



Title	A dynamic probabilistic principal components model for the analysis of longitudinal metabolomics data
Authors(s)	Nyamundanda, Gift, Gormley, Isobel Claire, Brennan, Lorraine
Publication date	2014-11
Publication information	Nyamundanda, Gift, Isobel Claire Gormley, and Lorraine Brennan. "A Dynamic Probabilistic Principal Components Model for the Analysis of Longitudinal Metabolomics Data." Wiley, November 2014. https://doi.org/10.1111/rssc.12060 .
Publisher	Wiley
Item record/more information	http://hdl.handle.net/10197/7107
Publisher's statement	This is the author's version of the following article: Gift Nyamundanda, Isobel Claire Gormley and Lorraine Brennan (2014) "A dynamic probabilistic principal components model for the analysis of longitudinal metabolomics data" Journal of the Royal Statistical Society: Series C (Applied Statistics), 63(5): 763-782 which has been published in final form at http://dx.doi.org/10.1111/rssc.12060 .
Publisher's version (DOI)	10.1111/rssc.12060

Downloaded 2026-05-01 23:38:14

The UCD community has made this article openly available. Please share how this access benefits you. Your story matters! (@ucd_oa)



© Some rights reserved. For more information

A dynamic probabilistic principal components model for the analysis of longitudinal metabolomics data.

Gift Nyamundanda

School of Mathematical Sciences, University College Dublin, Ireland.

Isobel Claire Gormley

School of Mathematical Sciences, University College Dublin, Ireland.

Lorraine Brennan

School of Agriculture and Food Science, Conway Institute, University College Dublin, Ireland.

1 **Summary.** In a longitudinal metabolomics study, multiple metabolites are measured from
2 several observations at many time points. Interest lies in reducing the dimensionality of such
3 data and in highlighting influential metabolites which change over time. A dynamic probabilistic
4 principal components analysis (DPPCA) model is proposed to achieve dimension reduction
5 while appropriately modelling the correlation due to repeated measurements. This is achieved
6 by assuming an autoregressive model for some of the model parameters. Linear mixed models
7 are subsequently used to identify influential metabolites which change over time. The proposed
8 model is used to analyse data from a longitudinal metabolomics animal study.

9 1. Introduction

10 Metabolomics is the study of low molecular weight compounds known as metabolites found
11 in biological samples; its application reveals information on metabolic pathways within an
12 organism. The number of areas in which metabolomics is applied has recently enjoyed rapid
13 growth and metabolomics is now employed in fields such as nutrition, toxicology and disease
14 diagnosis. In a typical metabolomics study large data sets are generated using analytical
15 technologies such as nuclear magnetic resonance spectroscopy (NMR) (Reo, 2002) and mass
16 spectrometry (MS) (Dettmer et al., 2007). With respect to NMR spectroscopy the resulting
17 spectrum consists of a series of peaks where the height of a peak is related to the relative
18 abundance of the associated metabolite. Studying such metabolomic profiles gives insight
19 to the metabolic state of a system.

20
21 Metabolomic data sets are usually high-dimensional, in that the resulting spectra con-
22 tain many peaks (i.e. variables), yet they are characterised by small sample sizes – hence
23 classical statistical approaches cannot be easily applied. The data sets contain variables
24 that are not independent in that metabolites can be represented by more than one peak
25 and metabolites can be highly correlated (van den Berg et al., 2006). In addition to corre-
26 lated variables, in longitudinal metabolomics data sets there is further correlation structure
27 due to the repeated measurements of observations over time. Hence, appropriate statistical
28 models are required in order to appropriately model the data and extract true, important
29 information.

30
31 Within the metabolomics literature, principal components analysis (PCA) (Jolliffe, 2002)
32 is often used for multivariate data exploration (Walsh et al., 2007; Smolinska et al., 2012;

33 Cassol et al., 2013; Carvalho et al., 2013; Bathen et al., 2013; Sachse et al., 2012). Methods
34 that improve and extend the application of this common statistical technique will prove
35 extremely useful to the metabolomics practitioner, and to scientists in other fields. The
36 application of PCA to longitudinal studies is limited however by the fact that PCA does
37 not take into account information about the experimental design i.e. if PCA is applied to
38 all time points simultaneously, measurements taken repeatedly over time are assumed inde-
39 pendent (Choi et al., 2006). In such a case, since PCA looks for directions in the data space
40 with maximum variation, time related variation will act as a confounding factor obscuring
41 potential differences due to treatment.

42
43 Several extensions to PCA have been developed to take into account the experimental
44 design of a study and therefore can be used to analyse longitudinal metabolomics data more
45 appropriately. These include weighted PCA (Jansen et al., 2004) which uses weights to ac-
46 count for variation due to repeated measurements and ASCA (Smilde et al., 2005) which
47 combines analysis of variance and simultaneous components analysis methods to deal with
48 complex multivariate datasets. Jansen et al. (2009) employ local PCA models at each time
49 point, and then link these local models to each other. Dynamic PCA (Smilde et al., 2010)
50 uses a back-shift matrix to analyse data from multiple time points simultaneously. The
51 main limitation of these approaches is that they do not have an associated generative prob-
52 abilistic model. Hence, it is difficult to assess the uncertainty in the fitted model estimates,
53 and model extensions are not feasible.

54
55 Mixed effects models have also been employed to model longitudinal metabolomics data.
56 Mei et al. (2009) employ a linear mixed-effects model (LMM) in the context of feature selec-
57 tion for longitudinal metabolomics data, but under the assumption that spectral peaks are
58 independent variables. The high levels of correlation between spectral peaks (i.e. metabo-
59 lites) is biologically important however, and such correlation structure should be explicitly
60 modeled. In a similar vein, Berk et al. (2011) employ smoothing splines mixed-effects models
61 to model longitudinal metabolomics data. While these models have a statistical modelling
62 basis and therefore appropriately model the longitudinal aspect of the data, multiple testing
63 issues (Dudoit et al., 2003) result as the chances of false positives increase with the dimen-
64 sionality of the data. While this problem can be controlled (Benjamini and Hochberg, 1995),
65 dimension reducing features of methods such as PCA are attractive.

66
67 Probabilistic PCA (PPCA) is an approach to PCA based on a Gaussian latent variable
68 model (Tipping and Bishop, 1999; Nyamundanda et al., 2010). PPCA retains the benefits
69 of PCA, such as dimension reduction, while facilitating model extensions through its basis
70 in a statistical model. Here an extension of PPCA called dynamic PPCA (DPPCA) is
71 proposed which allows PPCA to appropriately model the time dependencies in longitudinal
72 metabolomics data. This is achieved by assuming a stochastic volatility model for some
73 of the PPCA parameters. The proposed DPPCA model is closely related to the dynamic
74 factor analysis model (Aguilar and West, 2000) employed to model multivariate financial
75 time series data.

76
77 Data generated in longitudinal metabolomics studies form the basis for the development
78 of the proposed DPPCA model. Examples of such studies include, but are not limited
79 to, postprandial human studies and long term drug treatment studies (Wopereis et al.,
80 2009; Lin et al., 2011; Krug et al., 2012; Nicholson et al., 2012). Interest lies in reducing

81 the dimensionality of the data (for statistical and visualisation purposes) and subsequently
82 highlighting influential metabolites which change over time, while appropriately modelling
83 the longitudinal nature of the data. The proposed DPPCA model is employed to achieve di-
84 mension reduction and model the time dependencies; linear mixed models (LMM) are then
85 employed to identify the metabolites which change over time. The utility of the DPPCA
86 approach is demonstrated through the analysis of data from a longitudinal metabolomics
87 animal study.

88
89 The remainder of the article is structured as follows. An overview of longitudinal
90 metabolomics studies is presented in Section 2. The DPPCA model is introduced in Section
91 3 and the use of stochastic volatility models to account for the correlation due to repeated
92 measurements is detailed. The DPPCA model is estimated within the Bayesian paradigm;
93 accordingly Section 4 specifies the necessary prior distributions and describes the use of
94 Markov chain Monte Carlo (MCMC) techniques to fit the DPPCA model. Section 5 details
95 the application of the DPPCA model to a longitudinal metabolomics data set. Discussion
96 of the developed model and further avenues of research are deferred until the conclusion, in
97 Section 6.

98 2. Longitudinal metabolomics studies

99 In recent years, a number of longitudinal metabolomics datasets have emerged in the lit-
100 erature (Wopereis et al., 2009; Lin et al., 2011; Krug et al., 2012). With regard to human
101 applications, a number of studies employing metabolomics over time following acute chal-
102 lenges such as the oral glucose tolerance test have recently been published and shown to
103 be extremely powerful in studying subtle changes. Applying metabolomics to longitudinal
104 animal studies for determining long term drug toxicity and efficacy is also an important
105 emergent area. In such applications a number of key study aims typically exist which, in
106 general, can be described as follows:

- 107 (i) data visualisation
- 108 (ii) assessing the effect of time within each treatment group and
- 109 (iii) identifying metabolites which change over time within each treatment group.

110 The DPPCA model proposed here helps address these specific aims. In the case of (i) the
111 DPPCA model facilitates visualisation of the study participants in a reduced dimensional
112 space, while appropriately modelling the time course nature of the data. The effect of time
113 within each treatment group (aim (ii)) can be assessed by applying the DPPCA model to
114 the data from each treatment group. An additional output of the DPPCA model is a list of
115 the most influential metabolites within each group. To address aim (iii) univariate analyses
116 with LMM are then carried out to identify those influential metabolites which change over
117 time.

118
119 Metabolomics data from a longitudinal animal study motivate and illustrate the pro-
120 posed DPPCA model. The study has been described in detail in Carmody and Brennan
121 (2010). Briefly, an animal model of epilepsy was employed by repeated administration
122 of pentylenetetrazole (PTZ) which leads to the development of generalised tonic-clonic
123 seizures. Over the administration period (5 weeks) urine samples were collected from treated
124 animals (PTZ treated) and control animals (saline treated animals). The aim of the study

125 was to determine metabolic changes that occur over time during PTZ treatment.

126

127 NMR spectra were acquired from the urine samples and the spectra were integrated into
128 bin regions of 0.04 parts per million (ppm), excluding the water regions (4.0–6.0 ppm). For
129 the purposes of this work, the final acquired data set consists of NMR spectra for $n = 15$
130 animals (8 treated and 7 control), each containing $p = 189$ spectral bin regions, from $M = 8$
131 time points. The $p = 189$ peaks in the spectra at different chemical shift values (measured in
132 ppm) relate to specific metabolites; the height of a peak in any spectrum details the relative
133 abundance of the associated metabolite in the animal’s urine sample. Figure 1 illustrates
134 a metabolomic spectrum resulting from the urine sample collected at a single time point
135 from an animal in the study.

136

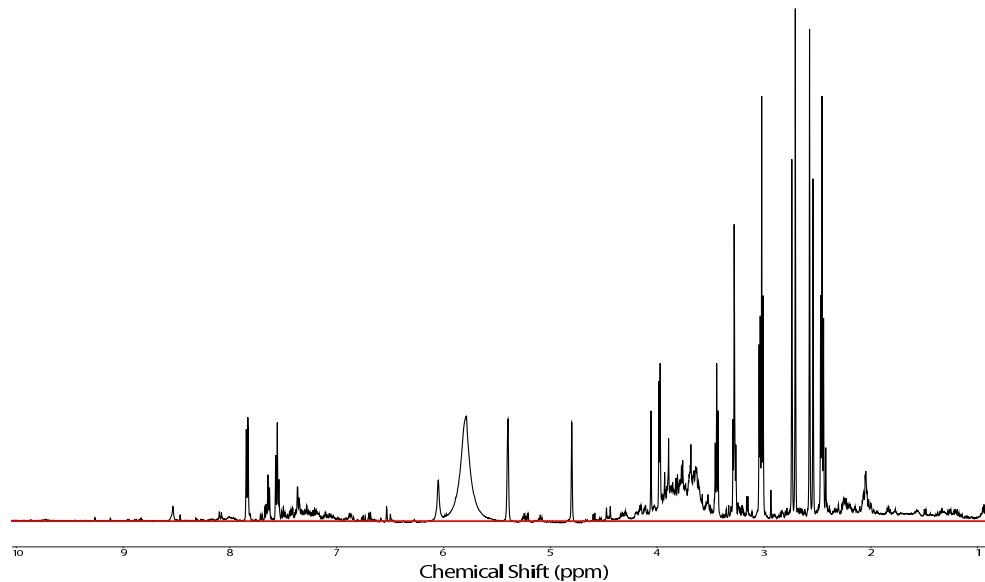


Fig. 1. A metabolomic profile resulting from the urine sample collected at a single time point from an animal in the longitudinal metabolomic study.

137 **3. Dynamic Probabilistic Principal Components Analysis**

138 Probabilistic principal components analysis (PPCA) is a latent factor model constrained
139 such that the maximum likelihood estimates of the parameters span the principal subspace
140 of conventional PCA. Given its underlying assumptions however, PPCA is only applicable
141 to data from a cross sectional study. Here an extension of PPCA to a dynamic PPCA
142 (DPPCA) model is developed; a brief introduction to PPCA, and its extension to the
143 DPPCA model, are detailed in what follows.

144 **3.1. Probabilistic Principal Components Analysis (PPCA)**

PPCA is a generative statistical model which models a high-dimensional observed data point as a linear function of a corresponding low-dimensional latent variable plus isotropic (full-dimensional) noise. For each of n animals, let $\mathbf{x}_i^T = (x_{i1}, \dots, x_{ip})$ denote the set of p observed variables for animal i (eg. an NMR spectrum with p spectral bins). The PPCA model relates each \mathbf{x}_i to a q -dimensional latent Gaussian variable \mathbf{u}_i (typically $q \ll p$) through the linear model:

$$\mathbf{x}_i = W\mathbf{u}_i + \boldsymbol{\epsilon}_i$$

145 where W is a $p \times q$ loadings matrix and the error term $\boldsymbol{\epsilon}_i$ is assumed to have a multivariate
 146 Gaussian distribution, centred at zero with covariance $\sigma^2 I$, where I denotes the identity
 147 matrix. The error term models the part of the observed data which cannot be accounted for
 148 by the q underlying latent variables, or principle components (PCs). Assuming a standard
 149 multivariate normal (MVN) distribution for \mathbf{u}_i , each data point has a zero mean multivariate
 150 normal distribution with covariance $WW^T + \sigma^2 I$.

151 Crucially, the likelihood of the PPCA model is maximized when the columns of W
 152 span the principal subspace of conventional PCA (Tipping and Bishop, 1999). Thus the
 153 maximum likelihood estimate of the loadings matrix in PPCA corresponds exactly to the
 154 loadings matrix in conventional PCA. Hence the model output in PPCA is exactly that
 155 obtained in conventional PCA, but with the additional advantages of uncertainty assessment
 156 and potential model extensions.

157 **3.2. Dynamic Probabilistic Principal Components Analysis (DPPCA)**

158 The derivation of PCA from a probabilistic framework facilitates the development of dy-
 159 namic PPCA as a tool for modelling longitudinal multivariate data. Under the DPPCA
 160 model, the set of p observed variables \mathbf{x}_{im} for animal i at time point m ($m = 1, \dots, M$) is
 161 modeled as:

$$\mathbf{x}_{im} = W_m \mathbf{u}_{im} + \boldsymbol{\epsilon}_{im} \quad (1)$$

162 where W_m , the loadings, and $\mathbf{u}_{im}^T = (u_{i1m}, \dots, u_{iqm})$, the latent scores, vary with time.
 163

164 Unlike the PPCA model which constrains the covariance of the multivariate Gaussian
 165 distribution of the latent variables to be an identity matrix, the DPPCA model eases the
 166 equal variance restriction such that

$$p(\mathbf{u}_{im}) = \text{MVN}_q(\mathbf{0}, H_m)$$

167 where $H_m = \text{diag}(h_{1m}, \dots, h_{qm})$. This assumption allows the variances of the underlying
 168 latent variables to differ across the latent dimensions and to depend on time.

169

The error, $\boldsymbol{\epsilon}_{im}$, for animal i at time m is also assumed to have a multivariate Gaussian distribution:

$$p(\boldsymbol{\epsilon}_{im}) = \text{MVN}_p(\mathbf{0}, \sigma_m^2 I).$$

170 Again, the variance parameter σ_m^2 varies with time. The errors, $\boldsymbol{\epsilon}_{im}$ and the latent variables
 171 (or scores), \mathbf{u}_{im} are assumed to be mutually independent for all $m = 1, \dots, M$.

172

173 While the variance parameter of the error terms σ_m^2 varies with time, it is constrained to
 174 be constant across all observed variables. This is in line with the assumptions of the under-
 175 lying PPCA model; should the variances be unconstrained across variables a dynamic factor
 176 analytic model results (McNicholas and Murphy, 2008; Aguilar and West, 2000). Thus the
 177 DPPCA model can be viewed as a constrained dynamic factor model.

178
 179 The choice of developing the DPPCA model, rather than employing an alternative dy-
 180 namic factor model to analyse the metabolomic data under study, deserves explanation.
 181 The manner in which time dependence is accounted for in the DPPCA model, and the
 182 constraints employed, are motivated by the explicit needs of the motivating metabolomics
 183 application. The metabolomics practitioners are interested in time evolving metabolites,
 184 hence the need for a different loadings matrix at each time point, leading to a highly pa-
 185 rameterised model. Further, strongly motivated by the ubiquitous use, understanding and
 186 acceptance of PCA in the metabolomics field (Smolinska et al., 2012; Cassol et al., 2013;
 187 Carvalho et al., 2013; Bathen et al., 2013; Sachse et al., 2012), maintaining a link to PPCA
 188 was deemed to be highly desirable. As the link to PPCA occurs by constraining the er-
 189 ror variances to be equal, this modelling decision satisfied the metabolomic scientists, and
 190 provided a more parsimonious model than a generic dynamic factor model. The appropri-
 191 ateness of the DPPCA model assumptions are assessed after model fitting in Section 5.4,
 192 using posterior predictive model checking.

193 3.3. Stochastic Volatility Models

194 Stochastic volatility models (Jacquier et al., 1994; Kim et al., 1998) are popular in econo-
 195 metrics and finance where they are typically employed to model the variance of returns over
 196 time, which are highly correlated. The DPPCA model accounts for the correlation due to
 197 repeated measurements through the use of stochastic volatility (SV) models. Specifically,
 198 the DPPCA model assumes that at time point m the variances h_{1m}, \dots, h_{qm} of the latent
 199 variables and the error variances σ_m^2 follow a latent stochastic process. These assumptions
 200 allow the DPPCA model to account for any potential time dependence in longitudinal mul-
 201 tivariate data.

202
 203 Again, the motivation behind the incorporation of SV models in DPPCA requires ex-
 204 planation. While SV models typically model settings with many time points (Aguilar and
 205 West, 2000), they have been employed when modelling longitudinal multivariate data, where
 206 the number of time points is low. Ramoni et al. (2002), Fang-Xiang et al. (2005) and Wang
 207 et al. (2008), for example, employ SV models for modelling high dimensional time course
 208 data where the number of time points ranges from 8 to 18. Hence the SV model was deemed
 209 suitable to model the evolution of the latent variables over time. The appropriateness of
 210 the SV model assumptions is assessed after model fitting in Section 5.4.

211 3.3.1. A stochastic volatility model for the latent variables

212 An SV model on the latent variable u_{ijm} of animal i ($i = 1, \dots, n$) for principal component
 213 j ($j = 1, \dots, q$) at time point m ($m = 1, \dots, M$) can be expressed as:

$$u_{ijm} = \exp(\lambda_{jm}/2)\zeta_{ijm}$$

214 where $\lambda_{jm} = \log(h_{jm})$ is known as the log volatility and ζ_{ijm} , which has a standard univari-
 215 ate Gaussian distribution, denotes the error term of the SV model. Thus the conditional
 216 distribution of the latent variable is $u_{ijm}|\lambda_{jm} \sim N[0, \exp(\lambda_{jm})]$. The q -vector of log volatil-
 217 ities, $\boldsymbol{\lambda}_m^T = (\lambda_{1m}, \dots, \lambda_{qm})$, is assumed to have a stationary first order vector autoregressive
 218 process VAR(1) centered around a mean $\boldsymbol{\mu}^T = (\mu_1, \dots, \mu_q)$:

$$\boldsymbol{\lambda}_m = \boldsymbol{\mu} + \Phi(\boldsymbol{\lambda}_{m-1} - \boldsymbol{\mu}) + \mathbf{R}_m$$

219 where Φ is a matrix of persistence parameters and $\mathbf{R}_m \sim \text{MVN}_q(\mathbf{0}, V)$ are independent in-
 220 novations. The model restricts dependencies across the principal dimensions by constraining
 221 the matrix of persistence parameters Φ and the covariance of the innovations V to be di-
 222 agonal i.e. $\Phi = \text{diag}(\phi_1, \dots, \phi_q)$ and $V = \text{diag}(v_1^2, \dots, v_q^2)$ respectively. The innovation
 223 variance v_j^2 is the uncertainty associated with predicting the current log volatility using
 224 the log volatility from the previous time point on component j . The persistence param-
 225 eter Φ is the parameter of interest; it measures the strength of the relationship between
 226 time points. For stationarity, the persistence parameter ϕ_j is constrained to lie between
 227 -1 and 1 (Kim et al., 1998). The initial state, by stationarity, is drawn from the model
 228 $p(\boldsymbol{\lambda}_1) = \text{MVN}_q[\boldsymbol{\mu}, \text{diag}(\frac{v_1^2}{1-\phi_1^2}, \dots, \frac{v_q^2}{1-\phi_q^2})]$. The distribution of the log volatilities $\boldsymbol{\lambda}_m$ given
 229 the log volatilities of the previous time point $\boldsymbol{\lambda}_{m-1}$ is given by $\text{MVN}_q[\boldsymbol{\mu} + \Phi(\boldsymbol{\lambda}_{m-1} - \boldsymbol{\mu}), V]$
 230 for $m > 1$.

231

232 Constraining the covariance matrix V to be diagonal is a modelling decision motivated
 233 by the fact that the PPCA model does not facilitate dependence across the principal com-
 234 ponents and PPCA underpins the DPPCA model, as detailed in Section 3.2. Such a model
 235 was considered by Harvey et al. (1994), Kim et al. (1998) and Jacquier et al. (1995) among
 236 others; Aguilar and West (2000) allow correlation across dimensions, motivated by their
 237 financial application area.

238 3.3.2. A stochastic volatility model for the errors

239 Additionally, another SV model is adopted to model the potential time dependence in the
 240 errors of the DPPCA model. The p -vector of errors of observation i at time m can be
 241 expressed as $\epsilon_{im} = \exp[\eta_m/2]\boldsymbol{\xi}_{im}$ where $\eta_m = \log(\sigma_m^2)$ is the log volatility at time m and
 242 $\boldsymbol{\xi}_{im} \sim \text{MVN}_p(\mathbf{0}, I)$. The log volatilities η_m on the errors are assumed to have a stationary
 243 first order autoregressive process AR(1):

$$\eta_m = \nu + \phi(\eta_{m-1} - \nu) + r_m$$

244 where the center of the AR(1) model is ν and the persistence parameter ϕ is constrained
 245 such that $\phi \in [-1, 1]$. The innovations of the AR(1) model are assumed to be normally
 246 distributed, $r_m \sim N(0, v^2)$. It follows that the initial state of the SV model is $p(\eta_1) =$
 247 $N(\nu, \frac{v^2}{1-\phi^2})$ and that $p(\eta_m|\eta_{m-1}) = N[\nu + \phi(\eta_{m-1} - \nu), v^2]$ for $m > 1$. Note that, as stated
 248 in Section 3.2, to maintain the link to PPCA and for reasons of parsimony, each of the p
 249 dimensions in the error ϵ_{im} are constrained to follow the same AR(1) model.

250 **4. Estimation of the DPPCA model**

251 Under the DPPCA model, the full augmented data likelihood function based on the data
 252 $X = (X_1, \dots, X_n)$ and the latent variables $U = (U_1, \dots, U_n)$, $\Lambda = (\boldsymbol{\lambda}_1, \dots, \boldsymbol{\lambda}_M)$ is:

$$p(X, U, \Lambda, \boldsymbol{\eta} | W, \theta_1, \theta_2) = \left[\prod_{m=1}^M \prod_{i=1}^n p(\mathbf{x}_{im} | W_m, \mathbf{u}_{im}, \eta_m) p(\mathbf{u}_{im} | \boldsymbol{\lambda}_m) \right] p(\boldsymbol{\eta} | \theta_1) p(\Lambda | \theta_2)$$

253 where $\theta_1 = (\nu, \phi, v^2)$ and $\theta_2 = (\boldsymbol{\mu}, \boldsymbol{\Phi}, V)$ denote the SV model parameters on the errors
 254 and latent scores respectively. The PPCA model on each time point $p(\mathbf{x}_{im} | W_m, \mathbf{u}_{im}, \eta_m)$
 255 is $\text{MVN}_p[W\mathbf{u}_{im}, \exp(\eta_m)I]$.

256

257 A Bayesian approach is taken when estimating the DPPCA model; this requires the
 258 specification of prior distributions for all the model parameters. The resulting posterior
 259 distribution is intricate and Markov chain Monte Carlo methods are necessary to produce
 260 realizations of the model parameters. Specifically, a Metropolis-within-Gibbs algorithm is
 261 required to sample from the full conditional distributions for all model parameters and la-
 262 tent variables.

263

264 **4.1. Prior distributions**

265 Prior distributions over the full set of the model parameters need to be specified. It
 266 is assumed that the prior distributions on the model parameters are independent. Under
 267 the PPCA part of the DPPCA model, the only parameters are the loadings matrices
 268 W_1, \dots, W_M . A q -dimensional multivariate normal prior distribution, centered at $\mathbf{0}$ with
 269 covariance Ω_m , is assumed for each row of the loadings matrix W_m at time m .

270

271 The remaining model parameters are all parameters of the SV part of the DPPCA model.
 272 Non-informative normal prior distributions are specified on the means of the SV models i.e.
 273 a $N(\mu_\nu, \sigma_\nu^2)$ distribution is specified for ν and a $N(\mu_\mu, \sigma_\mu^2)$ distribution is assumed on each
 274 of the univariate elements of $\boldsymbol{\mu}$, where the variance hyperparameter in each of these priors
 275 is large. A conjugate prior is assumed for the variances of the innovations in the SV models
 276 i.e. an inverse gamma $IG(\alpha/2, \beta/2)$ distribution is chosen for the prior distribution of v^2
 277 and for each of the diagonal elements of V . For stationarity, the persistence parameters of
 278 the SV models are constrained to lie in $[-1, 1]$; accordingly the prior distributions on ϕ and
 279 on the diagonal elements of $\boldsymbol{\Phi}$ are truncated normal distributions, $N_{[-1,1]}(\mu_\phi, \sigma_\phi^2)$.

280

281 As in any Bayesian setting, the choice of prior distribution can potentially influence
 282 parameter inference. Sensitivity analyses were conducted to assess the influence of different
 283 choices of priors on the resulting posterior distribution. Some sensitivity was observed in
 284 the case of the persistence parameters. Kim et al. (1998) employ a transformed beta prior
 285 for the persistence parameters, but sensitivity analyses here suggested that the posterior
 286 distribution strongly depended on the values of the hyperparameters used. In a similar
 287 setting to the DPPCA model, Aguilar and West (2000) employ a truncated (between ± 1)
 288 Gaussian prior for the persistence parameters; the posterior distributions were less sensitive
 289 to the parameter specification under this prior. Thus, a Gaussian prior, truncated (between
 290 ± 1), was employed here for the persistence parameters.

291 4.2. The Metropolis-within-Gibbs sampler

292 Given the specified prior distributions, the resulting posterior distribution is intricate and
 293 Markov chain Monte Carlo (MCMC) methods are required to produce realizations of the
 294 model parameters. The full conditional distributions for the loadings matrices W_m , the
 295 latent scores U_m , the SV model means ν and μ , and the SV model innovation variances
 296 v^2 and V exist in standard form, and a straightforward Gibbs sampler can be employed to
 297 draw samples. However, the full conditional distributions for the persistence parameters ϕ
 298 and Φ and for the log volatilities Λ and η are not available in closed form; values from these
 299 distributions are therefore sampled using a Metropolis Hastings step. Hence a Metropolis-
 300 within-Gibbs algorithm (Gilks et al., 1996) is required to sample from the full conditional
 301 distributions for all model parameters and latent variables. Carlin and Louis (2000) detail
 302 the conditions necessary for the convergence of such a hybrid algorithm.

303
 304 Detailed derivations of the full conditional distributions for the DPPCA model param-
 305 eters and latent variables are given in the Supplementary Material. For the Metropolis-
 306 Hastings steps to update the log volatilities, proposal distributions which are closely related
 307 to the shape and orientation of the target full conditional distributions provide an im-
 308 proved rate of convergence. To achieve this, second order Taylor expansions of the full
 309 conditional distributions for η and Λ are employed to guide the choice of an effective pro-
 310 posal distribution and its parameter values (Kim et al., 1998). A summary of one sweep of
 311 the Metropolis-within-Gibbs sampler for the DPPCA model is given in the Supplementary
 312 Material.

313 4.3. Model Identification

314 As with factor analytic models, the DPPCA model suffers from identification issues. Sub-
 315 jecting the loadings matrix and latent scores to an orthogonal rotation gives rise to the same
 316 distribution for the observed data. Thus it is not possible to identify the model parameters
 317 from the observed data unless restrictions are imposed.

318
 319 Many attempts to deal with non-identifiability of the related factor analytic models are
 320 detailed in the literature. Most commonly, a unique model is defined by constraining the
 321 loadings matrix such that the first q rows are lower-triangular with positive diagonal ele-
 322 ments (Geweke and Zhou, 1996). However imposing this structure also imposes structure
 323 on the ordering of the variables (Aguilar and West, 2000). Within the context of the mo-
 324 tivating metabolomics application, such a structure cannot be imposed on the variables as
 325 the ordering of the spectral peaks within a metabolomics spectrum is important.

326
 327 The approach taken here is to estimate a fully unconstrained loadings matrix using
 328 the Metropolis-within-Gibbs sampler detailed in the Supplementary Material. Procrustean
 329 techniques (Borg and Groenen, 2005) are then employed to post-process the sampled load-
 330 ings matrices to match them to the maximum likelihood estimate (MLE) of the loadings
 331 matrix resulting from fitting a PPCA model to data from the relevant time point. The
 332 MLE is used only as a template, to identify the model. The transformation required to
 333 match the loadings matrices is also applied to the latent scores. In practice, this has proved
 334 to be a fast and satisfactory approach to dealing with model non-identifiability.

335 **5. Results**

336 As detailed in Section 2, three specific issues associated with the longitudinal metabolomics
 337 study need to be addressed: (i) data visualisation, (ii) assessing the effect of time within
 338 each treatment group and (iii) identifying the specific metabolites which change over time
 339 within each treatment group. The DPPCA model, in combination with linear mixed models,
 340 is fitted to the longitudinal metabolomics data set to address these issues. For reasons of
 341 visual clarity, only models with $q = 2$ were considered. For each set of results detailed
 342 below, the prior distributions employed for the DPPCA model parameters were specifically:

$$\begin{aligned} \mathbf{w}_{km} &\sim \text{MVN}_q(\mathbf{0}, I) \quad \text{for } k = 1, \dots, p \text{ and } m = 1, \dots, M. \\ \nu &\sim N(0, 10) \\ v^2 &\sim \text{IG}(6/2, 0.5/2) \\ \phi &\sim N_{[-1,1]}(0.75, 0.1) \end{aligned}$$

343 The priors on the univariate entries of the set of parameters $\theta_2 = (\boldsymbol{\mu}, \Phi, V)$ were the
 344 same as those for $\theta_1 = (\nu, \phi, v^2)$. The Metropolis-within-Gibbs sampler was run for 500,000
 345 iterations, thinned every 500th iteration. The first 5,000 iterations were discarded as burn-in.
 346 The MCMC algorithm was initialized using estimates of the loading matrices from fitting
 347 a PPCA model to data from each time point independently; stochastic volatility model
 348 parameters were set equal to their prior means. Trace plots and autocorrelation function
 349 (ACF) plots for the MCMC samples of the parameters were used to assess convergence of
 350 the algorithm.

351 **5.1. Data Visualisation: Exploring Metabolomic Trajectories**

352 In longitudinal metabolomics studies, trajectories through the latent principal subspace can
 353 be used to gain visual insight to the response of animals during the study period. Examining
 354 the location, magnitude and direction of these metabolomic trajectories provides visual
 355 insight to the metabolomic changes over time.

356 Here metabolomic trajectories were estimated using the latent scores of animals resulting
 357 from collectively modelling data from both treatment groups using a DPPCA model. Such
 358 a model takes into account the covariation between the metabolites and any correlation
 359 across time; this facilitates visualisation of animals in a reduced dimensional space, while
 360 appropriately modelling the time course nature of the data. Trace plots for the estimated
 361 latent scores and loadings are given in the Supplementary Material.

362 The metabolomic trajectories of four randomly sampled animals are illustrated in Fig-
 363 ure 2. Under the DPPCA model, each time point m has a different principal subspace,
 364 defined by the columns of the relevant loadings matrix W_m . Hence the latent scores of
 365 animals at different time points lie in different subspaces. To visualise the metabolomic
 366 trajectories the latent scores must therefore be unified. This is achieved by again drawing
 367 on Procrustean ideas, where the loadings matrix from the first time point is used as the
 368 reference matrix. The loadings matrix from each subsequent time point m is rotated to best
 369 match the loadings matrix from the first time point; the same rotation is then applied to the
 370 associated set of scores from time point m . This facilitates illustration of the movement of
 371 the latent scores over time within the same principal subspace. Figure 2 therefore provides
 372 the latent scores over time within the same principal subspace. Figure 2 therefore provides
 373

374 visual insight to the animals' metabolomic trajectories in the principal subspace from the
 375 first time point.

376

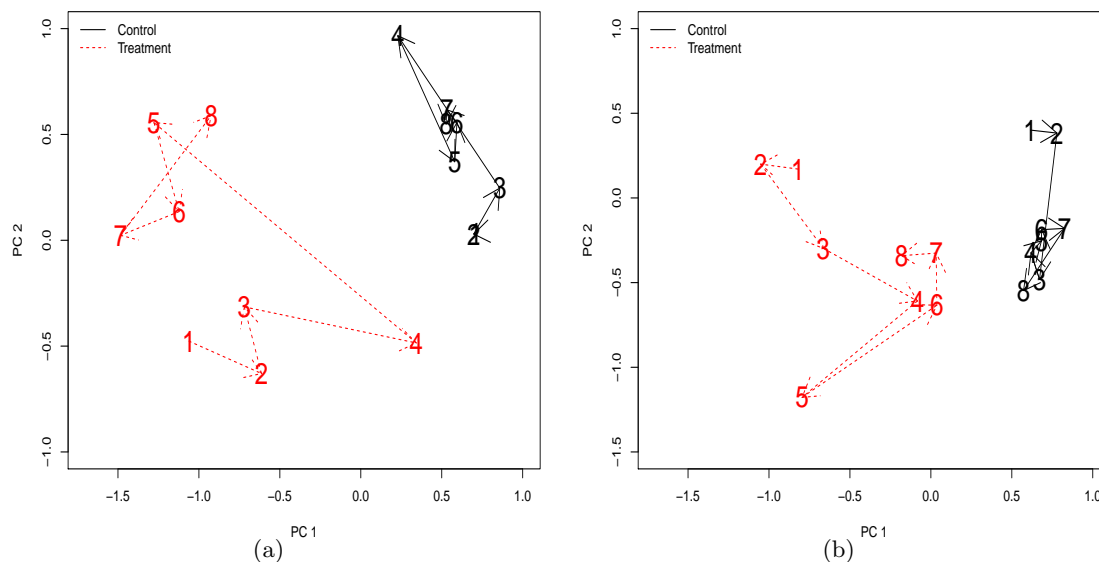


Fig. 2. Individual trajectories for four randomly sampled animals, in the principal subspace from the first time point. (a) An animal from the control group (black solid lines) and an animal from the treated group (red dashed lines) and (b) an animal from the control group (black solid lines) and an animal from the treated group (red dashed lines). The digits represent the time points of the study and arrows illustrate movement through time.

377 Figure 2 suggests the presence of a treatment effect through the visible separation of the
 378 locations of the treated and control animals in the principal subspace from the first time
 379 point. The difference in the biochemical composition of the urine due to treatment is high-
 380 lighted by the different 'metabolic starting positions' of the trajectories for the randomly
 381 selected animals from the control group and those from the treatment groups. This is due
 382 to the fact that the urine samples analysed at time point 1 actually resulted from day 3 of
 383 the study, at which stage the treatment is apparently having an effect.

384

385 The trajectories also demonstrate that the magnitude of the metabolic changes in the
 386 biochemical composition of the urine samples is much greater in the treatment group than in
 387 the control group, over time. This is evidenced by the larger movements between time points
 388 by the treated animals. This shows that the variability in the urinary composition of the
 389 treated animals over time is greater than that in the control group. Thus, the metabolomic
 390 trajectories provide a visual insight to the metabolomic changes occurring over time.

391

392 **5.2. Exploring the Effect of Time**

393 The second aim of the longitudinal study was to ascertain if there is a time effect within
 394 each treatment group. In an effort to quantify the effect of time, the DPPCA model was
 395 fitted separately to each treatment group. If a time effect is established, the task will then
 396 be to identify metabolites whose concentration level is significantly changing over time.

397 **5.2.1. Exploring the Effect of Time in the Treatment Group**

398 The DPPCA model was fitted to the metabolomic spectra from the animals in the treat-
 399 ment group. The persistence parameters in the SV models are the parameters of interest
 400 as they quantify the strength of the relationship between the time points. Figure 3(a) il-
 401 lustrates the posterior distribution of the persistence parameter (ϕ) of the SV model on
 402 the errors. The relevant trace and ACF plots are given in Figure 3(b) and Figure 3(c) re-
 403 spectively. The posterior mean of ϕ was large and positive ($\hat{\phi} = 0.69$) and significant (95%
 404 quantile based credible interval (CI) (0.15, 0.97)). The persistence parameters of the SV
 405 model on the latent variables for PC 1 and PC 2 were also estimated to be large and signifi-
 406 cant at $\hat{\phi}_1 = 0.64$ (0.07, 0.97) and $\hat{\phi}_2 = 0.66$ (0.08, 0.97), respectively. The posterior means
 407 suggest that a positive time dependency exists among the spectra from the treatment group.
 408

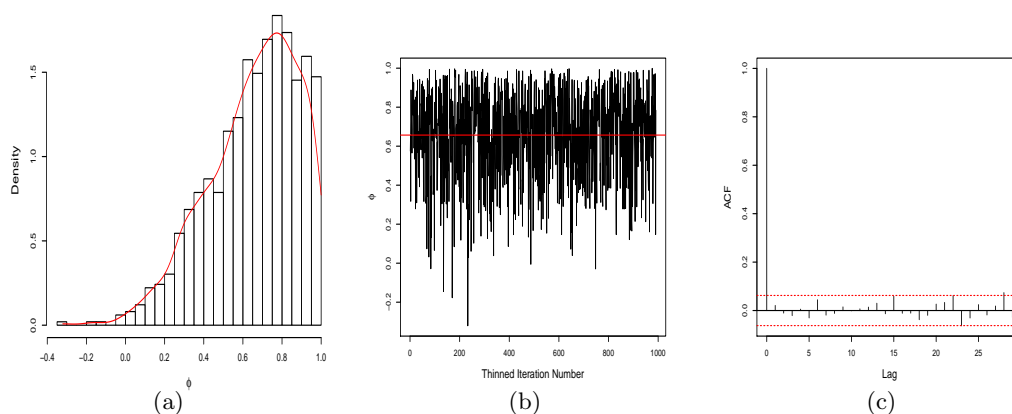


Fig. 3. The persistence parameter, ϕ , of the SV model on the error variances in the treatment group: (a) plot of the posterior density, (b) trace plot and (c) ACF plot. The horizontal line in (b) illustrates the posterior mean of ϕ .

409 Given that a time effect has been established, the third aim of the study was to identify
 410 the specific metabolites which change over time within the treatment group. This is achieved
 411 by first using the DPPCA model to expose those metabolites which influence the data struc-
 412 ture at each time point. Under the DPPCA model, this translates to identifying a subset
 413 of metabolites whose posterior mean loadings are largest (in terms of magnitude) at each
 414 time point. Standard linear mixed models are then fitted to these ‘influential metabolites’
 415 to identify those which change over time. This approach yields a panel of metabolites which
 416 evolve over time, while appropriately accounting for the covariation in the high-dimensional
 417 data, and the time related dependencies.

Table 1. Posterior means of the persistence parameters and the corresponding 95% CIs for the control group.

SV model	Estimate	(95% CI)
Errors (ϕ)	0.66	(0.09,0.98)
PC 1 (ϕ_1)	0.65	(0.10,0.98)
PC 2 (ϕ_2)	0.66	(0.07,0.97)

419 After fitting the DPPCA model to the spectra from animals in the treatment group,
 420 several spectral regions (corresponding to metabolites) were identified as influencing the
 421 underlying structure of the data. At each time point, the absolute values of the posterior
 422 mean loadings on PC1 were ranked in descending order. The top five influential spectral
 423 bins at each time point were determined and are shown in Figure 4. None of the 95%
 424 CIs associated with these spectral bins included zero. The set of the top five spectral bins
 425 across all $M = 8$ time points consists of only eight unique spectral bins (2.46ppm, 2.54ppm,
 426 2.58ppm, 2.66ppm, 2.7ppm, 2.74ppm, 3.02ppm and 3.26ppm).

427
 428 Bayesian linear mixed models were fitted to the data associated with the eight unique
 429 influential spectral bins to determine which, if any, have concentrations which evolve over
 430 time. A random intercept model with cubic time effect was the most complex model consid-
 431 ered; no interaction terms were considered. A backwards selection type approach was taken
 432 to model selection for each spectral bin considered. Of the eight spectral bins considered,
 433 six were deemed to have significantly fluctuating concentration levels over time. Figure 5
 434 illustrates the predicted average intensity levels for each of the six spectral bins.

435
 436 The metabolites identified to be evolving over time include the metabolite 2-oxoglutarate,
 437 represented by the spectral bins 2.46ppm and 3.02ppm. The concentration level of 2-
 438 oxoglutarate decreases initially during the study and increases at later time points, as
 439 illustrated by the similar behaviour of the predicted intensities of 2.46ppm and 3.02ppm in
 440 Figure 5. The model also predicts a linear decreasing metabolic time profile for spectral bin
 441 2.7ppm. Spectral bin 2.54ppm has a positive quadratic time effect in the treated animals
 442 i.e. the concentration level decreases and then increases over time. Spectral bins 2.58ppm
 443 and 3.26ppm have a positive linear time trend. Individual animal and predicted profiles for
 444 three of the six evolving spectral bins are given in the Supplementary Material.

445 5.2.2. Exploring the Effect of Time in the Control Group

446 To establish the presence or absence of a time effect in the control group of animals, and
 447 to subsequently highlight those metabolites which evolve over time, the same approach as
 448 that taken in Section 5.2.1 was followed. That is, the DPPCA model was fitted to the
 449 spectra of animals in the control group only; Table 1 details the posterior means of the
 450 persistence parameters of the SV model on the errors and on the latent variables, with their
 451 corresponding 95% CIs. Table 1 shows that the persistence parameters of the SV models
 452 are large and significant, suggesting that there is a relationship across time.

453
 454 Given that a time effect has been established in the control group, interest then lies in
 455 highlighting those metabolites which evolve over time. The posterior mean PC1 loadings of

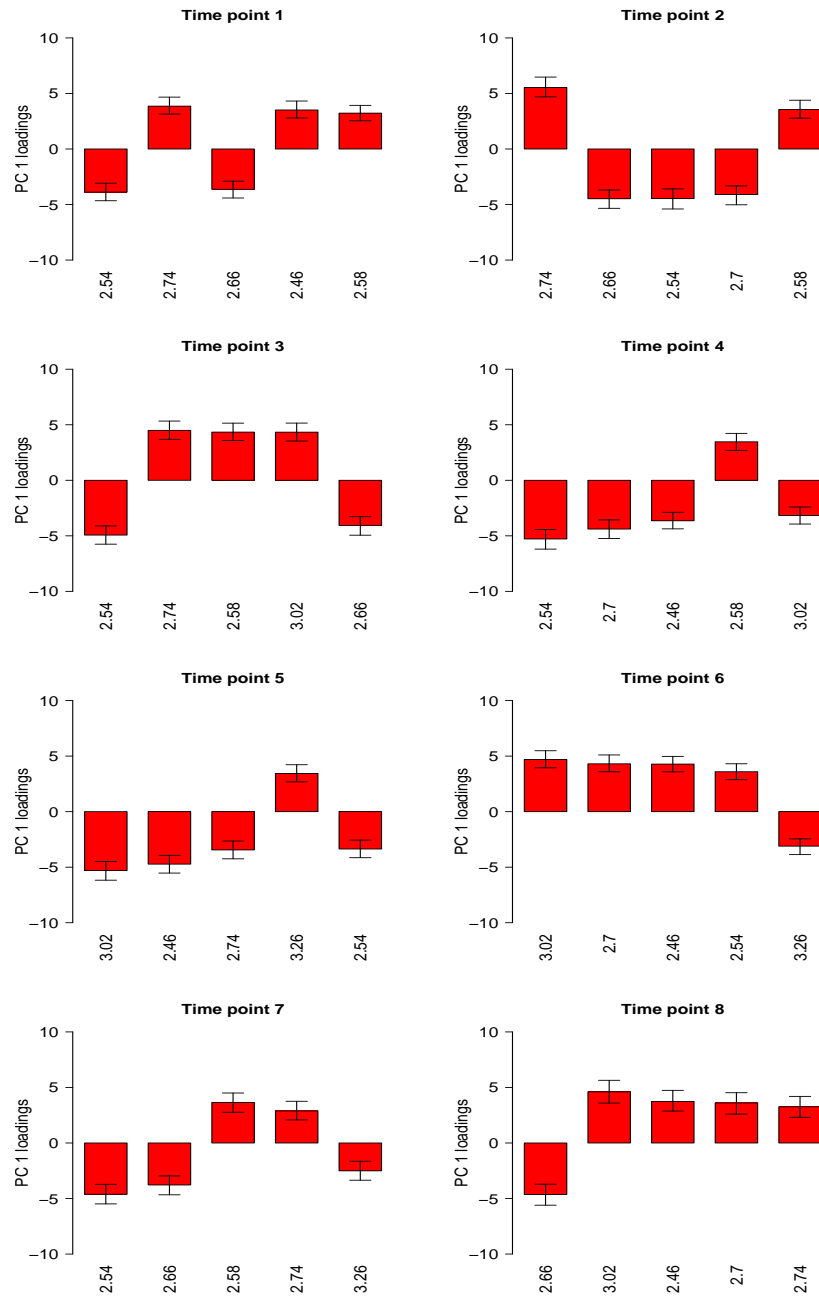


Fig. 4. Barplots of the posterior mean loadings for the top five influential spectral bins, which correspond to metabolites, in the treatment group. The error bars are the corresponding 95% quantile based credible intervals.

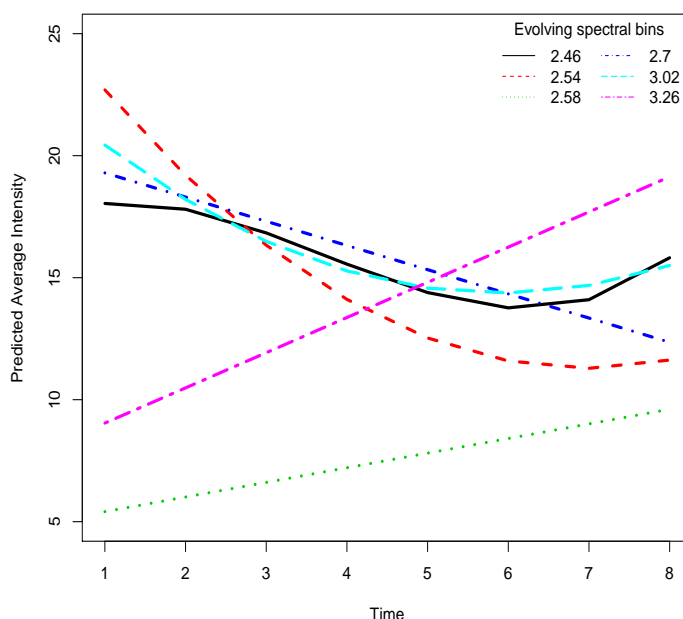


Fig. 5. The LMM predicted average intensities of the six influential spectral bins which evolve over time in the treatment group.

456 the DPPCA model were ranked to select the top five influential spectral bins at each time
 457 point; again, none of the associated 95% CIs included zero. From this list of spectral bins,
 458 those which evolve over time in the control group were identified. Seven unique influential
 459 spectral bins were ranked in the top five over the eight time points; Bayesian LMM models
 460 were fitted to the profiles for each of these and all seven were identified as evolving over
 461 time. Figure 6 illustrates the predicted average intensity levels over the eight time points,
 462 under the selected LMM for each of the seven evolving spectral bins.

463

464 The metabolite 2-oxoglutarate (with corresponding spectral bins 2.46ppm and 3.02ppm)
 465 was predicted by the Bayesian LMM to have a negative quadratic time effect in the control
 466 group i.e. its concentration increases and then decreases over time (see Figure 6). Spectral
 467 bins 2.54ppm and 3.42ppm have positive quadratic time effects. The remaining evolving
 468 spectral bins (2.58ppm, 2.7ppm and 3.26ppm) have cubic time effects. Individual animal and
 469 predicted profiles for three of the seven evolving spectral bins are given in the Supplementary
 470 Material.

471 5.3. Comparing evolving metabolites in the two treatment groups

472 As the aim of the longitudinal metabolomics study was to determine metabolic changes that
 473 occur over time during PTZ treatment, of interest are the similarities and differences be-
 474 tween the set of evolving metabolites in the treatment group and the set in the control group.

475

476 A total of six spectral bins were highlighted as evolving in the treatment group and
 477 seven in the control group. There is considerable overlap between the two sets of evolving

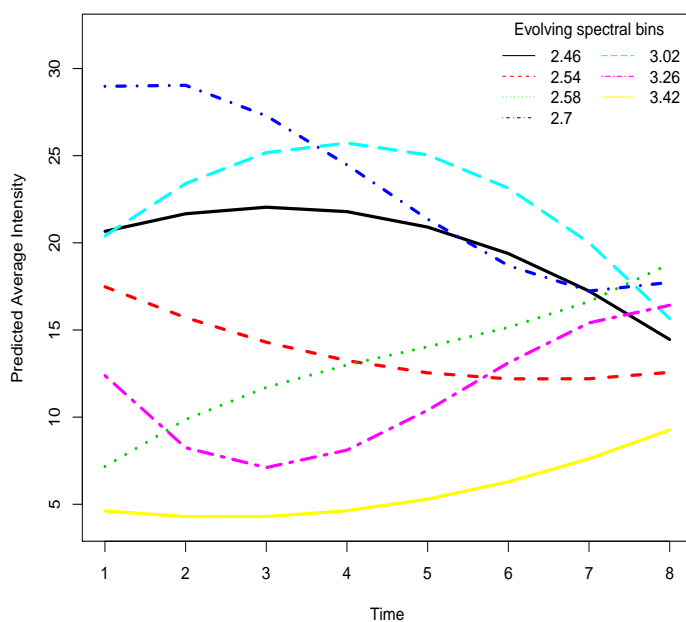


Fig. 6. The LMM predicted average intensities of the seven influential spectral bins which evolve over time in the control group.

478 bins, with 3.42ppm evolving in the control group only. While some of the common spec-
 479 tral bins had the same evolution pattern, some differed. In particular, the spectral bins
 480 2.46ppm and 3.02ppm relating to the 2-oxoglutarate metabolite were predicted to have op-
 481 posite quadratic effects in the treatment group and in the control group. Figure 7, which
 482 shows the predicted average intensities for these two spectral bins only in both treatment
 483 groups, clearly illustrates this phenomenon. The biological basis of the diverse response of
 484 this metabolite will be investigated in future metabolomic experiments.

485

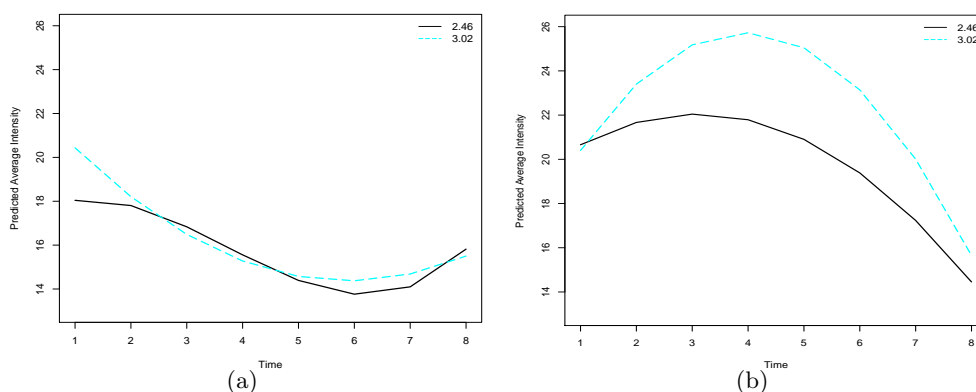


Fig. 7. The LMM predicted average intensities of the two spectral bins 2.46ppm and 3.02ppm which relate to the metabolite 2-oxoglutarate in (a) the treatment group and (b) the control group.

486 5.4. Assessing model fit

487 As with any applied statistical analysis, the modelling assumptions employed need to be
488 assessed to ensure valid inference. In the case of the DPPCA model, the modelling assump-
489 tions are the multivariate Gaussian distribution for the latent variables and the error terms,
490 and the stochastic volatility model assumed to control the evolution of the latent variables
491 over time. Posterior predictive model checking (Gelman et al., 2003) was employed to assess
492 these modelling assumptions. Replicated data were simulated from the posterior predictive
493 distribution and compared to the observed data from each treatment group. Given the
494 multivariate nature of the data, the replicated and observed data were compared by exam-
495 ining the mean absolute deviations (MADs) between the covariance matrix of the observed
496 data and the covariance matrix of the replicated data at each time point (Ansari et al.
497 (2002)). The resulting MADs suggested that the DPPCA model fits well since the vast
498 majority of the deviations were close to zero. A histogram of the MADs is available in the
499 Supplementary Material. There were some large MADs (6% of MADs were > 1 for the
500 treatment group data and 4% for the control group data) but given the large number of
501 covariance parameters being compared, this was not viewed as sufficient evidence of invalid
502 assumptions and poor model fit. The few large MADs may arise due to the fact that the
503 number of latent dimensions was fixed at 2 (for visual substantive reasons), and that some
504 parameters were constrained (for reasons of parsimony). Fitting a higher dimensional and
505 less parsimonious model to the time course metabolomic data is an area of further research.

506 6. Discussion

507 analysing longitudinal data from metabolomics studies is problematic due to the dimen-
508 sionality of the data, the correlated metabolites and correlation structure due to repeated
509 measurements over time. Many currently existing approaches to analysing such data sets
510 either have the limitation of confounding treatment variation with variability due to the
511 longitudinal nature of the data or they ignore the fact that metabolites do not work inde-
512 pendently of each other. Here the DPPCA methodology has been proposed which combines
513 probabilistic PCA and stochastic volatility models to disentangle the two types of variation
514 in the data, while also accounting for its high-dimensionality.

515
516 The DPPCA model successfully addressed the aims of the metabolomic study i.e. vi-
517 sualising the metabolomic trajectories through time, quantifying the effect of time, and
518 highlighting metabolites which evolve over time. Importantly, the DPPCA model high-
519 lighted the contrasting behaviour of the 2-oxoglutarate metabolite between the two treat-
520 ment groups under study. Future work will examine further this contrasting behaviour.

521
522 Many areas of further research naturally arise from the DPPCA model. From a practi-
523 cal viewpoint, fitting the DPPCA model is computationally expensive, mostly due to the
524 costly sampling of the log volatilities. Several approaches to sampling log volatilities for
525 SV models are suggested and reviewed by Jacquier et al. (1994); Kim et al. (1998) and
526 Platanioti et al. (2005). Further work in this area would expedite the convergence of the
527 MCMC chain. Also, while data from 16 time points were collected, only 8 time points were
528 analysed here, due to missing data. Imputation of such data would potentially be feasible
529 within the model fitting algorithm.

530

531 Motivated by the real application area, only principal subspaces of dimension 2 were
532 considered here; clearly the choice of dimensionality can be viewed as a model selection
533 issue and any of the myriad of approaches to model selection in the Bayesian paradigm
534 by evaluating the marginal likelihood could be employed; Friel and Wyse (2012) provide
535 a review of such approaches. However, it is anticipated that such approaches would be
536 computationally expensive in the setting of the DPPCA model. Minka (2000) proposes a
537 computationally efficient approach to selecting the optimal dimensionality in PCA, which
538 might also provide a possible solution to the model selection problem here.

539
540 In terms of the DPPCA model itself, the manner in which the dynamics are modelled in
541 the DPPCA model raises further research questions. Alternative approaches to modelling
542 the time dynamics should be examined, for example (as suggested by a referee) using state-
543 space models for the loadings matrix. Further, research into a random effects PPCA model
544 to model such longitudinal metabolomics data is underway (Nyamundanda et al., 2013).
545 The DPPCA approach proposed here can be thought of as an approach to identifying
546 the subset of influential variables, which are then analysed via LMMs to highlight those
547 which are time evolving. Hence, the issue of multiple testing is reduced but not eradicated
548 under the DPPCA model; this could be addressed by employing a hierarchical modelling
549 framework (Gelman et al., 2003). Further, the proposed DPPCA approach to highlighting
550 time evolving metabolites requires a two step process: fitting a DPPCA model, followed by
551 fitting LMMs. A more elegant approach would combine the ideas underlying both models
552 into a single model. Clearly the development of the DPPCA model gives rise to many and
553 varied areas of future work.

554 References

- 555 Aguilar, O. and M. West (2000). Bayesian dynamic factor models and portfolio allocation.
556 *Business and Economic Statistics* 18(3), 338–357.
- 557 Ansari, A., K. Jedidi, and L. Dube (2002). Heterogeneous factor analysis model: a Bayesian
558 approach. *Psychometrika* 67(1), 49 – 78.
- 559 Bathen, T. F., B. Geurts, B. Sitter, H. E. Fjøsne, S. Lundgren, L. M. Buydens, I. S.
560 Gribbestad, G. Postma, and G. F. Giskeødegård (2013). Feasibility of MR metabolomics
561 for immediate analysis of resection margins during breast cancer surgery. *PloS one* 8(4),
562 e61578.
- 563 Benjamini, Y. and Y. Hochberg (1995). Controlling false discovery rate: a practical and
564 powerful approach to multiple testing. *Journal of the Royal Statistical Society, Series*
565 *B* 57, 289–300.
- 566 Berk, M., T. Ebbels, and G. Montana (2011). A statistical framework for biomarker dis-
567 covery in metabolomic time course data. *Bioinformatics* 27(14), 1979–1985.
- 568 Borg, I. and P. J. F. Groenen (2005). *Modern Multidimensional Scaling. Theory and Ap-*
569 *plications*. New York: Springer.
- 570 Carlin, B. P. and T. A. Louis (2000). *Bayes and empirical Bayes methods for data analysis*.
571 New York: Chapman and Hall.

- 572 Carmody, S. and L. Brennan (2010). Effects of pentylenetetrazole-induced seizures on
573 metabolomic profiles of rat brain. *Neurochemistry International* 56(2), 340–344.
- 574 Carvalho, E., P. Franceschi, A. Feller, L. Palmieri, R. Wehrens, and S. Martens (2013). A
575 targeted metabolomics approach to understand differences in flavonoid biosynthesis in
576 red and yellow raspberries. *Plant Physiology and Biochemistry* 72, 79 – 86.
- 577 Cassol, E., V. Misra, A. Holman, A. Kamat, S. Morgello, and D. Gabuzda (2013). Plasma
578 metabolomics identifies lipid abnormalities linked to markers of inflammation, microbial
579 translocation, and hepatic function in HIV patients receiving protease inhibitors. *BMC*
580 *Infectious Diseases* 13(1), 203.
- 581 Choi, Y., H. Kim, H. Linthorst, J. Hollander, A. Lefeber, C. Erkelens, J. Nuzillard, and
582 R. Verpoorte (2006). NMR metabolomics to revisit the tobacco mosaic virus infection in
583 nicotiana tabacum leaves. *Journal of Natural Products* 69(5), 742–748.
- 584 Dettmer, K., P. A. Aronov, and B. D. Hammock (2007). Mass spectrometry-based
585 metabolomics. *Mass Spectrometry Reviews* 26(1), 51–78.
- 586 Dudoit, S., J. P. Shaffer, and J. C. Boldrick (2003). Multiple hypothesis testing in microarray
587 experiments. *Statistical Science* 18(1), 71–103.
- 588 Fang-Xiang, W., W. J. Zhang, and A. J. Kusalik (2005). Dynamic model-based cluster-
589 ing for time-course gene expression data. *Journal of Bioinformatics and Computational*
590 *Biology* 3(4), 821 – 836.
- 591 Friel, N. and J. Wyse (2012). Estimating the evidence – a review. *Statistica Neerlandica* 6,
592 288–308.
- 593 Gelman, A., J. B. Carlin, H. S. Stern, and D. B. Rubin (2003). *Bayesian Data Analysis*.
594 Chapman and Hall/CRC.
- 595 Geweke, J. and G. Zhou (1996). Measuring the price of the arbitrage pricing theory. *The*
596 *Review of Financial Studies* 9(2), pp. 557–587.
- 597 Gilks, W. R., S. Richardson, and D. J. Spiegelhalter (1996). *Markov Chain Monte Carlo in*
598 *Practice*. London: Chapman and Hall.
- 599 Harvey, A., E. Ruiz, and N. Shephard (1994). Multivariate stochastic variance models. *The*
600 *Review of Economic Studies* 61(2), 247–264.
- 601 Jacquier, E., N. G. Polson, and P. E. Rossi (1994). Bayesian analysis of stochastic volatility
602 models. *Journal of Business and Economic Statistics* 12, 371–389.
- 603 Jacquier, É., N. G. Polson, and P. E. Rossi (1995). Models and priors for multivariate
604 stochastic volatility. Technical report, CIRANO.
- 605 Jansen, J., N. van Dam, H. Hoefsloot, and A. Smilde (2009). Crossfit analysis: a novel
606 method to characterize the dynamics of induced plant responses. *BMC Bioinformat-*
607 *ics* 10(1), 425.
- 608 Jansen, J. J., H. C. Hoefsloot, H. F. Boelens, J. van der Greef, and A. K. Smilde (2004).
609 Analysis of longitudinal metabolomics data. *Bioinformatics* 30(15), 2438–2446.

- 610 Jolliffe, I. T. (2002). *Principal Component Analysis, 2nd edition*. New York: Springer.
- 611 Kim, S., N. Shephard, and S. Chibb (1998). Stochastic volatility: likelihood inference and
612 comparison with arch models. *Review of economic studies* 65, 361–393.
- 613 Krug, S., G. Kastenmuller, F. Stuckler, M. J. Rist, T. Skurk, M. Sailer, J. Raffler,
614 W. Romisch-Margl, J. Adamski, C. Prehn, T. Frank, K. H. Engel, T. Hofmann, B. Luy,
615 R. Zimmermann, F. Moritz, P. Schmitt-Kopplin, J. Krumsiek, W. Kremer, F. Huber,
616 U. Oeh, F. J. Theis, W. Szymczak, H. Hauner, K. Suhre, and H. Daniel (2012). The
617 dynamic range of the human metabolome revealed by challenges. *The Journal of the*
618 *Federation of American Societies for Experimental Biology* 26(6), 2607 – 2619.
- 619 Lin, S., Z. Yang, H. Liu, L. Tang, and Z. Cai (2011). Beyond glucose: metabolic shifts in
620 responses to the effects of the oral glucose tolerance test and the high-fructose diet in
621 rats. *Molecular BioSystems* 7(5), 1537–1548.
- 622 McNicholas, P. D. and T. B. Murphy (2008). Parsimonious Gaussian mixture models.
623 *Statistics and Computing* 18(3), 285–296.
- 624 Mei, Y., B. S. Kim, and K. Tsui (2009). Linear mixed effects models for feature selection
625 in high dimensional NMR spectra. *Expert Systems with Applications* 36(3), 4703–4708.
- 626 Minka, T. P. (2000). Automatic choice of dimensionality for PCA. In *NIPS*, Volume 13,
627 pp. 598–604.
- 628 Nicholson, J. K., J. R. Everett, and J. C. Lindon (2012). Longitudinal pharmacometabo-
629 nomics for predicting patient responses to therapy: drug metabolism, toxicity and efficacy.
630 *Expert Opinion on Drug Metabolism & Toxicology* 8(2), 135–139.
- 631 Nyamundanda, G., L. Brennan, and I. Gormley (2010). Probabilistic principal component
632 analysis for metabolomic data. *BMC Bioinformatics* 11(1), 571.
- 633 Nyamundanda, G., L. Brennan, and I. C. Gormley (2013). A random effects probabilistic
634 principal components model for longitudinal metabolomic data. Technical report, School
635 of Mathematical Sciences, University College Dublin.
- 636 Platanioti, K., E. McCoy, and D. Stephens (2005). A review of stochastic volatility: uni-
637 variate and multivariate models. Technical report, Imperial College London.
- 638 Ramoni, M. F., P. Sebastiani, and I. S. Kohane (2002). Cluster analysis of gene expression
639 dynamics. *PNAS* 99(14), 9121 – 9126.
- 640 Reo, N. V. (2002). Metabonomics based on NMR spectroscopy. *Drug and Chemical Toxi-*
641 *cology* 25(4), 375–382.
- 642 Sachse, D., L. Sletner, K. Mørkrid, A. K. Jennum, K. I. Birkeland, F. Rise, A. P. Piehler,
643 and J. P. Berg (2012). Metabolic changes in urine during and after pregnancy in a
644 large, multiethnic population-based cohort study of gestational diabetes. *PloS one* 7(12),
645 e52399.
- 646 Smilde, A., J. Jansen, H. Hoefsloot, S. Lamers R N, J. Greef, and M. Timmerman (2005).
647 ANOVA-simultaneous component analysis (ASCA): a new tool for analyzing designed
648 metabolomics data. *Bioinformatics* 21(13), 3043–3048.

- 649 Smilde, A., J. Westerhuis, H. Hoefsloot, S. Bijlsma, C. Rubingh, D. Vis, R. Jellema,
650 H. Pijl, and F. Roelfsema (2010). Dynamic metabolomic data analysis: a tutorial re-
651 view. *Metabolomics* 6(2), 3–17.
- 652 Smolinska, A., L. Blanchet, L. Buydens, and S. S. Wijmenga (2012). NMR and pattern
653 recognition methods in metabolomics: from data acquisition to biomarker discovery: a
654 review. *Analytica chimica acta* 750, 82–97.
- 655 Tipping, M. E. and C. M. Bishop (1999). Probabilistic principal component analysis. *Jour-
656 nal of the Royal Statistical Society, Series B* 61(3), 611–622.
- 657 van den Berg, R. A., H. C. J. Hoefsloot, J. A. Westerhuis, A. K. Smilde, and M. J. van der
658 Werf (2006). Centering, scaling, and transformations: improving the biological informa-
659 tion content of metabolomics data. *BMC Genomics* 7(1), 142.
- 660 Walsh, M., L. Brennan, E. Pujos-Guillot, J. Sébédio, A. Scalbert, A. Fagan, D. Hig-
661 gins, and M. Gibney (2007). Influence of acute phytochemical intake on human urinary
662 metabolomic profiles. *The American Journal of Clinical Nutrition* 86(6), 1687–1693.
- 663 Wang, Z., F. Yang, D. W. C. Ho, S. Swift, A. Tucker, and X. Liu (2008). Stochastic dynamic
664 modeling of short gene expression time-series data. *NanoBioscience, IEEE Transactions
665 on* 7(1), 44–55.
- 666 Wopereis, S., C. M. Rubingh, M. J. van Erk, E. R. Verheij, T. van Vliet, N. H. P. Cnubben,
667 A. K. Smilde, J. van der Greef, B. van Ommen, and H. F. J. Hendriks (2009). Metabolic
668 profiling of the response to an oral glucose tolerance test detects subtle metabolic changes.
669 *PLoS ONE* 4(2), e4525.