

BAYESIAN PARTIAL LEAST SQUARES (BPLS)

Szymon Urbaś

Maynooth University, Ireland

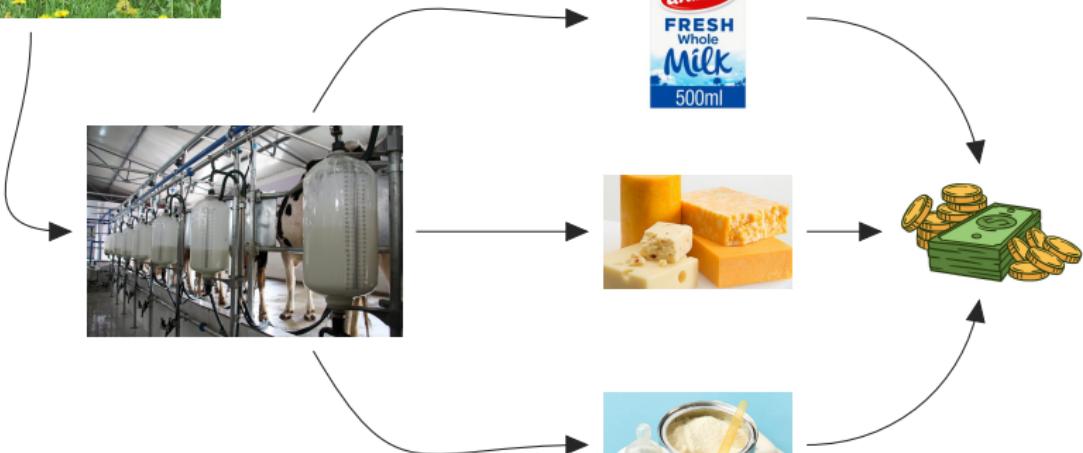
:: THANKS AND ACKNOWLEDGEMENTS ::

This joint work under VistaMilk SFI Research Centre.

- Prof Claire Gormley (*UCD*)
- Prof Donagh Berry (*Teagasc*)
- Dr Pierre Lovera, Dr Rob Daly
and Prof Alan O'Riordan
(*Tyndall Institute, Cork*)



MOTIVATING PROBLEM—PREDICTING TRAITS OF MILK



MOTIVATING PROBLEM—PREDICTING TRAITS OF MILK

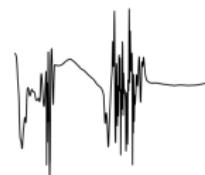
The quality and the value of a food product (e.g. milk) are determined by its chemical and technological traits (properties).

The process of measuring each of these in a lab can be **costly and time-consuming**.

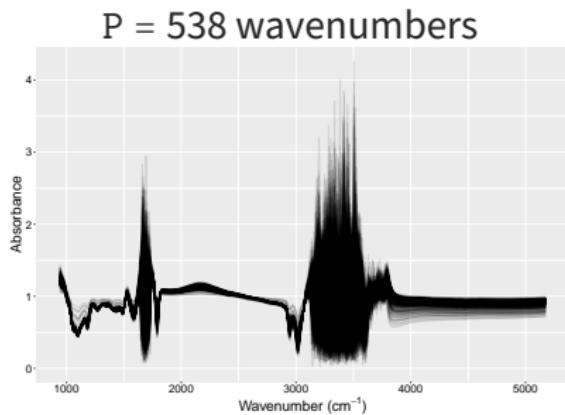
Spectrometry: examining how light at different wavelengths passes through the substance (**cheaper, quicker, non-destructive**)

Aim: Devise prediction models which can reliably predict the traits from a given spectral reading.

MOTIVATING PROBLEM—PREDICTING TRAITS OF MILK



MOTIVATING PROBLEM — MID-INFRARED SPECTRA OF MILK



Traits ($R = 3$) are measured in a lab:

- **heat stability**
 - for making infant formula;
- **rennet coagulation time (RCT)**
 - for making cheese;
- **casein content**
 - major protein.

We have $N = 363$ complete observations. (large P, small N problem)

PARTIAL LEAST SQUARES

Predictors: $(\mathbf{x}_1, \dots, \mathbf{x}_N)^\top = \mathbf{X} \in \mathbb{R}^{N \times P}$ **Responses:** $(\mathbf{y}_1, \dots, \mathbf{y}_N)^\top = \mathbf{Y} \in \mathbb{R}^{N \times R}$

Suppose there is a common projection onto a Q-dimensional latent space with variables $\mathbf{Z} \in \mathbb{R}^{N \times Q}$,

$$\mathbf{X} = \mathbf{Z}W^\top + \mathbf{E}_x,$$

$$\mathbf{Y} = \mathbf{Z}C^\top + \mathbf{E}_y,$$

where W and C are loading matrices, and \mathbf{E}_x and \mathbf{E}_y are residuals.

PLS: Iteratively find $\text{argmax } \|\text{Cov}(\mathbf{X}, \mathbf{Y})\|$.

- Use W and C estimates to predict traits from new sample spectra.

PARTIAL LEAST SQUARES



Fast, easy, accurate*, easy to modify the “regression” part

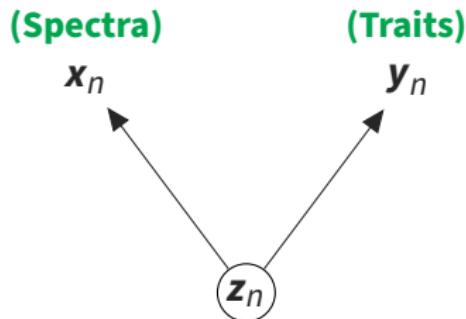
— Not a statistical model \implies no uncertainty, difficult to introduce sample correlations



Choosing Q can be tricky, strongly depends on the quality of data

BAYESIAN PARTIAL LEAST SQUARES

Each sample is generated through:



$$\mathbf{x}_n = W\mathbf{z}_n + \varepsilon_n, \quad n = 1, \dots, N,$$

$$\mathbf{y}_n = C\mathbf{z}_n + \boldsymbol{\eta}_n,$$

where $\mathbf{z}_n \stackrel{iid}{\sim} N(\mathbf{0}, I_Q)$, $\varepsilon_n \stackrel{iid}{\sim} N(\mathbf{0}, \Sigma)$ and $\eta_n \stackrel{iid}{\sim} N(\mathbf{0}, \Psi)$.

This allows for likelihood-based inference of model parameters:

$$p(\mathbf{X}, \mathbf{Y} | \Theta) = \prod_{n=1}^N f(\mathbf{x}_n, \mathbf{y}_n | \Theta),$$

where $\Theta = (\mathbf{Z}, W, C, \Sigma, \Psi, \dots)$ are all (unknown) model parameters.

BAYESIAN PARTIAL LEAST SQUARES

This allows for likelihood-based inference of the model:

- Frequentist methods for these types of models typically require $Q \leq R$ which greatly limits prediction utility.
- Bayesian inference with appropriate priors bypasses this constraint.
→ Bayesian partial least squares (**BPLS**)

Priors: Here we assume vague(ish) priors: apart from regularising model fitting, they won't heavily contribute to the final predictions.

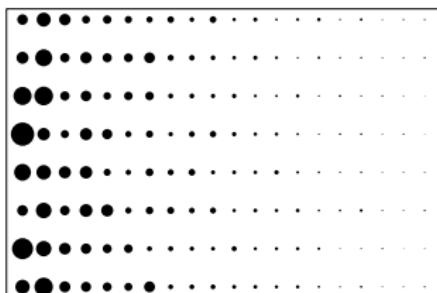
*Identifiability of Θ is not required for predictions.

SHRINKAGE ON THE LATENT VARIABLES

Structure of loading matrix:

→ Suppose Q is infinite:

$W =$



- Take a **stochastically increasing** sequence $\tau = (\tau_1, \tau_2, \dots)$.
- Elements of W are given conjugate normal priors such that

$$\mathbb{V}[w_{pq}] \propto \frac{1}{\tau_q}, \quad p=1, \dots, P, q=1, 2, \dots.$$

⇒ Subsequent latent components contribute less to the signal.

- **Multiplicative gamma process** prior:

$$\tau_q = \prod_{k=1}^q \delta_k, \text{ where } \delta_k \sim \text{Gamma}(\alpha, \beta).$$

SPARSITY, PRIORS AND INFERENCE

Large parts of the spectrum may tell you nothing about the traits.

Response part of BPLS: $\mathbf{y}_n = C\mathbf{z}_n + \boldsymbol{\eta}_n, n = 1, \dots, N$

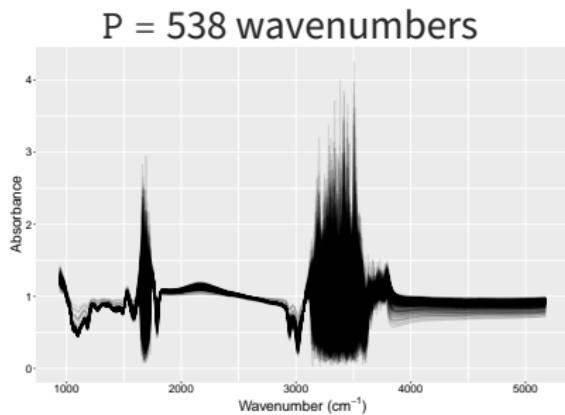
We consider two sparse variants:

- **Spike-and-slab** — sets some columns of C to be exactly zero (**ss-BPLS**)
→ Corresponds to Two-Way Orthogonal PLS (O2-PLS)
- **Bayesian LASSO** — emulates the ℓ_1 -penalty on elements of C (**L-BPLS**)
→ Corresponds to sparse PLS (sPLS)

Can assign conjugate priors everywhere → Gibbs sampling

Output: Posterior predictive distribution of $\mathbf{Y}^{\text{new}} | \mathbf{X}^{\text{new}}, \mathcal{D}$
(marginalised over parameter posterior)

MOTIVATING PROBLEM — MID-INFRARED SPECTRA OF MILK

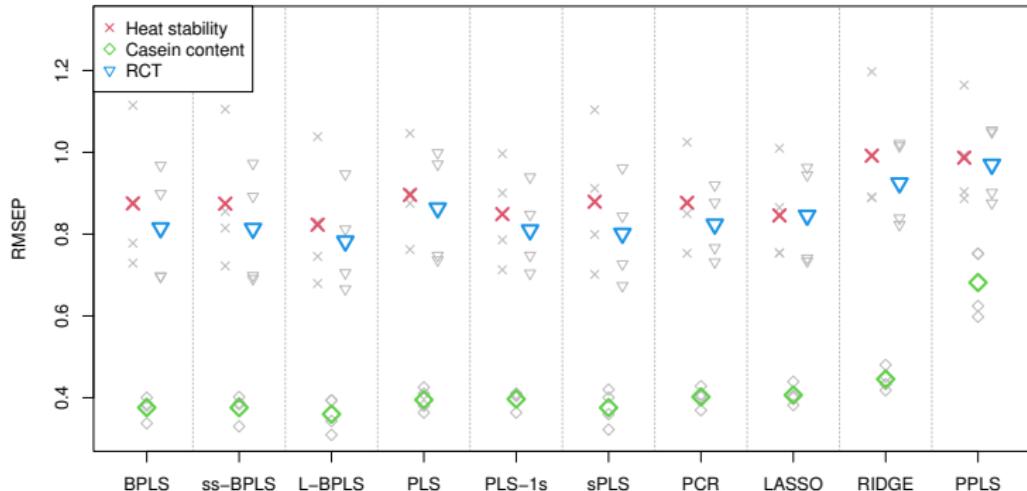


Traits ($R = 3$) are measured in a lab:

- **heat stability**
 - for making infant formula;
- **rennet coagulation time (RCT)**
 - for making cheese;
- **casein content**
 - major protein.

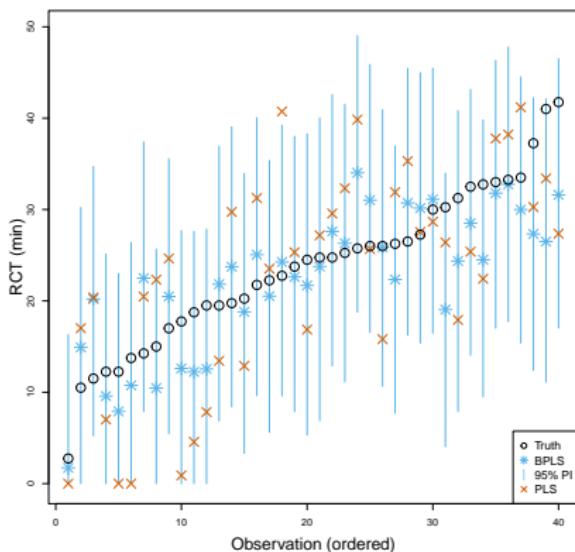
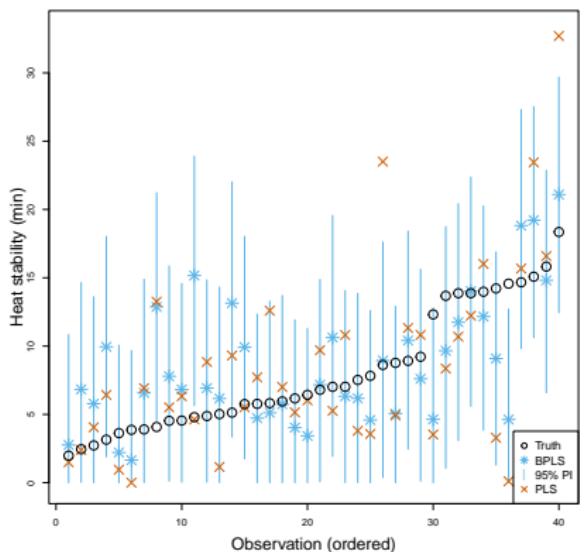
We have $N = 363$ complete observations. (large P, small N problem)

MILK MIR SPECTRAL DATA RESULTS



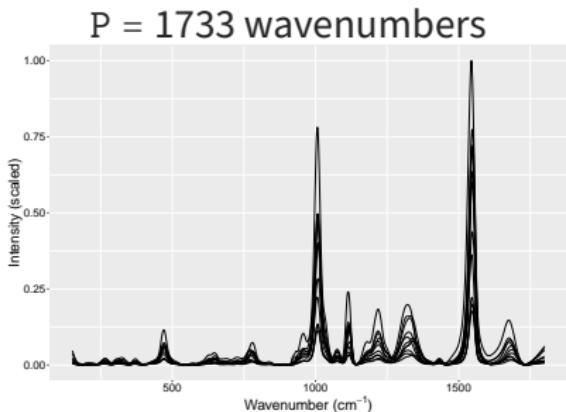
- **L-BPLS** is consistently the most accurate method.
- Statistical model outputs have a lot more utility.

MILK MIR SPECTRAL DATA RESULTS



MOTIVATING PROBLEM — SERS-BASED PH SENSORS

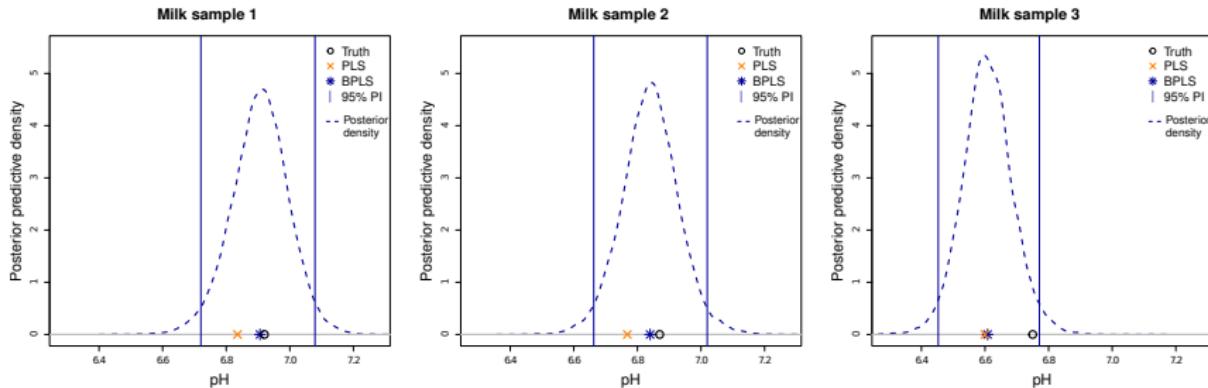
pH sensors using surface-enhanced Raman spectroscopy
→ detect product spoilage or indicate mastitis.



The pH of two cartons of milk is measured in a lab over 6 days:

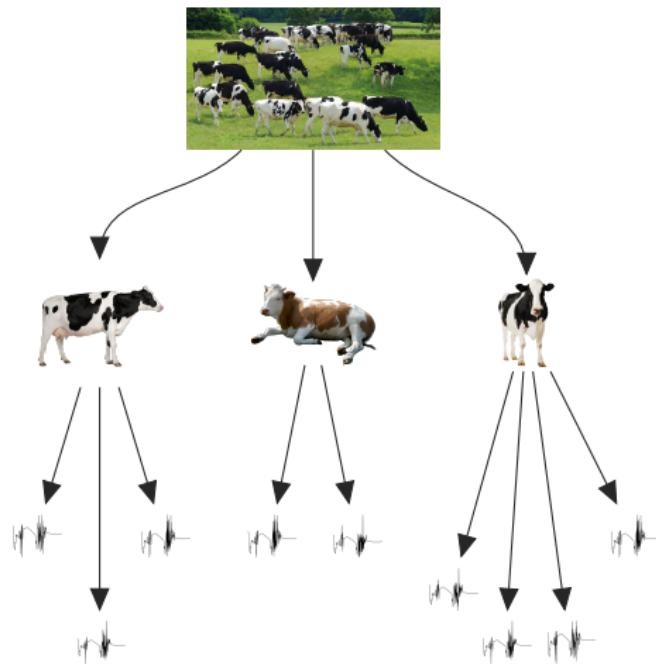
- Very small dataset ($N = 11$) — can we recover any signal?

MILK SERS DATA RESULTS



- Standard PLS methods fail to identify signal + cross-validation methods used for finding Q really struggle here.
- BPLS methods manage to produce reasonable point predictions but highlight the uncertainty due to the small dataset.

TOWARDS HIERACHICAL MODELLING



Paper:

Urbas, S., Lovera, P., Daly, R., O'Riordan, A., Berry, D., & Gormley, I. C. (2024). Predicting milk traits from spectral data using Bayesian probabilistic partial least squares regression. *The Annals of Applied Statistics*, 18(4), 3486–3506.

Code:

- `bplsr` package available on CRAN.

SPARSITY — SPIKE-AND-SLAB

- **Idea:** large components explaining variation in \mathbf{X} may play no part in explaining the variation in \mathbf{Y}
 - In the model, this implies that some columns of \mathbf{C} may be zero.
 - **Spike and slab approach:** For each column introduce Bernoulli variables b_q , $q = 1, 2, \dots$ which can “switch” the latent variable on and off.

With $B = \text{diag}(b_1, b_2, \dots)$, the response part of the model becomes

$$\mathbf{y}_n = CB\mathbf{z}_n + \boldsymbol{\varepsilon}_n.$$

Elements of B can be inferred via posterior Gibbs sampling.

C =

SPARSITY — BAYESIAN LASSO

- Explicit prior of the form $\pi_0(C) \propto \exp(-\lambda \sum_{r,q} |c_{rq}|)$, $\lambda > 0$ results in an intractable and non-differentiable posterior
- Park and Casella (2009) uses a scale mixture of normal distributions to get a similar form:

$$\frac{\lambda}{2} e^{-\lambda|c|} = \int_0^\infty \frac{1}{\sqrt{2\pi\nu}} e^{c^2/(2\nu)} \times \frac{\lambda^2}{2} e^{\lambda^2\nu/2} d\nu.$$

- So

$$c_{rq} | \tau_q, \nu_{rq} \sim N(0, \nu_{rq}/\tau_q), \quad \nu_{rq} | \lambda \sim \text{Exp}(\lambda^2/2), \quad r = 1, \dots, R, \quad q = 1, 2, \dots$$

gives us what we want.

- Posterior conditionals of ν_{rq} are inverse-Gaussian so can still use Gibbs