and a SinhArcsinh warped BLR. Figure 7 shows the comparative results in a Bland-Altman plot for the FA dataset (which were similar for the MD dataset). The figure presents the EV, MSLL and the BIC for the B-spline BLR and the warped BLR. These results are consistent with the IDPs in that according to the EV and MSLL, the models perform quite similarly for most voxels. Although, we would argue that these measures are not necessarily sensitive for the added benefit of the warping of the likelihood, which will mostly affect the predictions in the outer centiles. For the BIC the results demonstrate that the warped BLR is preferred for certain voxels. The voxels where a warped model is favoured generally showed more non-Gaussian behaviour.

Finally, We used a paired-sample t-test, pairing the whole brain results (EV, MSLL and BIC) of the different models to estimate the difference between performance measures of the warped and non-warped BLR. For MD the following effect sizes were found: *EV*: $d = 0.33$, *MSLL*: $d = 0.003$ and *BIC*: $d = -0.79$. For FA the following effect sizes were found: *EV*: $d = 0.028$, *MSLL*: $d = 0.017$ and *BIC*: $d = 0.55$. We can see that the difference between the methods is small. Indicating that the B-spline BLR and the warped BLR model are quite similar in their model fit for MD and FA.

3.2.1 Correlation deviance scores DTI and cognitive phenotypes—Finally, we correlated the Z-scores of the whole brain warped BLR model for the MD dataset to the cognitive phenotypes. First, we scaled the cognitive data and performed a principal component analysis. We selected the first component, which explained 29% of the variance in the data. Afterwards, we made a summary score of the Z-scores for each participant by looking at the largest deviations, which in the limit should follow an extreme value distribution [32]. We fitted a generalized extreme value distribution to the top 1% of the absolute Z-scores of each subject. Subsequently, we computed a Spearman correlation between the extreme values and the first principal component of the cognitive phenotypes, which gave $\rho = 0.158$, $p < 0.001$. The results demonstrate a clear correlation between the warped deviations from normal development and the cognitive phenotypes. This relationship will be explored further in future studies.

4 Discussion

In this paper, we presented a next-generation framework to scale normative models for large population-sized datasets based on warped Bayesian linear regression (BLR). Normative models can capture the heterogeneity in the population and model individual deviations from normal brain development. We demonstrated that the shift in normative modelling to a B-spline BLR with a likelihood warping gives several benefits. In this study we showed that: (i) Compared to Gaussian process regression, it is computationally much less demanding and is therefore scalable to big datasets. (ii) The non-linearity of the model, incorporated by the B-spline, improves the fit and out of sample predictions for most variables. (iii) Non-Gaussianity of the data can be naturally included due to the incorporation of the likelihood warping in the algorithm, which allows for a wider range of datasets to be accurately modelled. (iv) Model selection criteria based on the marginal likelihood, such as the BIC, can be calculated in closed form and therefore a trade-off between model fit and model complexity can be chosen optimally from the training data, without cross-validation. (v) The deviations scores from normal brain development can be meaningfully related to behaviour. Furthermore, we demonstrated the use of the normative model with the warped BLR on different datasets from the UK Biobank, including image-derived phenotypes (IDPs); focusing on white matter hyperintensities (WMHs) as an example of non-Gaussianity and a diffusion tensor imaging (DTI) modality for a whole-brain model.

Our proposed method makes it possible to apply normative modelling to considerably larger samples than was feasible before [7], [8]. The results from the computational experiments on the whole brain model showed that the BLR method is scalable to population-sized data sets and fine-grained voxel-level data. In comparison, most normative models used Gaussian process regression, which due to its high computational complexity could only be used in studies with a relatively low sample size. This improvement is mainly because the approximation of the covariance matrix by a set of basis functions allowed us to account for non-linearity in a computationally less demanding way than the Gaussian process regression method, therefore making the B-spline BLR scalable for big datasets. Computationally scalable modelling of nonlinear effects is important since our experiments showed that a cubic B-spline transformation of the age covariate improved model fit compared to linear models for most neuroimaging modalities.

Another major benefit of our method is the possibility of modelling non-Gaussian distribution by the use of the likelihood warping technique. This is important in general, as the aim of normative modelling is to accurately model the centiles of variation in addition to modelling the mean and is especially important for normative modelling of variables that are not approximately Gaussian distributed. For example, we showed that the WMHs show non-Gaussian behaviour that is well suited to uncover the benefits of the warped model over the standard model. We demonstrated the improved fit of the WMHs by including a B-spline transformation and a SinhArcsinh likelihood warping in the normative model, which was also exemplified for the longitudinal data. The same improvement in fit for other data modalities that showed more non-Gaussianity in their residuals was also demonstrated by comparing the warped BLR to the B-spline BLR for all the IDPs. Furthermore, it was

shown on a whole-brain model of a DTI modality that for several voxels the warped BLR gives a better model performance than a B-spline BLR.

We emphasize that the addition of non-linear effects and non-gaussianity makes the model more flexible which increase the need for model selection in order to avoid possible overfitting. We presented several model selection criteria that can be used to choose the optimal model settings for different neuroimaging modalities. It should be recognized that for some IDPs and voxels the B-spline BLR gives a better fit, showing that a more flexible model is not always needed. Therefore, we recommend carefully examining the type of data one wants to model and based on the data trends found for the residuals (Gaussian or non-Gaussian) to decide if a more flexible model is preferred. This can easily be checked by looking at the skewness and kurtosis of the distribution or making a QQ-plot. Additionally, different model selection criteria can sometimes contradict each other, as they are sensitive to different parts of the data. As we showed above, classical metrics such as EV and MSL are not very sensitive to the shape of the predictive distribution. The consequence is that per task, we have to decide if we want a better EV, most sensitive to the mean fit and dependent on the flexibility of the model, or a better MSL/BIC, which is more sensitive to the variance and penalizes the flexibility of the model. The variability in model selection criteria demonstrates that for different imaging modalities, different normative modelling settings are preferred and the added flexibility is confirmed to only give an advantage for response variables that show non-Gaussianity in their residuals.

We confirmed that the deviations from the normative modelling framework can be meaningfully related to behaviour. We established a significant correlation between the warped deviance scores from the IDPs and several dimensions of the intelligence phenotype. These tests give a first indication of the possible relationships between the deviations and behaviour. For the whole brain model, the relationship with behaviour was shown with a significant correlation between an approximation to the g-factor in the form of the first principal component of the cognitive phenotypes and the warped deviance scores. This study demonstrates that the model could be extended to make predictive scores not only in the brain domain, but also for the behavioural phenotype. In the future, the neurobiological markers of deviation from normal development can be extended to become markers of psychiatric disorders. This has already been done on a smaller scale, using normative modelling [9], [10], [13], [30], [33], [34], but we would like to extend these studies to bigger data models, which include a wide variety of neuroimaging data modalities.

In conclusion, the current study suggests that non-linearity and non-Gaussianity are two parameters of big neuroimaging datasets that need to be captured to make accurate predictions for normal brain development. In this paper, we have done that through a warped BLR normative model. We have shown using several neuroimaging modalities the benefit of this model over more conservative models. Caution is essential when applying non-Gaussian models, as they can overfit and should mainly be used in the presence of non-normally distributed residuals. We recommend carefully assessing the distribution of residuals and the model selection parameters using the different model selection criteria mentioned in this paper that give a balance between model complexity and model fit.

Supplementary Material

Refer to Web version on PubMed Central for supplementary material.

References

- [1]. Sudlow C, Gallacher J, Allen N, Beral V, Burton P, Danesh J, Downey P, Elliott P, Green J, Landray M, et al. Uk biobank: an open access resource for identifying the causes of a wide range of complex diseases of middle and old age. *Plos med.* 2015; 12 e1001779 doi: 10.1371/journal.pmed.1001779 [PubMed: 25826379]
- [2]. Thompson PM, Stein JL, Medland SE, Hibar DP, Vasquez AA, Renteria ME, Toro R, Jahanshad N, Schumann G, Franke B, et al. The enigma consortium: large-scale collaborative analyses of neuroimaging and genetic data. *Brain imaging and behavior.* 2014; 8: 153–182. DOI: 10.1007/s11682-013-9269-5 [PubMed: 24399358]
- [3]. Casey B, Cannonier T, Conley MI, Cohen AO, Barch DM, Heitzeg MM, Soules ME, Teslovich T, Dellarco DV, Garavan H, et al. The adolescent brain cognitive development (abcd) study: imaging acquisition across 21 sites. *Developmental cognitive neuroscience.* 2018; 32: 43–54. DOI: 10.1016/j.dcn.2018.03.001 [PubMed: 29567376]
- [4]. Satterthwaite TD, Connolly JJ, Ruparel K, Calkins ME, Jackson C, Elliott MA, Roalf DR, Hopson R, Prabhakaran K, Behr M, et al. The philadelphia neurodevelopmental cohort: A publicly available resource for the study of normal and abnormal brain development in youth. *Neuroimage.* 2016; 124: 1115–1119. DOI: 10.1016/j.neuroimage.2015.03.056 [PubMed: 25840117]
- [5]. Insel TR, Cuthbert BN. Brain disorders? Precisely: Precision medicine comes to psychiatry. *Science.* 2015; 348: 499–500. DOI: 10.1126/science.aab2358 [PubMed: 25931539]
- [6]. Wolfers, T, Buitelaar, JK, Beckmann, CF, Franke, B, Marquand, AF. From estimating activation locality to predicting disorder: A review of pattern recognition for neuroimaging-based psychiatric diagnostics. 2015.
- [7]. Marquand, AF, Kia, SM, Zabihi, M, Wolfers, T, Buitelaar, JK, Beckmann, CF. Conceptualizing mental disorders as deviations from normative functioning. 2019.
- [8]. Marquand AF, Rezek I, Buitelaar J, Beckmann CF. Understanding Heterogeneity in Clinical Cohorts Using Normative Models: Beyond Case-Control Studies. *Biological Psychiatry.* 2016a; 80: 552–561. DOI: 10.1016/j.biopsych.2015.12.023 [PubMed: 26927419]
- [9]. Wolfers T, Doan NT, Kaufmann T, Alnæs D, Moberget T, Agartz I, Buitelaar JK, Ueland T, Melle I, Franke B, Andreassen OA, et al. Mapping the Heterogeneous Phenotype of Schizophrenia and Bipolar Disorder Using Normative Models. *JAMA Psychiatry.* 2018; 75: 1146–1155. DOI: 10.1001/jamapsychiatry.2018.2467 [PubMed: 30304337]
- [10]. Zabihi M, Oldehinkel M, Wolfers T, Frouin V, Goyard D, Loth E, Charman T, Tillmann J, Banaschewski T, Dumas G, Holt R, et al. Dissecting the Heterogeneous Cortical Anatomy of Autism Spectrum Disorder Using Normative Models. *Biological Psychiatry: Cognitive Neuroscience and Neuroimaging.* 2019; 4: 567–578. DOI: 10.1016/j.bpsc.2018.11.013 [PubMed: 30799285]
- [11]. Kaufmann T, van der Meer D, Doan NT, Schwarz E, Lund MJ, Agartz I, Alnæs D, Barch DM, Baur-Streubel R, Bertolino A, Bettella F, et al. Common brain disorders are associated with heritable patterns of apparent aging of the brain. *Nature Neuroscience.* 2019; 22: 1617–1623. DOI: 10.1038/s41593-019-0471-7 [PubMed: 31551603]
- [12]. Marquand, AF, Wolfers, T, Mennes, M, Buitelaar, J, Beckmann, CF. Beyond Lumping and Splitting: A Review of Computational Approaches for Stratifying Psychiatric Disorders. 2016b.
- [13]. Lv J, Di Biase M, Cash RF, Cocchi L, Croyley V, Klauser P, Tian Y, Bayer J, Schmaal L, Cetin-Karayumak S, et al. Individual deviations from normative models of brain structure in a large cross-sectional schizophrenia cohort. *bioRxiv.* 2020; doi: 10.1038/s41380-020-00882-5
- [14]. Kia SM, Marquand A. Normative Modeling of Neuroimaging Data using Scalable Multi-Task Gaussian Processes, *Lecture Notes in Computer Science (including subseries Lecture Notes in*

Artificial Intelligence and Lecture Notes in Bioinformatics) 11072 LNCS. 2018. 127–135. URL: <http://arxiv.org/abs/1806.01047> arXiv:1806.01047

- [15]. Rasmussen, CE, Williams, CK. Approximation methods for large datasets. 2005.
- [16]. Huertas I, Oldehinkel M, van Oort ES, Garcia-Solis D, Mir P, Beckmann CF, Marquand AF. A Bayesian spatial model for neuroimaging data based on biologically informed basis functions. *NeuroImage*. 2017; 161: 134–148. DOI: 10.1016/j.neuroimage.2017.08.009 [PubMed: 28782681]
- [17]. Alfaro-Almagro F, Jenkinson M, Bangerter NK, Andersson JL, Griffanti L, Douaud G, Sotiropoulos SN, Jbabdi S, Hernandez-Fernandez M, Vallee E, et al. Image processing and quality control for the first 10,000 brain imaging datasets from uk biobank. *Neuroimage*. 2018; 166: 400–424. DOI: 10.1016/j.neuroimage.2017.10.034 [PubMed: 29079522]
- [18]. Habes M, Pomponio R, Shou H, Doshi J, Mamourian E, Erus G, Nasrallah I, Launer LJ, Rashid T, Bilgel M, et al. The brain chart of aging: Machine-learning analytics reveals links between brain aging, white matter disease, amyloid burden, and cognition in the istaging consortium of 10,216 harmonized mr scans. *Alzheimer's & Dementia*. 2020; doi: 10.1002/alz.12178
- [19]. Cox SR, Ritchie SJ, Tucker-Drob EM, Liewald DC, Hagenaars SP, Davies G, Wardlaw JM, Gale CR, Bastin ME, Deary IJ. Ageing and brain white matter structure in 3,513 uk biobank participants. *Nature communications*. 2016; 7: 1–13. DOI: 10.1038/ncomms13629
- [20]. Miller KL, Alfaro-Almagro F, Bangerter NK, Thomas DL, Yacoub E, Xu J, Bartsch AJ, Jbabdi S, Sotiropoulos SN, Andersson JL, Griffanti L, et al. Multimodal population brain imaging in the UK Biobank prospective epidemiological study. *Nature Neuroscience*. 2016; 19: 1523–1536. DOI: 10.1038/nn.4393 [PubMed: 27643430]
- [21]. McCarthy, P. funpack. 2020.
- [22]. Fawns-Ritchie C, Deary IJ. Reliability and validity of the UK Biobank cognitive tests. *PLoS ONE*. 2020; 15 doi: 10.1371/journal.pone.0231627
- [23]. Lyall DM, Cullen B, Allerhand M, Smith DJ, Mackay D, Evans J, Anderson J, Fawns-Ritchie C, McIntosh AM, Deary IJ, Pell JP. Cognitive Test Scores in UK Biobank: Data Reduction in 480,416 Participants and Longitudinal Stability in 20,346 Participants. *PLOS ONE*. 2016; 11 e0154222 doi: 10.1371/journal.pone.0154222 [PubMed: 27110937]
- [24]. Fjell AM, Walhovd KB, Westlye LT, Østby Y, Tamnes CK, Jernigan TL, Gamst A, Dale AM. When does brain aging accelerate? dangers of quadratic fits in cross-sectional studies. *Neuroimage*. 2010; 50: 1376–1383. DOI: 10.1016/j.neuroimage.2010.01.061 [PubMed: 20109562]
- [25]. Kia SM, Huijsdens H, Dinga R, Wolfers T, Mennes M, Andreassen OA, Westlye LT, Beckmann CF, Marquand AF. Hierarchical bayesian regression for multi-site normative modeling of neuroimaging data. arXiv preprint. 2020. arXiv:2005.12055
- [26]. Bishop, CM. Pattern recognition and machine learning. springer; 2006.
- [27]. Snelson E, Ghahramani Z, Rasmussen CE. Warped gaussian processes. *Advances in neural information processing systems*. 2004. 337–344.
- [28]. Rios G, Tobar F. Compositionally-warped gaussian processes. *Neural Networks*. 2019; 118: 235–246. DOI: 10.1016/j.neunet.2019.06.012 [PubMed: 31319321]
- [29]. Jones MC, Pewsey A. Sinh-arcsinh distributions. *Biometrika*. 2009; 96: 761–780. DOI: 10.1093/biomet/asp053
- [30]. Marquand AF, Rezek I, Buitelaar J, Beckmann CF. Understanding heterogeneity in clinical cohorts using normative models: beyond case-control studies. *Biological psychiatry*. 2016; 80: 552–561. DOI: 10.1016/j.biopsych.2015.12.023 [PubMed: 26927419]
- [31]. Nave G, Jung WH, Karlsson Linnér R, Kable JW, Koellinger PD. Are Bigger Brains Smarter? Evidence From a Large-Scale Preregistered Study. *Psychological Science*. 2019; 30: 43–54. DOI: 10.1177/0956797618808470 [PubMed: 30499747]
- [32]. Fisher, RA, Tippett, LHC. *Mathematical Proceedings of the Cambridge Philosophical Society*. Vol. 24. Cambridge University Press; 1928. 180–190.
- [33]. Wolfers T, Beckmann CF, Hoogman M, Buitelaar JK, Franke B, Marquand AF. Individual differences v. the average patient: Mapping the heterogeneity in ADHD using normative models.

Psychological Medicine. 2019; 50: 314–323. DOI: 10.1017/S0033291719000084 [PubMed: 30782224]

- [34]. Zabihi M, Floris DL, Kia SM, Wolfers T, Tillmann J, Arenas AL, Moessnang C, Banaschewski T, Holt R, Baron-Cohen S, Loth E, et al. Fractionating autism based on neuroanatomical normative modeling. *Translational Psychiatry*. 2020; 10: 1–10. DOI: 10.1038/s41398-020-01057-0 [PubMed: 32066695]

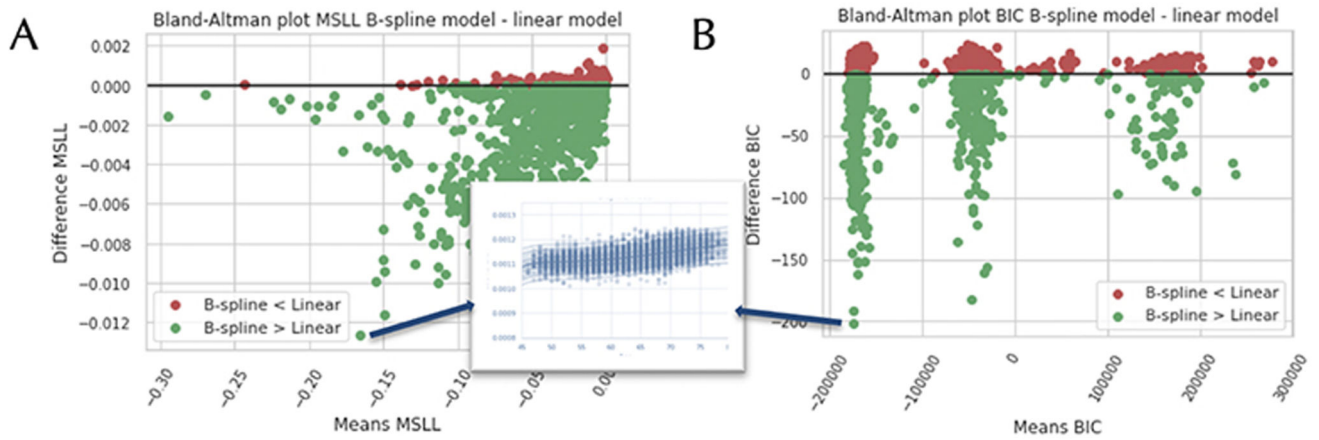


Figure 1. Bland-Altman plots comparing the standard and B-spline Bayesian Linear Regression (BLR) models, using Image-Derived Phenotypes (IDPs).

Each dot indicates one IDP. The models are compared according to the following model selection criteria: the Mean Standardized Log Loss (MSLL) (A) and the Bayesian Information Criteria (BIC) (B). The green colour indicates a better fit for the non-linear B-spline model compared to the linear model. We also plotted a zoomed-in view of the model fit for one of the IDPs.

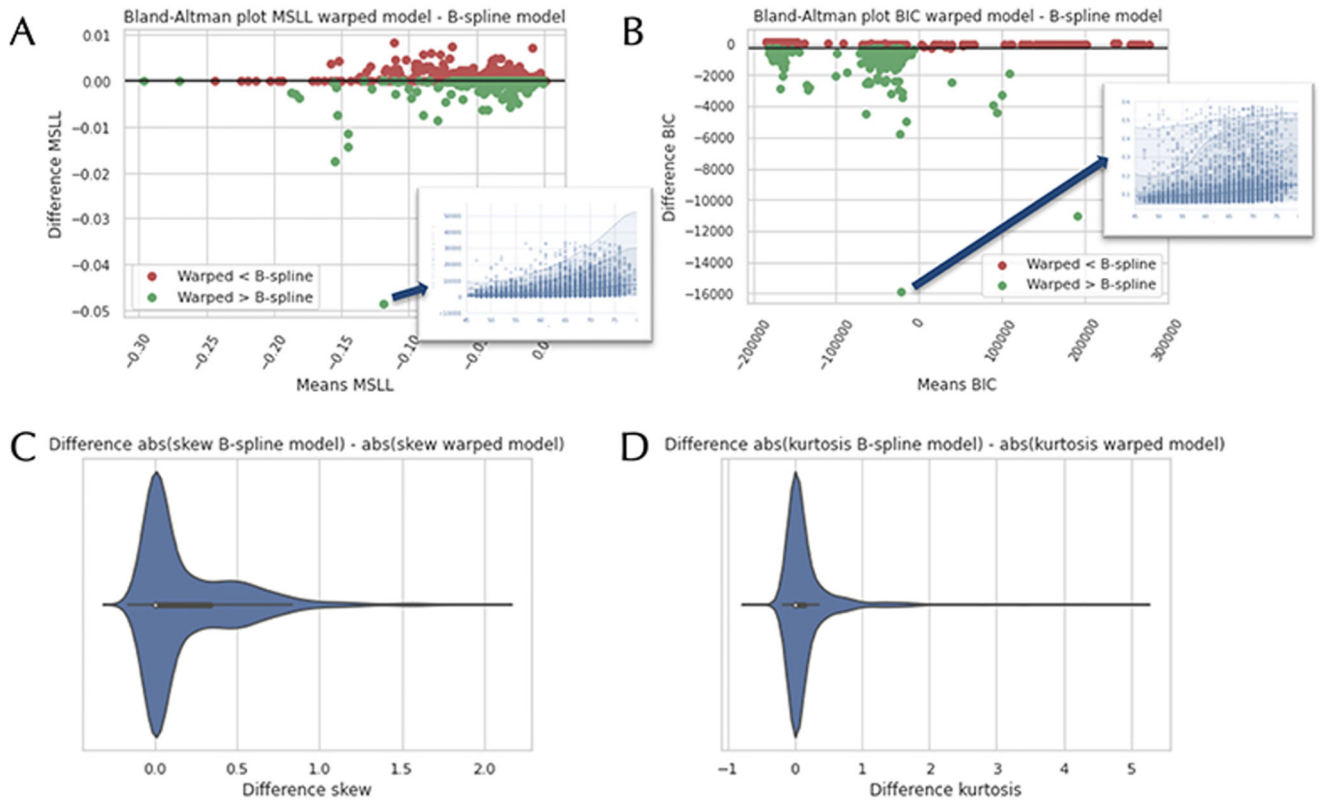


Figure 2. Bland-Altman plots comparing the B-spline and warped Bayesian Linear Regression (BLR) models, using Image-Derived Phenotypes (IDPs).

The models are compared according to the following model selection criteria: the Mean Standardized Log Loss (MSLL) (A) and the Bayesian Information Criteria (BIC) (B). The green colour indicates a better fit for the warped model compared to the B-spline model. We also plotted a zoomed-in view of the model fit for two of the IDPs. On images C and D, we show the difference in absolute values of the skewness and kurtosis between the B-spline and warped model. A more positive value indicates that the B-spline model had a higher skewness or kurtosis than the warped model.

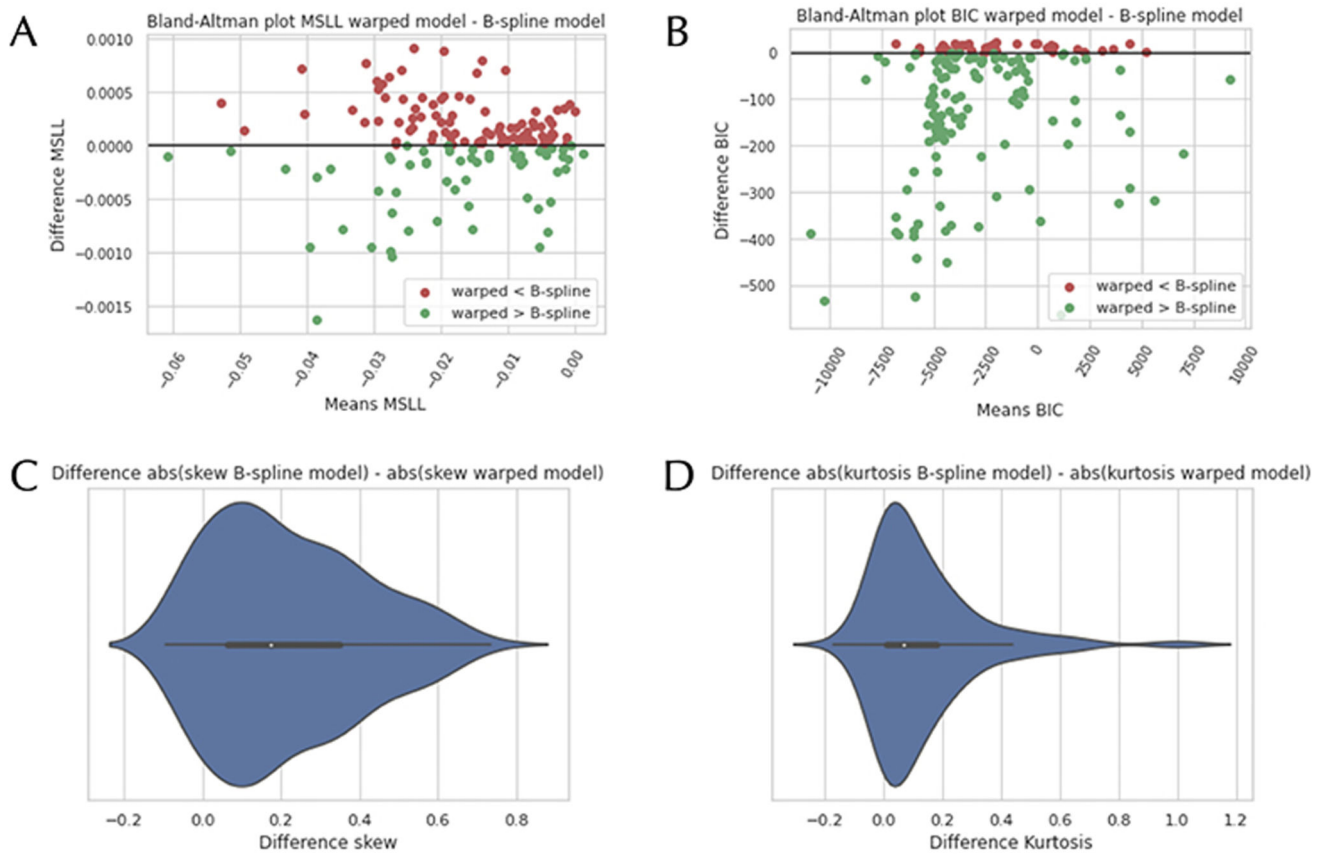


Figure 3. Bland-Altman plots comparing the B-spline and warped Bayesian Linear Regression (BLR) models, using the FreeSurfer measurements.

The models are compared according to the following model selection criteria: the Mean Standardized Log Loss (MSLL) (A) and the Bayesian Information Criteria (BIC) (B). We also plotted a zoomed-in view of the model fit for one of the IDPs. On images C and D, we show the difference in absolute values of the skewness and kurtosis between the B-spline and warped model. A more positive number means a better fit for the warped model compared to the B-spline model.

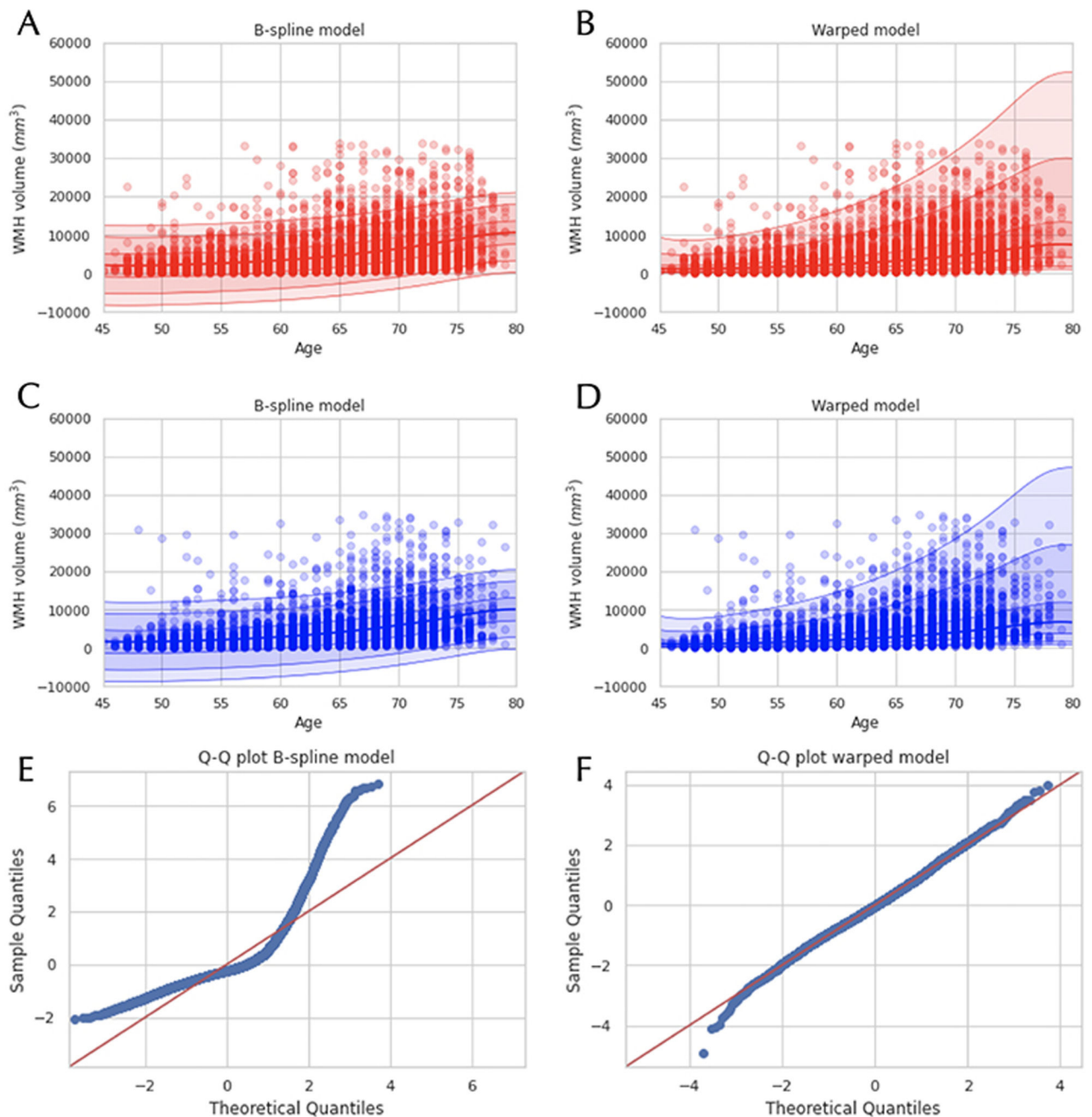


Figure 4. White matter hyperintensities (WMHs) modelled as a function of age using a Bayesian Linear Regression (BLR) model.

Images A and C demonstrate the model fit using a regular Gaussian B-spline BLR, for the female and male cohorts respectively, both visualizing the mean prediction and the centiles of variation for the WMHs. Images B and D show comparable fits for a SinArcsinh warped BLR, for the female and male cohorts respectively. In images E and F quantile-quantile (QQ) plots of the two models are shown, demonstrating a better fit for the data using a warped BLR model.

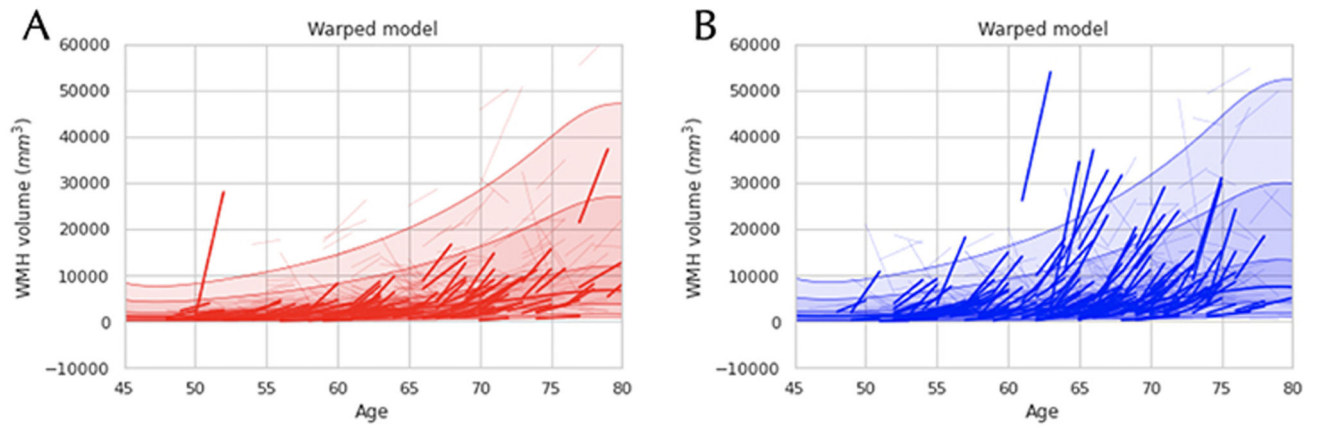


Figure 5. Here the longitudinal follow-up data of the WMHs is plotted for females (A) and males (B), using a SinhArcsinh warped BLR model.

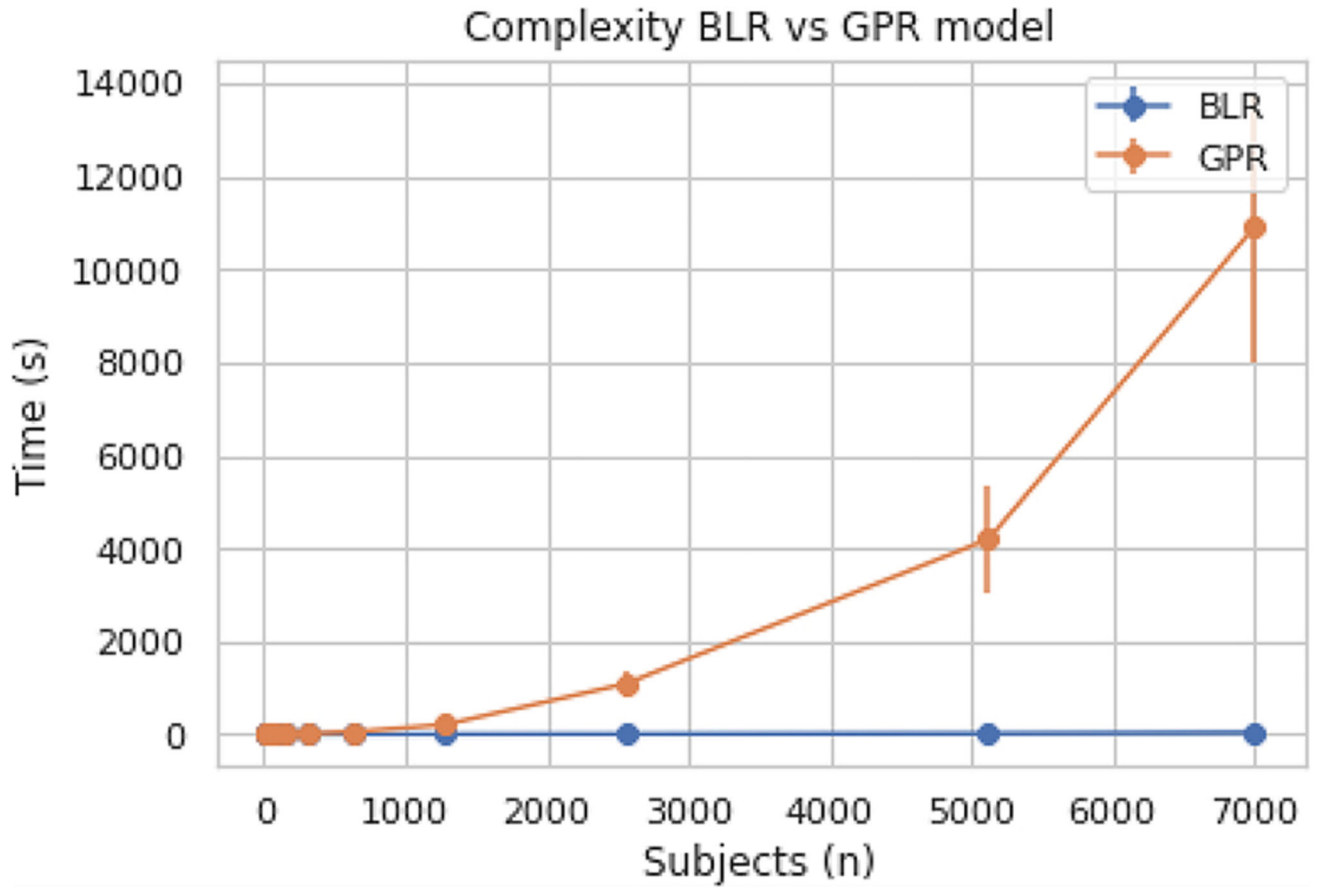


Figure 6. Computational complexity comparison between the Bayesian linear regression (BLR) model setting and the Gaussian process regression (GPR) model setting, giving the mean and the standard error (SE) over ten runs.

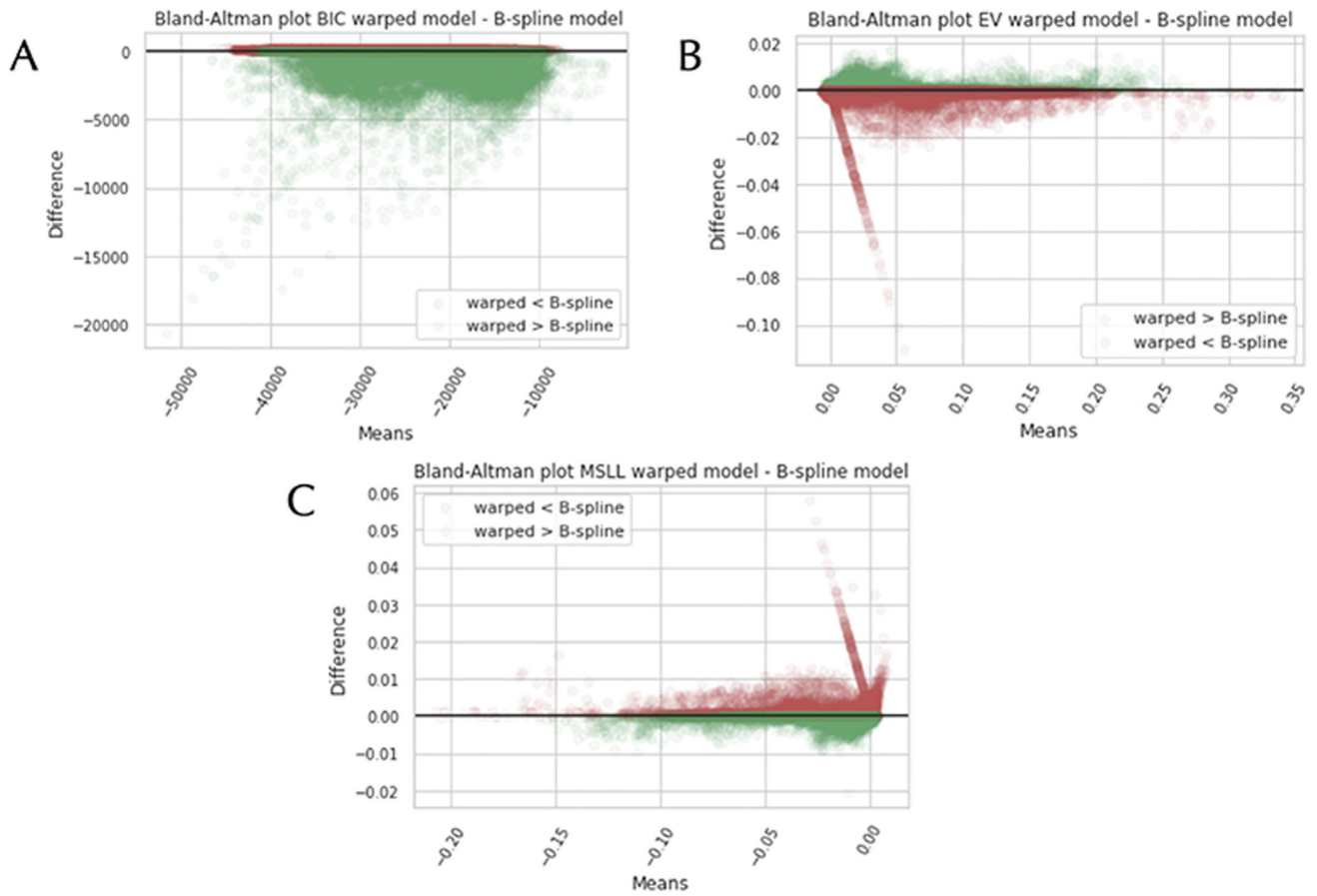


Figure 7. Bland-Altman plots comparing the warped Bayesian Linear Regression (BLR) model to the B-spline BLR model, using Fractional Anisotropy (FA) data.

The comparison is done according to the following model selection criteria: The Bayesian Information Criteria (BIC) (A), the Explained Variance (EV) (B), and the Mean Standardized Log Loss (MSLL) (C). The green colour indicates a better fit for the warped BLR.