

Gaussian Processes with Bayesian Inference of Covariate Couplings

Mattia Rosso *KAUST, Saudi Arabia*

Juho Ylä-Jääski *Aalto University, Finland*

Zheyang Shen *Newcastle University, UK*

Markus Heinonen *Aalto University, Finland*

Maurizio Filippone *KAUST, Saudi Arabia*

Reviewed on OpenReview: <https://openreview.net/forum?id=fameEAljo3>

Abstract

Gaussian processes are powerful probabilistic models that are often coupled with Automatic Relevance Determination (ARD) capable of uncovering the importance of individual covariates. We develop covariances characterized by affine transformations of the inputs, formalized via a precision matrix between covariates, which can uncover covariate couplings for enhanced interpretability. We study a range of couplings priors from Wishart to Horseshoe and present fully Bayesian inference of such precision matrices within sparse Gaussian processes. We empirically demonstrate the efficacy and interpretability of this approach.

1 Introduction

Statistical models based on Gaussian Processes (GPs) offer attractive modeling choices for various quantitative sciences due to their ability to impose functional priors with certain desired characteristics and to carry out principled uncertainty quantification (Rasmussen & Williams, 2006). Modeling and inference of GPs has evolved significantly in the directions of scalability for large data (Cutajar et al., 2017; Hensman et al., 2013; Wilson & Nickisch, 2015), deep learning (Damianou & Lawrence, 2013; Wilson et al., 2016; Salimbeni & Deisenroth, 2017), and generality with autodiff frameworks (Krauth et al., 2017; Matthews et al., 2017).

The choice of the covariance (kernel) function plays a crucial role in specifying the function space induced by GPs. This choice is often overlooked by opting for the reputable “default” exponential ARD covariances (Neal, 1996), which capture the importance of each covariate, but also assumes an *axis-aligned* anisotropic data structure, blind to covariate couplings (Matérn, 1960).

In contrast, *affine* anisotropic covariances are able to explicitly consider the linear dependencies between covariates (Matérn, 1960; Poggio & Girosi, 1990), which is a common feature of real-world data, via the *precision* matrix $\mathbf{\Lambda}$ of the Mahalanobis distance $(\mathbf{x} - \mathbf{x}')^\top \mathbf{\Lambda} (\mathbf{x} - \mathbf{x}')$. A seminal work of Vivarelli & Williams (1998) proposes a parameterization of the precision based on Principal Component Analysis (PCA), while Titsias & Lazaro-Gredilla (2013) apply mean-field variational inference over factors of such a precision matrix $\mathbf{\Lambda}$. Relevant works on affine-covariances GPs include non-stationary extensions (Paciorek & Schervish, 2003), and applications to imaging (Kalaitzis, 2009) and material sciences (Noack et al., 2020).

In this paper, our goal is to revitalise Mahalanobis distance-based covariances as a more interpretable and general alternative to “diagonal” ARD covariances, whereby we are able to uncover covariate couplings. This is illustrated in Fig. 1, where we refer to these more general types of covariance functions as Automatic Coupling Determination (ACD) covariances. Unlike previous works considering full precision matrices $\mathbf{\Lambda}$,

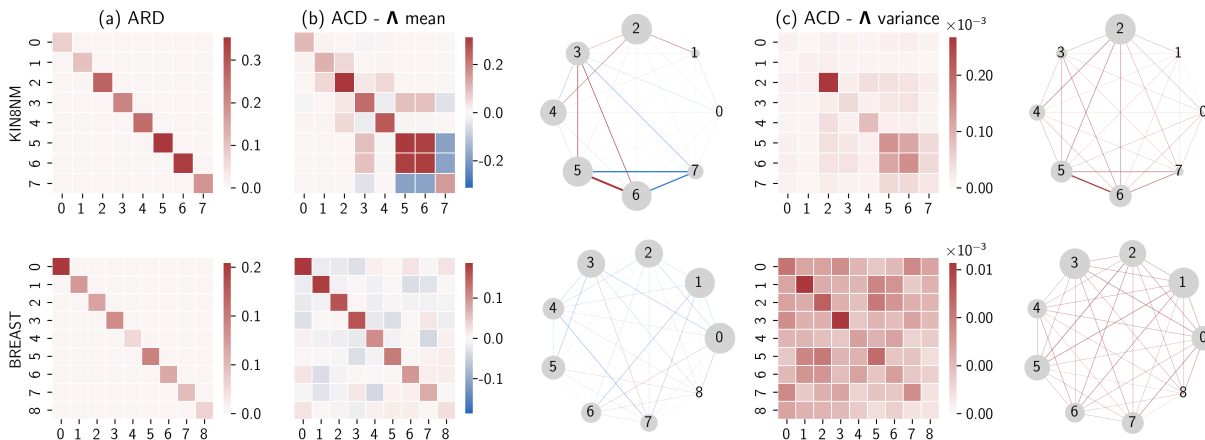


Figure 1: The Automatic Coupling Determination (ACD) covariance reveals rich predictive covariate couplings. Comparison between ARD diagonal precisions $\Sigma^{-1} = \text{diag}(\ell^{-2})$ (a) and ACD precision matrix Λ mean (b) and variance (c) with graph illustrations. We assume an element-wise Normal prior on Λ . The ACD covariance detects that the covariates (5,6,7) are close to redundant on the `kin8nm` dataset.

we study a fully Bayesian scalable formulation of GPs, where we carry out inference over the matrix Λ , thus obtaining posterior distributions over covariate couplings. An attractive feature of this approach in supervised learning problems is the possibility of obtaining information about covariate couplings that are instrumental in yielding accurate modeling of the labels.

Our contributions are as follows: (i) a GP model that determines covariate couplings through the analysis of the matrix Λ ; (ii) an analysis of sparsity-inducing priors for the matrix Λ from Wishart, Laplace and Horseshoe families; (iii) a demonstration of the enhanced explainability of ACD covariances compared to ARD covariances; (iv) the development of a fully Bayesian Markov chain Monte Carlo (MCMC) inference scheme of the couplings, while operating in a scalable sparse GP framework; and (v) an empirical demonstration of the usefulness of ACD covariances.

2 Background

We consider supervised learning problems with N input-label pairs $\{\mathbf{X}, \mathbf{y}\} = \{(\mathbf{x}_n, y_n)\}_{n=1}^N$, with $\mathbf{x}_n \in \mathbb{R}^D$ and $y_n \in \mathbb{R}$. We model the labels through a latent function $f(\mathbf{x})$, for which we assume a GP prior (Rasmussen & Williams, 2006). Let $\mathbf{f} = (f(\mathbf{x}_1), \dots, f(\mathbf{x}_N))^T \in \mathbb{R}^N$ denote the collection of latent variables associated with inputs $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)$, and $\prod_{n=1}^N p(y_n | f(\mathbf{x}_n))$ be the likelihood function, for which we made an i.i.d assumption on the realization of the labels y_n conditioned on the corresponding $f(\mathbf{x}_n)$.

2.1 Gaussian process priors

By imposing a GP prior $f(\mathbf{x}) \sim \mathcal{GP}(0, k)$ on the latent function $f(\mathbf{x})$, we are assuming that any subset of these random variables are jointly Gaussian (Rasmussen & Williams, 2006). The kernel function $k(\mathbf{x}, \mathbf{x}'; \boldsymbol{\theta})$ determines the properties of the functions that can be drawn from the GP prior, where $\boldsymbol{\theta}$ are hyper-parameters. The prior over the realizations of $f(\mathbf{x})$ at the inputs \mathbf{X} is then $p(\mathbf{f} | \boldsymbol{\theta}) = \mathcal{N}(\mathbf{0}, \mathbf{K}_{\mathbf{xx} | \boldsymbol{\theta}})$, where $\mathbf{K}_{\mathbf{xx} | \boldsymbol{\theta}}$ is the $N \times N$ covariance matrix obtained by evaluating $k(\mathbf{x}, \mathbf{x}'; \boldsymbol{\theta})$ at all input pairs $\{\mathbf{x}, \mathbf{x}'\}$. For simplicity, we assume zero-mean GPs and omit the conditioning on \mathbf{X} .

The posteriors over \mathbf{f} at inputs \mathbf{x}_* , and inference or optimization over $\boldsymbol{\theta}$ is based on the analysis of the joint

$$p(\mathbf{y}, \mathbf{f}, \boldsymbol{\theta}) = p(\mathbf{y} | \mathbf{f}) p(\mathbf{f} | \boldsymbol{\theta}) p(\boldsymbol{\theta}). \quad (1)$$

With Gaussian likelihoods it is possible to marginalize out \mathbf{f} leading to a tractable $p(\mathbf{y}, \boldsymbol{\theta}) = p(\mathbf{y} | \boldsymbol{\theta}) p(\boldsymbol{\theta})$. With non-Gaussian likelihoods, further complications arise due to the lack of conjugacy (Williams & Barber,

1998; Opper & Winther, 2000). For non-conjugate GP models, popular approaches to sample from the posterior distribution of covariance parameters and latent variables include Murray & Adams (2010) (slice sampling) and Filippone & Girolami (2014) (pseudo-marginal MCMC).

An overarching issue with GP models is scalability, as these models generally require costly $\mathcal{O}(N^3)$ operations involving $\mathbf{K}_{\mathbf{xx}|\boldsymbol{\theta}}$ inverses. Linearization techniques based on random features (Rahimi & Recht, 2008) were proposed in Lázaro-Gredilla et al. (2010), and they were later developed to operate with mini-batches within stochastic gradient optimization and to deep GPs (Cutajar et al., 2017). In a parallel line of work, sparsification techniques based on inducing points (Williams & Seeger, 2000; Snelson & Ghahramani, 2005) were later embedded within a variational formulation (Titsias, 2009), and they were extended to mini-batching (Hensman et al., 2013; Krauth et al., 2017). In this paper, we consider sparse GPs, and in particular their fully Bayesian version presented in Rossi et al. (2021), where all variables are treated in a Bayesian way and inference is carried out using stochastic gradient MCMC (Chen & Zhang, 2004).

2.2 Fully Bayesian sparse GPs

We focus on the Bayesian sparse Gaussian process (BSGP) framework (Rossi et al., 2021), but the ACD covariance specifications generally apply to any GP implementations. In sparse GPs, we introduce a set of M inducing variables $\mathbf{u} = (u_1, \dots, u_M)$, which denote the realization of the latent function $f(\mathbf{x})$ at inducing inputs $\mathbf{Z} = \{\mathbf{z}_1, \dots, \mathbf{z}_M\}$, such that $u_m = f(\mathbf{z}_m)$ (Candela & Rasmussen, 2005). The definition of the inducing variables \mathbf{u} implies that they are multivariate Gaussian with covariance depending on the inducing inputs \mathbf{Z} . In summary, the GP prior assumption on the latent function yields the following prior specifications for the latent variables \mathbf{f} and \mathbf{u} :

$$p(\mathbf{f}, \mathbf{u} | \mathbf{X}, \mathbf{Z}, \boldsymbol{\theta}) = p(\mathbf{f} | \mathbf{u}, \mathbf{X}, \mathbf{Z}, \boldsymbol{\theta}) p(\mathbf{u} | \mathbf{Z}, \boldsymbol{\theta}) \quad (2)$$

$$p(\mathbf{u} | \mathbf{Z}, \boldsymbol{\theta}) \sim \mathcal{N}(\mathbf{0}, \mathbf{K}_{\mathbf{zz}|\boldsymbol{\theta}}) \quad (3)$$

$$p(\mathbf{f} | \mathbf{u}, \mathbf{X}, \mathbf{Z}, \boldsymbol{\theta}) \sim \mathcal{N}(\mathbf{A}\mathbf{u}, \mathbf{K}_{\mathbf{xx}|\boldsymbol{\theta}} - \mathbf{A}\mathbf{K}_{\mathbf{zz}|\boldsymbol{\theta}}^\top), \quad (4)$$

where $\mathbf{A} = \mathbf{K}_{\mathbf{xz}|\boldsymbol{\theta}} \mathbf{K}_{\mathbf{zz}|\boldsymbol{\theta}}^{-1}$. This augmented model can be used for modeling tasks by introducing a likelihood $p(\mathbf{y} | \mathbf{f})$. We assign priors over all remaining variables $p_\psi(\boldsymbol{\theta})$ and $p_\xi(\mathbf{Z})$, notably including inducing locations \mathbf{Z} and covariance parameters $\boldsymbol{\theta}$ (Rossi et al., 2021). Here we assume that the priors over $\boldsymbol{\theta}$ and \mathbf{Z} are conditionally independent, but this can be relaxed if one intends to couple \mathbf{Z} and $\boldsymbol{\theta}$ in the prior. The joint becomes

$$p(\boldsymbol{\theta}, \mathbf{Z}, \mathbf{u}, \mathbf{f}, \mathbf{y} | \mathbf{X}) = p_\psi(\boldsymbol{\theta}) p_\xi(\mathbf{Z}) p(\mathbf{f}, \mathbf{u} | \mathbf{X}, \mathbf{Z}, \boldsymbol{\theta}) p(\mathbf{y} | \mathbf{f}). \quad (5)$$

Inference over latent variables and hyper-parameters is intractable, and variational inference allows one to recover tractability. The introduction of inducing variables \mathbf{u} in the definition of the model enables a scalable formulation leading to an objective that factorizes across data, while requiring algebraic operations with the $M \times M$ matrix $\mathbf{K}_{\mathbf{zz}|\boldsymbol{\theta}}$ instead of $\mathbf{K}_{\mathbf{xx}|\boldsymbol{\theta}}$. Following Rossi et al. (2021), who study the differences between a fully Bayesian treatment of the variational free energy (VFE) and fully independent training conditional (FITC) approximations, we adopt the latter to approximate the conditional $p(\mathbf{f} | \mathbf{u}, \mathbf{X}, \mathbf{Z}, \boldsymbol{\theta})$. This allows for parameter inference over $\boldsymbol{\Psi} \stackrel{\text{def}}{=} \{\mathbf{u}, \mathbf{Z}, \boldsymbol{\theta}\}$ with scalable MCMC based on stochastic gradients (Chen & Zhang, 2004).

2.3 Computational Complexity

In the sparse GP approach the matrix factorization of \mathbf{K}_{uu} has complexity $\mathcal{O}(M^3)$, where M is the number of inducing points. Additionally, multiplying \mathbf{K}_{fu} with the inverse (or performing forward/backward substitution with the Cholesky decomposition) of \mathbf{K}_{uu} incurs a complexity of $\mathcal{O}(NM^2)$ where N is the number of data points. In the MCMC sampling method proposed by Rossi et al. (2021), mini-batching reduces the complexity of the $\mathcal{O}(NM^2)$ operation to $\mathcal{O}(N'M^2)$, where $N' \ll N$, significantly improving computational efficiency. To be more specific, the complexity for $\mathbf{K}_{uu} \in \mathbb{R}^{M \times M}$ is $\mathcal{O}(M^2 D^2)$ since each of the M^2 entries involves matrix-vector products with the $\boldsymbol{\Lambda}$ matrix. Similarly, for the $\mathbf{K}_{uf} \in \mathbb{R}^{M \times N}$ matrix the complexity is $\mathcal{O}(MND^2)$.

3 Related Works

3.1 Gaussian processes with Automatic Coupling Determination

The possibility of performing kernel-based modeling with the determination of the importance of inputs dates at least back to work on Automatic Relevance Determination ARD (MacKay, 1995; Neal, 1996). This is usually implemented by scaling input covariates within the calculation of the covariance function by some coefficients which are treated as hyper-parameters and optimized through marginal likelihood optimization.

An extension of this idea involves the use of affine transformations (rotation and stretching) of the covariates; in distance-based covariance functions, the affine transformation implies the calculation of the so-called Mahalanobis distance (Vivarelli & Williams, 1998; Titsias & Lazaro-Gredilla, 2013). In (Vivarelli & Williams, 1998), Λ is made positive definite by construction through the parameterization $\mathbf{U}\mathbf{U}^\top$ with \mathbf{U} upper triangular, and it is factorized to gain insights into the dimensionality of a possible low-dimensional latent representation of the inputs. In our work, we consider a similar parameterization for Λ , but instead of optimizing its factors, we carry out a Bayesian treatment, for which we study sparsity-inducing priors. Also, we propose a PCA-based decomposition of Λ , which allows us to operate in large-dimensional input regimes; this is done through the first d principal components \mathbf{P}_d of the input covariance, which we use to express $\Lambda = \mathbf{P}_d\Lambda_d\mathbf{P}_d$.

The work by (Titsias & Lazaro-Gredilla, 2013) considers the parameterization $\Lambda = \mathbf{W}^\top\mathbf{W}$ without imposing any structure on \mathbf{W} except for imposing that \mathbf{W} maps the input covariates to a lower dimensional input space. Their work proposes a variational formulation to obtain an approximation to the posterior over Λ but it does not extensively study priors over \mathbf{W} , which is the focus of our work.

3.2 Sparsity-Inducing priors for covariance/precision matrices

The literature on Gaussian Graphical Model (GGM) provides studies on sparsity-inducing priors for covariance and precision matrices. Sparsity can be imposed in a structured fashion by considering graph decomposability (Banerjee & Ghosal, 2014; Lee & Lee, 2021; Xiang et al., 2015; Banerjee et al., 2021), or through the G-Wishart prior, which has been introduced as a conjugate prior for the precision matrix in a Gaussian framework and it is also suitable for cases where graph decomposability does not apply (Roverato, 2000; 2002; Khare & Rajaratnam, 2011; Silva & Ghahramani, 2009; van den Boom et al., 2022).

Various approaches have been developed to carry out inference in GGMS, including Gibbs sampling (Khare & Rajaratnam, 2011; Wang, 2012) and Laplace approximations (Banerjee & Ghosal, 2015). Other approaches, such as those by Gan et al. (2018); Wang (2015), propose spike-and-slab sparsity-inducing priors which typically complicate posterior sampling. Castillo et al. (2015); Li et al. (2017); Sagar et al. (2024) study horseshoe priors, which perform well in practice.

Our work differs from this literature, given that we propose a model for the labels given the inputs, while attempting to uncover some couplings among covariates. The way this is done is by a parameterization of the covariance function akin to the precision matrix in a GGM, and in our work we explore both matrix-variate and element-wise priors for such model parameters. In addition, we consider scalable sampling-based approaches to obtain samples from the posterior distribution over these parameters.

4 Bayesian inference of covariate couplings

In this section, after briefly discussing covariances with Automatic Relevance Determination (ARD) (MacKay, 1995; Neal, 1996), which induce some scaling of individual covariates, we present an extension involving an affine transformation of the covariates revealing couplings among these. We discuss how this is achieved by introducing a Mahalanobis distance among inputs with a precision matrix Λ , and we show how to treat this in a Bayesian way by imposing matrix-variate and sparsity-inducing element-wise priors. We term this type of covariance Automatic Coupling Determination (ACD).

4.1 Automatic relevance determination

The design of covariance functions for GP models is an important part of the modeling process. Considering the space of functions $f : \mathbb{R}^D \mapsto \mathbb{R}$, the choice of a covariance $\text{cov}[f(\mathbf{x}), f(\mathbf{x}')] = k(\mathbf{x}, \mathbf{x}'; \boldsymbol{\theta})$ determines the prior distribution over f before observing data. A common choice is the Gaussian covariance function (Radial Basis Function (RBF)):

$$k_{\text{RBF}}(\mathbf{x}, \mathbf{x}'; \boldsymbol{\theta}) \propto \exp\left(-\frac{1}{2}d^2(\mathbf{x}, \mathbf{x}'; \boldsymbol{\theta})\right), \quad (6)$$

where $d(\mathbf{x}, \mathbf{x}'; \boldsymbol{\theta})$ is a parametric distance function between inputs \mathbf{x} and \mathbf{x}' . This covariance imposes a prior over infinitely differentiable (smooth) functions. Other common covariance functions based on the distance $d(\mathbf{x}, \mathbf{x}'; \boldsymbol{\theta})$ include the Matérn covariance, exponential, arc-cosine; see, e.g., Shawe-Taylor & Cristianini (2004) for an in-depth treatment.

The simplest distance form

$$d_{\text{ISOTROPIC}}^2(\mathbf{x}, \mathbf{x}'; \boldsymbol{\theta}) = \frac{1}{\ell^2}(\mathbf{x} - \mathbf{x}')^\top (\mathbf{x} - \mathbf{x}') \quad (7)$$

induces an isotropic covariance, as all input features are scaled by the same length-scale parameter ℓ and contribute equally to the distance, which assumes *spherical* data.

Another choice increasing model flexibility introduces covariate-specific length-scales parameters,

$$d_{\text{ARD}}^2(\mathbf{x}, \mathbf{x}'; \boldsymbol{\theta}) = (\mathbf{x} - \mathbf{x}')^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \mathbf{x}') \quad (8)$$

with $\boldsymbol{\Sigma} = \text{diag}(\ell_1^2, \dots, \ell_D^2)$. This choice gives rise to covariances suitable for ARD (MacKay, 1995; Neal, 1996).

Intuitively, this definition is built on the assumption that, if a dimension d has a small value of the associated length-scale ℓ_d , small changes in the covariate lead to large responses in the target. The covariance induced by this choice is anisotropic with an axis-aligned metric acting as a scaling of individual covariates.

4.2 Automatic coupling determination

The family of ARD covariances allows GP models to yield non-parametric and probabilistic mappings from inputs to labels, while simultaneously determining the importance of each covariate if ℓ_d 's are optimized or inferred. In this paper, we do not limit ourselves to assessing the relevance of each input covariate, but to automatically discover couplings among these in a general way which can be readily applied to any distance-based covariance function.

We replace the diagonal matrix $\boldsymbol{\Sigma}^{-1}$ containing the inverse length-scales with a full Positive Semi-Definite (PSD) precision matrix $\boldsymbol{\Lambda} = \boldsymbol{\Sigma}^{-1}$ in the calculation of distances,

$$d_{\text{ACD}}^2(\mathbf{x}, \mathbf{x}'; \boldsymbol{\theta}) = (\mathbf{x} - \mathbf{x}')^\top \boldsymbol{\Lambda}(\mathbf{x} - \mathbf{x}') = \sum_{i,j}^D \Lambda_{ij}(x_i - x'_i)(x_j - x'_j), \quad (9)$$

yielding the so-called Mahalanobis distance (Titsias & Lazaro-Gredilla, 2013), which can be interpreted as a distance obtained after an affine transformation (rotation and stretching) of the inputs by the identity (Matérn, 1960; Vivarelli & Williams, 1998; Kalaitzis, 2009)

$$d(\boldsymbol{\Lambda}^{\frac{1}{2}}\mathbf{x}, \boldsymbol{\Lambda}^{\frac{1}{2}}\mathbf{x}'; \mathbf{I}) = d(\tilde{\mathbf{x}}, \tilde{\mathbf{x}}'; \boldsymbol{\Lambda}). \quad (10)$$

If the underlying distribution of the inputs \mathbf{x} is Gaussian, this operation produces an implicit *whitening* of the input data yielding $\tilde{\mathbf{x}}$. While the quadratic form in Eq. 9 has an additive form, the induced functions do not lend themselves to an additive function interpretation (Vivarelli & Williams, 1998; Duvenaud et al., 2011). Note that if the precision matrix has zero elements $\Lambda_{ij} = 0$, the distance function ignores the coupling between covariates i and j in the calculation of pairwise distances among inputs.

Discriminative vs Generative modeling The parameterization of ACD covariances has apparent connections with Markov Random Fields (MRFs) (see, e.g., Murphy (2023)), where the matrix $\mathbf{\Lambda}$ is used to specify an adjacency structure for a set of D random variables $\{X_1, \dots, X_D\}$. MRFs offer the possibility to verify conditional independence properties of groups of random variables based on the analysis of $\mathbf{\Lambda}$, while placing no other assumptions on their underlying distribution. While it is tempting to think of the ACD parameterization of the covariance function as something to be used to draw conclusions on conditional independence among covariates, we are effectively not modeling their distribution. Instead, we are pushing $\mathbf{\Lambda}$ directly in the definition of the GP prior $p(\mathbf{f}|\mathbf{\Lambda})$. Therefore $\mathbf{\Lambda}$ assumes the interpretation of a precision matrix inducing an affine transformation of the input, which is optimized or inferred based on the marginal likelihood (or a lower bound thereof). Thus the focus is on performing optimization or inference of $\mathbf{\Lambda}$ to accurately modeling the labels, with the intention of obtaining some indication of the predictive power of couplings of covariates. We leave the modeling of the input through MRFs as an interesting avenue for future work.

4.3 Precision parameterizations

The precision matrix $\mathbf{\Lambda}$ in the ACD covariance needs to be symmetric and PSD. The PSD constraint in the ARD covariance is easy to satisfy, since working with a diagonal version of $\mathbf{\Lambda}$ only requires to have non-negative elements on its diagonal and consequently, a log-transformation of the length-scales is sufficient.

4.3.1 Lower triangular factorization

In the case of the ACD covariance, optimization or inference of $\mathbf{\Lambda}$ needs to be performed while preserving the PSD constraint, so that it is straightforward to operate with unconstrained optimization/MCMC sampling. Among all factorizations that can be used to express $\mathbf{\Lambda}$, following (Kalaitzis, 2009), a natural parameterization is via the lower-triangular matrix \mathbf{L} ,

$$\mathbf{\Lambda} = \mathbf{L}\mathbf{L}^\top. \quad (11)$$

This allows us to directly optimize or sample \mathbf{L} element-wise and recover $\mathbf{\Lambda}$. In addition, this parameterization has some computational advantages in calculating Jacobians which are useful within MCMC, as discussed shortly.

Note that, due to rotation invariance, the Cholesky factorization is one of infinitely many possible decompositions of $\mathbf{\Lambda}$. Rotation invariance implies that the product of two factors, \mathbf{L} and \mathbf{L}^\top in our case, gives the same $\mathbf{\Lambda}$ regardless of any rotation (multiplication with a unitary matrix) of the factors. In other words, given any unitary matrix \mathbf{U} , $\mathbf{L}\mathbf{U}$ is also a factor of $\mathbf{\Lambda}$. This argument indicates that even if we pick a given ordering of the covariates and the Cholesky factor \mathbf{L} is a factor of $\mathbf{\Lambda}$, then we would obtain the same $\mathbf{\Lambda}$ if we applied any unitary transformation (which includes permutations) to \mathbf{U} .

4.3.2 Computational Complexity of the ACD parameterization

Evaluating each element of \mathbf{K}_{uu} and \mathbf{K}_{uf} in the sparse GP framework requires $O(D^2)$ computations for ACD covariance functions. This is the case when evaluating element-wise priors and matrix variate priors on $\mathbf{\Lambda}$. To be more specific, the complexity for $\mathbf{K}_{uu} \in \mathbb{R}^{M \times M}$ is $O(M^2 D^2)$ since each of the M^2 entries involves matrix-vector products with the $\mathbf{\Lambda}$ matrix. Similarly, for the $\mathbf{K}_{uf} \in \mathbb{R}^{M \times N}$ matrix the complexity is $O(MND^2)$. The quadratic scaling in D motivates the adoption of low-rank parameterizations of $\mathbf{\Lambda}$, which we discuss in the next subsection. Note also that for ARD covariances these costs reduce to linear in D .

4.3.3 Low-rank factorizations

The increased flexibility offered by the ACD formulation comes at a computational cost, which we need to deal with: going from learning $\ell \in \mathbb{R}^D$ length-scales in the ARD covariance to learning a full $\mathbf{\Lambda} \in \mathbb{R}^{D \times D}$ matrix. This is why, for problems where the dimensionality D is high but we are at the same time interested in obtaining an informative precision matrix recovering the underlying structure among the D features, we tackle this problem with PCA, similarly to Vivarelli & Williams (1998) and Paciorek & Schervish (2003). In this context, PCA serves as a natural way to obtain a low rank transformation of the covariates by retaining a

Table 1: Summary of precision priors and the range of hyperparameters studied.

| Prior | Definition | Parameters | Log pdf |
|-----------------|---|--|--|
| Wishart | $p(\mathbf{\Lambda}) = \mathcal{W}(\mathbf{V}, K)$ | $K = D, \mathbf{V} = K^{-1}\mathbf{I}_D$ | $\log C - \sum_d \log \ \mathbf{L}_{dd}\ - \frac{1}{2}\text{Tr}[K\mathbf{\Lambda}]$ |
| Inverse Wishart | $p(\mathbf{\Lambda}) = \mathcal{IW}(\mathbf{V}, K)$ | $K = D, \mathbf{V} = \mathbf{I}_D$ | $\log C - (2K + 1) \sum_d \log \mathbf{L}_{dd} - \frac{1}{2}\text{Tr}[\mathbf{V}\mathbf{\Lambda}^{-1}]$ |
| Laplace | $p(\mathbf{\Lambda}_{ij}) = \mathcal{L}(m, b)$ | $m = 0, b \in \{0.01, 0.1, 1\}$ | $\log C - \frac{1}{b} \ \mathbf{\Lambda}_{ij} - m\ _1$ |
| Horseshoe | $p(\mathbf{\Lambda}_{ij}) = \mathcal{HS}(\tau)$ | $\tau \in \{0.01, 0.1, 1\}$ | $\log C + \frac{1}{2\tau^2} \mathbf{\Lambda}_{ij}^2 + \log E_1(\frac{1}{2\tau^2} \mathbf{\Lambda}_{ij}^2)$ |

given amount of explained variance. Focusing on the ACD distance, by applying a projection to the difference between data samples, we obtain:

$$(\mathbf{x} - \mathbf{x}')^\top \mathbf{P}_d \mathbf{\Lambda}_d \mathbf{P}_d^\top (\mathbf{x} - \mathbf{x}') \quad (12)$$

where \mathbf{P}_d is the $\mathbb{R}^{D \times d}$ matrix obtained from the eigendecomposition of the empirical covariance matrix

$$\mathbf{\Sigma} = \frac{1}{N} \mathbf{X}_c^\top \mathbf{X}_c = \mathbf{P} \mathbf{S} \mathbf{P}^\top, \quad (13)$$

and \mathbf{X}_c is the centered $\mathbb{R}^{N \times D}$ input matrix. To obtain \mathbf{P}_d we select the $d < D$ columns of \mathbf{P} corresponding to the d highest eigenvalues from \mathbf{S} . A sample $\mathbf{x} \in \mathbb{R}^D$ can be projected down to \mathbb{R}^d through $\mathbf{P}_d^\top \mathbf{x}$. As a result, we learn a projected version $\mathbf{\Lambda}_d$ in this latent representation of the full precision matrix. By applying the projection in Eq. 12 we recover the precision matrix in the original space. Note that in this parameterization, even an ARD model for $\mathbf{\Lambda}_d$ would lead to a full precision $\mathbf{P}_d \mathbf{\Lambda}_d \mathbf{P}_d^\top$, but given the favorable computational scaling of this representation, we prefer the added flexibility offered by an ACD. In Fig. 6 we report a comparison between the ACD and ARD parameterizations for $\mathbf{\Lambda}_d$.

5 Priors over $\mathbf{\Lambda}$

As a consequence of adopting the BSGP framework we need to specify a prior $p_\psi(\boldsymbol{\theta})$ over covariance hyperparameters $\boldsymbol{\theta}$. Dealing with the ACD covariance, the prior is separately placed over both the marginal variance parameter σ_f^2 and on the precision matrix $\mathbf{\Lambda}$. While the first is simply a LogNormal distribution with a fixed mean and variance, the prior distribution over the precision matrix $\mathbf{\Lambda}$ requires a deeper understanding. First of all, the Cholesky parameterization $\mathbf{\Lambda} = \mathbf{L}\mathbf{L}^\top$ in the context of MCMC sampling introduces a change of variable. We impose a prior probability over a non-linear transformation of \mathbf{L} , while this is the variable that is actually sampled together with \mathbf{U} , \mathbf{Z} and σ_f^2 .

The change of measure induced by the change of variables, requires the determinant of the Jacobian \mathcal{J} :

$$p(\text{vec } \mathbf{L}) = p(\text{vec } \mathbf{\Lambda}) |\mathcal{J}(\text{vec } \mathbf{\Lambda}, \text{vec } \mathbf{L})| \quad (14)$$

where, for simplicity, we vectorized the matrices $\mathbf{\Lambda}$ and \mathbf{L} so that the Jacobian of the transformation is simpler to define. By vectorizing the matrices $\mathbf{\Lambda}$ and \mathbf{L} , its determinant takes a particularly convenient form, which can be computed linearly in D (Magnus & Neudecker, 1980):

$$\log |\mathcal{J}(\text{vec } \mathbf{\Lambda}, \text{vec } \mathbf{L})| = \log 2^D \prod_d (\mathbf{L}_{dd})^{D-d+1}. \quad (15)$$

We have identified two different families of priors $p(\mathbf{\Lambda})$: (1) matrix-variate distributions over $\mathbf{\Lambda}$ and (2) factorized scalar distributions over the single entries of the precision $\mathbf{\Lambda}$ defined as $p(\mathbf{\Lambda}) = \prod_{ij} p(\mathbf{\Lambda}_{ij})$.

5.1 Matrix-variate priors

Wishart prior When dealing with matrix-valued distributions over PSD matrices, a natural probability distribution to consider is the Wishart distribution. Beside being defined over symmetric PSD matrices,

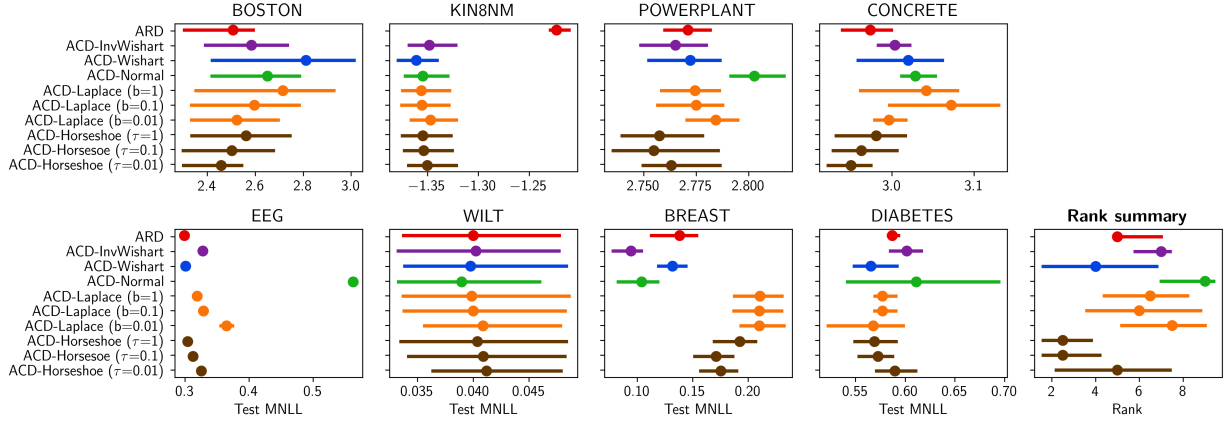


Figure 2: The ACD covariances significantly outperform ARD ones on select datasets, while being competitive throughout. Test mean negative loglikelihood (MNLL) on both UCI regression benchmarks (*top*) and classification (*bottom*) benchmarks with 20% – 80% error quantiles (lower is better), and rank summaries (*bottom right*).

the Wishart prior is a conjugate distribution of precision matrices. Considering $\mathbf{\Lambda} \in \mathbb{R}^{D \times D}$ the probability density function can be expressed as:

$$\begin{aligned} p(\mathbf{\Lambda}) &= \mathcal{W}(\mathbf{\Lambda}|\mathbf{V}, K) \\ &= C|\mathbf{\Lambda}|^{-\frac{1}{2}(K-D-1)} \exp\left(-\frac{1}{2}\text{Tr}(\mathbf{V}^{-1}\mathbf{\Lambda})\right), \end{aligned} \quad (16)$$

where $C = (2^{KD}|\mathbf{V}|^{K/2}\Gamma_D(K/2))^{-1}$ is a constant term, \mathbf{V} is the scale matrix and $K \geq D$ is the degrees of freedom parameter. The Bartlett decomposition proves that imposing independent Gaussian priors on the columns of the lower-triangular matrix $\mathbf{L} = (\mathbf{l}_1, \dots, \mathbf{l}_D)$ as $p(\mathbf{l}_i) = \mathcal{N}(\mathbf{0}, \mathbf{V})$ is equivalent to a Wishart distribution over $\mathbf{L}\mathbf{L}^T$ as $\mathcal{W}(\lambda\mathbf{I}_D, K)$. We choose $K = D$ degrees of freedom and $\mathbf{V} = D^{-1}\mathbf{I}_D$, such that the expected precision $\mathbb{E}[\mathbf{\Lambda}] = \mathbf{I}_D$ is identity.

Inverse Wishart Another prior over PSD matrices related to the Wishart distribution is the inverse Wishart. An interesting interpretation stems from the interpretation of $\mathbf{\Lambda}$ as a covariance matrix in the Fourier domain when Bochner’s theorem is applied:

$$\begin{aligned} k_{\text{RBF-ACD}}(\mathbf{x}_i, \mathbf{x}_j; \sigma_f^2, \mathbf{\Lambda}) &= \mathbb{E}_{\boldsymbol{\mu}, b} \sqrt{2}\sigma_f \cos(\boldsymbol{\mu}^T \mathbf{x}_i + b) \cdot \sqrt{2}\sigma_f \cos(\boldsymbol{\mu}^T \mathbf{x}_j + b), \\ \boldsymbol{\mu} &\sim \mathcal{N}(\mathbf{0}, \mathbf{\Lambda}), b \sim \text{Unif}[0, 2\pi]. \end{aligned} \quad (17)$$

Therefore, apart from viewing $\mathbf{\Lambda}$ as the precision matrix of the kernel, it can also be seen as a covariance matrix in the frequency domain, which offers a motivation for using such a prior

$$\begin{aligned} p(\mathbf{\Lambda}) &= \mathcal{IW}(\mathbf{\Lambda}|\mathbf{V}, K) \\ &= C|\mathbf{\Lambda}|^{-\frac{1}{2}(K+D+1)} \exp\left(-\frac{1}{2}\text{Tr}(\mathbf{V}\mathbf{\Lambda}^{-1})\right), \end{aligned} \quad (19)$$

where $C = (2^{KD/2}|\mathbf{V}|^{-(K/2)}\Gamma_D(K/2))^{-1}$. We set $K = D$ and $\mathbf{V} = \mathbf{I}_D$. The inverse Wishart view can translate into more efficient random Fourier approximations.

5.2 Sparsity-inducing priors

Moving away from matrix-variate distributions, it is possible to encourage sparsity in $\mathbf{\Lambda}$ with an element-wise prior. Since we might be interested in promoting sparsity in recovering covariance couplings to a different

degree than in the contribution of individual covariates, we separate the prior over the elements of $\mathbf{\Lambda}$ as follows:

$$\begin{aligned} p(\mathbf{\Lambda}) &= p(\mathbf{\Lambda}^\perp) \cdot p(\text{diag } \mathbf{\Lambda}) \\ &= \prod_{i,j|i \neq j} p(\mathbf{\Lambda}_{ij}) \prod_i p(\mathbf{\Lambda}_{ii}), \end{aligned} \quad (20)$$

where $\mathbf{\Lambda}^\perp$ and $\text{diag } \mathbf{\Lambda}$ are the off-diagonal elements and the \mathbb{R}^D array of the diagonal elements of $\mathbf{\Lambda}$, respectively. In this work, we assume a weakly informative Gaussian prior on the diagonal of $\mathbf{\Lambda}$, while we study different sparsity-promoting prior distributions for $\mathbf{\Lambda}^\perp$, as discussed next.

Laplace. A natural way to promote sparse solutions is L1-regularization (cf. graphical lasso in Friedman et al. (2008)), which is equivalent to a Laplace prior. The expression in (20) becomes:

$$p(\mathbf{\Lambda}) = \prod_{i,j|i \neq j} \mathcal{L}(\mathbf{\Lambda}_{ij}|m, b) \prod_i \mathcal{N}(\mathbf{\Lambda}_{ii}|\mu, \sigma^2), \quad (21)$$

where

$$\mathcal{L}(\mathbf{\Lambda}_{ij}|m, b) = C_1 \exp\left(-\frac{1}{b} \|\mathbf{\Lambda}_{ij} - m\|_1\right) \quad (22)$$

$$\mathcal{N}(\mathbf{\Lambda}_{ii}|\mu, \sigma^2) = C_2 \exp\left(-\frac{1}{2\sigma^2} (\mathbf{\Lambda}_{ii} - \mu)^2\right), \quad (23)$$

where C_1 and C_2 are the normalizing constants. We fix $m = \mu = 0$, $\sigma^2 = 1$, and analyze the resulting posteriors for different sparsity coefficients b (lower b increases sparsity).

Horseshoe. The Horseshoe prior has become a popular probabilistic sparsity-inducing prior (Carvalho et al., 2009),

$$\mathbf{\Lambda}_{ij}|\sigma, \tau \sim \mathcal{N}(0, \sigma^2 \tau^2), \quad \sigma \sim C^+(0, 1) \quad (24)$$

where $C^+(0, 1)$ is a Half-Cauchy distribution for the local shrinkage σ , while τ is the global shrinkage parameter. The Horseshoe density of a single entry $\mathbf{\Lambda}_{ij}$ is

$$\pi_\tau(\mathbf{\Lambda}_{ij}) = \frac{1}{\sqrt{2\pi^3 \tau^2}} \exp\left(\frac{\mathbf{\Lambda}_{ij}^2}{2\tau^2}\right) E_1\left(\frac{\mathbf{\Lambda}_{ij}^2}{2\tau^2}\right), \quad (25)$$

where $E_1(\cdot)$ is the exponential integral function that can be approximated by elementary functions. These priors for covariance parameters have been considered, e.g., in Oh et al. (2019) and Eriksson & Jankowiak (2021).

6 Experiments

We consider eight UCI datasets as a benchmark to evaluate the performance of GP models for regression and classification tasks. We standardize all datasets to zero mean and unit variance, and report all results with five-fold cross-validation. Following previous works (e.g., Rossi et al. (2021)), we report test MNLL for all data, and normalized root mean square error (RMSE) for regression and error for classification tasks.

In all experiments, we chose to approximate GPs with 500 inducing points. We ran BSGP for 10,000 iterations with a step-size of 0.01 and mini-batch of 1,000 data points. We evaluate performance on test data from 50 samples collected during training after 1,500 burn-in iterations and using a thinning of 180. We adopt gradient clipping for numerical stability and to avoid exploding gradients, which we experienced when working with the Horseshoe priors. Note that augmentation techniques might be useful in improving the numerical stability of the sampler for this prior and, more generally, to avoid pathologies introduced by hierarchical models (e.g., (Eriksson & Jankowiak, 2021; Betancourt & Girolami, 2013)).

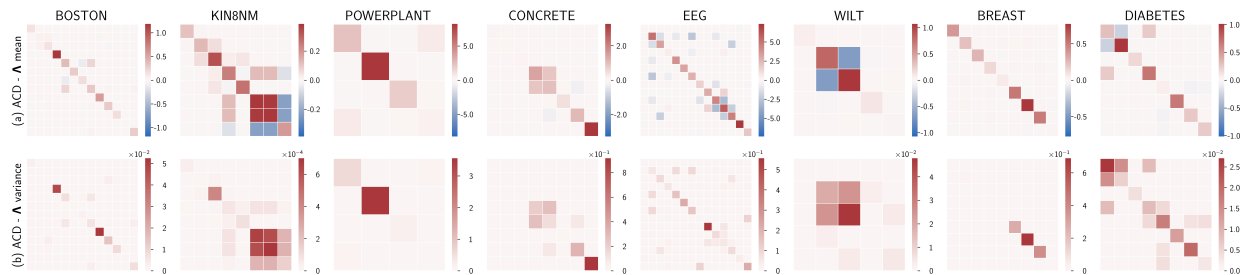


Figure 3: The precision matrices reveal couplings, redundancies and separabilities. The posterior mean (a) and variance (b) of precision matrices Λ of UCI benchmark datasets with Horseshoe prior ($\tau = 0.1$).

6.1 UCI benchmarks

With the above setup, we report results on UCI benchmarks by considering various choices of priors (See Table 1). For the kernel variance σ_f^2 we placed a Lognormal prior with unit variance and mean 0.05 as in Rossi et al. (2021). The proposed MCMC scheme yields good convergence and sampling efficiency, as illustrated in Appendix C in the supplement; see also Fig. 20 and Fig. 21 for insights on the multimodality of the posterior. Fig. 2 shows the comparative performance for the UCI benchmark datasets, including the range between the 20th and 80th percentiles over the different folds, together with a rank summary. For the small data sets, we could also run full GPs and we observed a similar trend; we refer the reader to Fig. 15 and Fig. 16 for a direct comparison between full GPs and BSGPs.

Interestingly, different choices of prior and prior hyper-parameters yield comparable performance. A closer inspection indicates that the element-wise Laplace prior performs worst overall, and this might be due to the heavy sparsity promoted by this prior (or the lack thereof) for some hyper-parameter settings (Fig. 5). The element-wise Horseshoe prior, while promoting sparsity, fares slightly better than the Laplace prior. It is interesting how the inverse Wishart prior, which operates directly on Λ , promotes some sparsity after all, while offering relatively competitive performance.

6.2 Sparse couplings

Next, we study the precision matrices themselves. Fig. 3 shows the posterior precisions of all benchmark datasets. Notably we see strong dependencies emerging in `kin8nm`, `eeg` and `wilt` datasets, while `powerplant`, `concrete`, `breast`, and `diabetes` are sparsely diagonal. For the `concrete` dataset, the posterior precision reveals both a coupling between *water* and *superplasticizer* (a chemical component that improves the usability of concrete without increasing water content) and the importance of the covariate *age*. In the `wilt` dataset, beyond the importance of individual covariates, our approach reveals a coupling between mean *green* and *red* values in multi-spectral remote sensing images used to classify diseased tree and other land covers. Finally, in the `diabetes` dataset, we observe a stronger coupling between *pregnancies*, *glucose*, and *skin thickness*, while unveiling the importance of individual covariates such as *glucose*, *skin thickness*, and *body mass index*. We notice that the standard deviation of the elements on Λ is larger for covariate pairs with large positive/negative partial covariance, while it is generally small for pairs that have small partial covariance. This indicates both the relative scaling of uncertainty, and the flexibility in coupling magnitudes.

We provide a more in-depth look into the dependencies in Fig. 1 (page 2) that contrasts the precision matrix of the ARD covariance of `kin8nm` and `breast` datasets to the posterior mean and standard deviations of the precision matrix of the ACD covariances. The ACD detects that 5th and 6th covariates of `kin8nm` are close to redundant, and negatively coupled to 7th covariate. Less evidently, in `breast` we detect coupling chains over covariates such as (0,3,6,8) and (1,6,7), indicating predictive dependencies in the data. Fig. 9 shows for this dataset how a different choice of the prior distribution over Λ can reveal a different and sparser structure of the couplings. We visualize these as circular graphs along with the standard deviations of the elements of the precision matrices.

Fig. 4 shows an ablation of comparing the posterior mean precision structures from Horseshoe prior with $\tau = \{0.01, 0.1, 1\}$ on the `concrete` dataset. The Horseshoe is able to sparsify the entire structure into an ARD-like structure, while higher $\tau = 1$ reveals off-diagonal dependencies. To obtain more intuition into the couplings, we also visualize the covariate graphs in the bottom panel of Fig. 4 that indicate for instance the strong dependence between the 3rd and the 4th covariate. Further illustrations on all UCI data sets for the Wishart prior Fig. 9, inverse Wishart Fig. 10, and Laplace prior with $b = 0.1$ Fig. 11 can be found in the supplement.

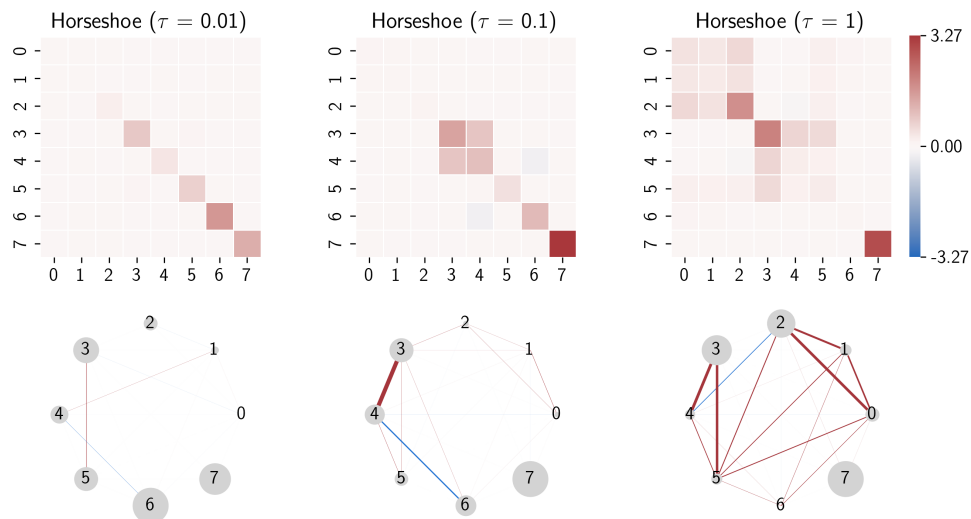


Figure 4: The sparsity control of Horseshoe prior. The posterior mean precision matrices of Horseshoe priors on `concrete` dataset with high (left) to low (right) sparsity.

6.3 Sparsification effect

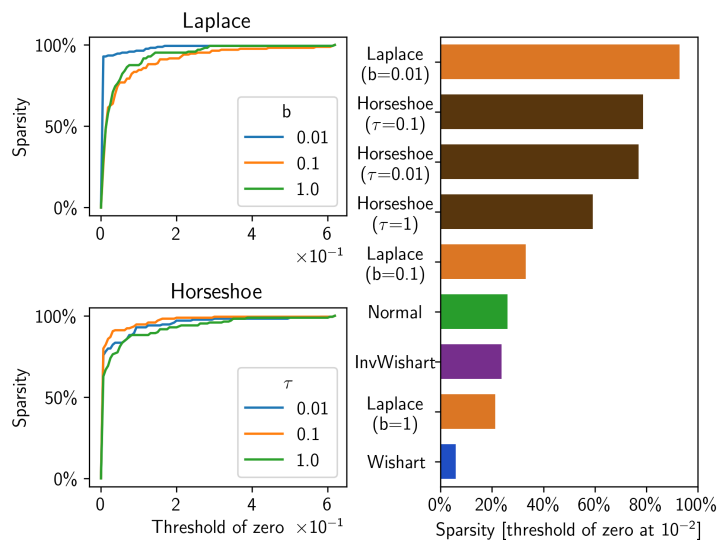


Figure 5: The sparsification of boston dataset. Left: Relationship between precision sparsity and hyperparameters. Right: Posterior mean sparsity from different priors.

Fig. 5 shows the sparsity of the posterior precision matrices $\mathbf{\Lambda}$ in the `boston` dataset. Surprisingly, the Inverse Wishart prior has an intrinsic sparsifying effect. The Laplace prior sparsifies according to its hyperparameter b while, for this dataset, the Horseshoe prior with $\tau = 0.1$ achieves slightly more sparsity than the Horseshoe prior with $\tau = 0.01$.

As a conclusion of these experiments on UCI, we observe that the Horseshoe prior obtains better performance compared to the Laplace prior and it is competitive with matrix-variate priors. Also, these generally outperform the ARD covariance. Interestingly, there seems to be some data-dependent effect connecting sparsity and performance; in data sets such as `boston`, high sparsity seems to be associated with good performance, while for others such as `eeg` it is the opposite. This indicates that sparsity should perhaps be treated as a hyper-parameter and learned together with the model. We leave this interesting development for future work.

6.4 Low-rank precision matrices

We also look at the effect of low-rank precision matrices. Fig. 6 shows the posterior precision patterns learned by the Wishart prior using a PCA with rank 11, 7 or 3 in contrast to the full rank 13. The performance degrades strongly at ranks lower than 11, which is likely indicative of the intrinsic rank of the dataset for this task. Also, it is interesting to observe how removing components, which reduces the variance explained by the PCA parameterization, affects the ability of the model to recover covariate couplings, indicating that such information is somewhat contained in the removed components. In the same figure, we also report the results obtained by an ARD modeling for the matrix $\mathbf{\Lambda}_d$ in the PCA-based parameterization. As shown in the figure, some structure emerges even in this case, but the parameterization is less flexible, and retaining a large number of components leads to nearly diagonal precisions, preventing the possibility to reveal interesting covariate couplings.

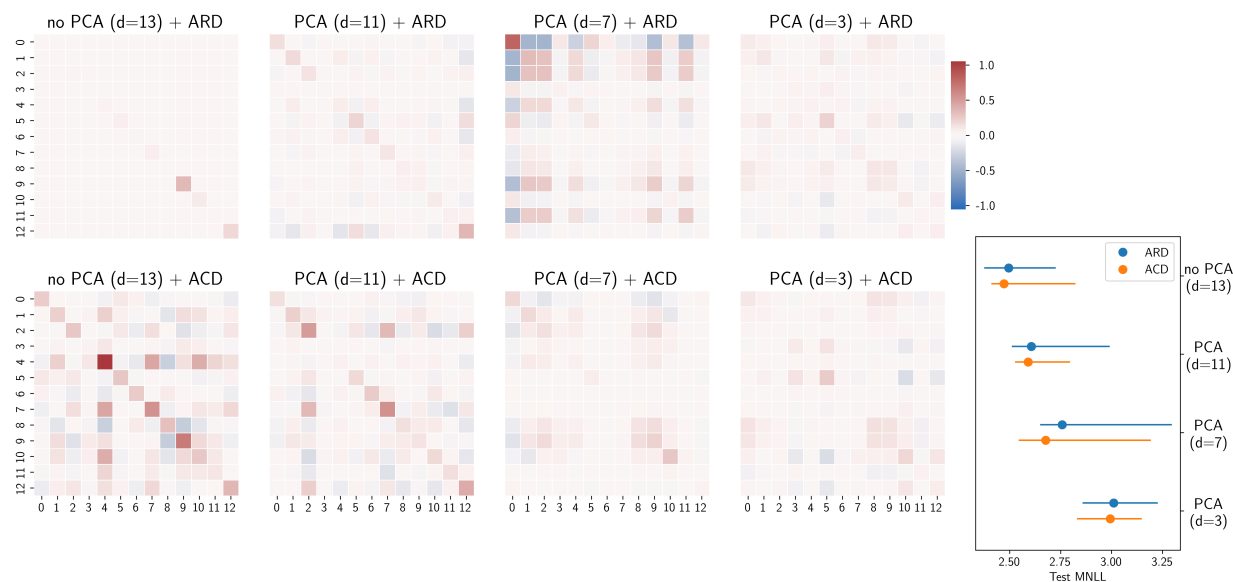


Figure 6: Overly low rank degrades performance. Posterior precisions with Wishart prior of varying rank d of Eq. (12) (*top*) and corresponding performances (*bottom*) on the `boston` dataset. The matrices displayed are of the form $\mathbf{P}_d \mathbf{\Lambda}_d \mathbf{P}_d^\top$ where $\mathbf{\Lambda}_d$ is diagonal and full in the ARD and ACD parameterizations, respectively. \mathbf{P}_d is obtained by pre-processing the training data with PCA (d components). This setup uses 200 inducing points and 3-fold cross validation.

Table 2: MoCap results on subject 09 using GP-ODE with ARD and ACD kernels. The data are projected in a 5-dimensional latent space and the model is trained with dimension-wise kernels: ARD (5 lengthscales per dimension) and ACD (5 × 5 full precision matrix per dimension). We report Test MNLL and Test Mean Squared Error (MSE) over 5 different folds.

| Metric | Method | Subject 09 (short) |
|---------|--------------------|--------------------|
| MNLL(↓) | GP-ODE-vanilla ARD | 1.17 ± 0.02 |
| | GP-ODE-vanilla ACD | 1.27 ± 0.17 |
| MSE(↓) | GP-ODE-vanilla ARD | 10.64 ± 1.58 |
| | GP-ODE-vanilla ACD | 14.72 ± 6.95 |

6.5 Dependencies of motion capture data

We illustrate the capability of ACD covariances to reveal dependencies in a motion capture task, where the subjects internal connectivity is known (Fig. 7). We observe a trajectory $\mathbf{Y} = (\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_N)^T \in \mathbb{R}^{N \times D}$ over N timepoints, where $\mathbf{y}_i \in \mathbb{R}^D$ represents the noisy observation of subject state $\mathbf{x}(t_i) \in \mathbb{R}^D$ at time t_i . The state consists of a total of $D = 50$ measurements across 21 body parts (Fig. 7). We follow the GP-ODE model (Heinonen et al., 2018; Hegde et al., 2022), where the state follows an ordinary differential equation $\dot{\mathbf{x}}(t) = \mathbf{f}(\mathbf{x}(t))$ with a vector-valued GP prior on the differential $\mathbf{f} : \mathbb{R}^D \mapsto \mathbb{R}^D$,

$$\mathbf{f}(\mathbf{x}) \sim \mathcal{GP}(\mathbf{0}, K_{\theta}(\mathbf{x}, \mathbf{x}')), \quad (26)$$

where $K_{\theta} \in \mathbb{R}^{D \times D}$ is an operator-valued kernel. The most straightforward covariance function is a separable one $K(\mathbf{x}, \mathbf{x}'; \theta) = k(\mathbf{x}, \mathbf{x}'; \theta) \mathbf{I}_D$, where we learn a shared precision matrix for all outputs. As an alternative, we also consider a variant $K(\mathbf{x}, \mathbf{x}'; \theta) = \text{diag}\{k(\mathbf{x}, \mathbf{x}'; \theta_1), \dots, k(\mathbf{z}, \mathbf{z}'; \theta_D)\}$, where each diagonal entry has its own kernel $k(\mathbf{z}, \mathbf{z}'; \theta_i)$ and its own precision matrix Λ_i associated with output x_i .

Fig. 7 shows the posterior shared precision mean pooled over the body parts in a human walk cycle. A rich pattern of dependencies emerges. For instance, right and left wrists are strongly coupled across the body, while being negatively coupled to each other. The wrists move in large, cyclic and synchronised patterns, while the back and root have little relevance, indicating their smaller movement ranges during walking. Finally, many adjacent body parts are coupled, such as foot and tibia, and wrist and radius. Table 2 shows the performance between ARD and ACD on subject 09, where the likelihoods are similar, but ACD does perform worse in mean square error. The purpose of the experiment was to demonstrate structure learning with standard inference runs, and we did not focus on performance tuning, which ODE models are known to be finicky about (Hegde et al., 2022). As a final note, we emphasize the importance of these results in light of possible simpler alternative analyses to determine covariate couplings; for instance, we could study the inverse of the empirical covariance of the input covariates directly. For completeness, we report this analysis in the Appendix in Fig. 17, where the sample covariance is also obtained by retaining the first 15 principal components. Unlike this simplistic approach, which one could view as a preliminary analysis of the covariates, the results in Fig. 7 illustrate the emergence of covariate couplings informed by the supervised learning task.

7 Conclusions

In the literature of GPs, covariances equipped with ARD are popular. These materialize with the definition of a set of length-scale parameters scaling the inputs, which are then optimized or inferred based on the marginal likelihood (or an approximation/bound). In this work, we revisited a more general definition of anisotropic covariances, where the distance metric among inputs is determined by a PSD precision matrix. We showed that this extension provides a framework for metric learning and we discussed some interesting insights on the determination of couplings among covariates. Crucially, thanks to a fully Bayesian scalable formulation of GPs, we can operate with virtually any number of data points and obtain samples from the posterior distribution over such covariate couplings, which can be used to determine the level of confidence in their predictive power.

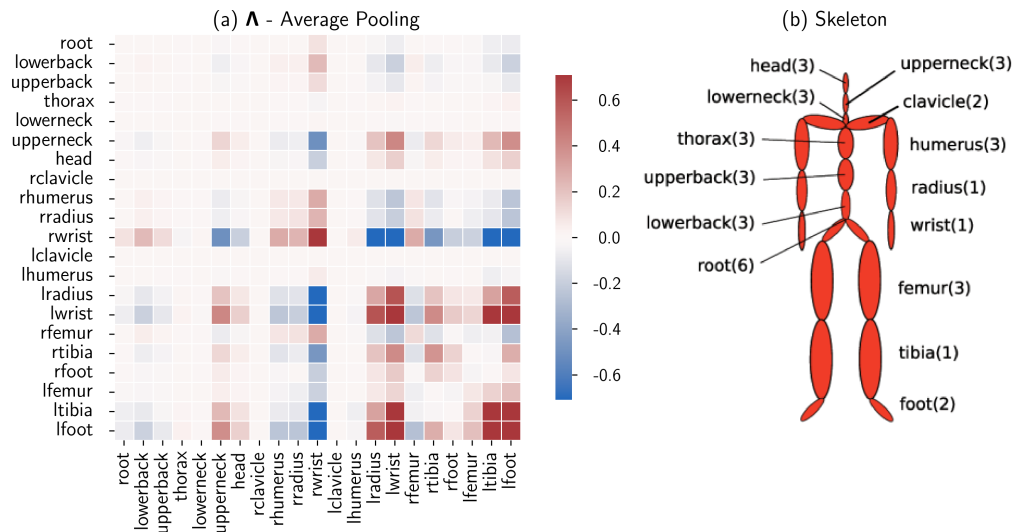


Figure 7: The ACD covariance reveals a highly regular coupling structure from human motion. GP-ODE model trained with shared ACD covariance in a latent space of 15 dimensions. Panel (a) shows the precision matrix Λ reporting just the average value for each group of sensors. Panel (b) shows the reference skeleton connectivity.

We also studied priors for the precision matrix Λ determining the input metric. We showed that element-wise Laplace and Horseshoe priors provide the highest level of sparsity, while Horseshoe priors seem to offer better performance. Interestingly, the inverse Wishart prior offers higher sparsity than the Wishart prior with overall comparable performance.

In order to address the quadratic scalability with respect to the number of covariates, we also revisited the work by Vivarelli & Williams (1998), which proposes a low-dimensional projection of the inputs through PCA, in light of modern scalable GPs and inference.

We are currently investigating an extension of our approach whereby the conclusions we can draw from the analysis of Λ are in terms of conditional independence statements. In order to do this, we plan to extend our model to target the modeling of both labels and inputs, including a prior over the inputs $p(\{\mathbf{x}_n\}|\Lambda)$ in the form of a Markov Random Field, where Λ now determines the conditional independence among covariates.

As future work, it would also be interesting to consider ways in which we could learn the adequate level of sparsity from data by inferring relevant prior hyper-parameters.

References

- S. Banerjee and S. Ghosal. Posterior convergence rates for estimating large precision matrices using graphical models. *Electronic Journal of Statistics*, 8(2):2111 – 2137, 2014.
- S. Banerjee and S. Ghosal. Bayesian structure learning in graphical models. *Journal of Multivariate Analysis*, 136(C):147–162, 2015.
- S. Banerjee, I. Castillo, and S. Ghosal. Bayesian inference in high-dimensional models. *arxiv:2101.04491*, 2021.
- D. Barry, J.-Y. Parlange, and L. Li. Approximation for the exponential integral (Theis well function). *Journal of Hydrology*, 227(1):287–291, 2000.
- M. J. Betancourt and M. Girolami. Hamiltonian Monte Carlo for hierarchical models. *arXiv:1312.0906*, 2013.
- J. Q. Candela and C. E. Rasmussen. A Unifying View of Sparse Approximate Gaussian Process Regression. *Journal of Machine Learning Research*, 6:1939–1959, 2005.

- C. Carvalho, N. Polson, and J. Scott. Handling sparsity via the Horseshoe. In *AISTATS*, 2009.
- I. Castillo, J. Schmidt-Hieber, and A. van der Vaart. Bayesian linear regression with sparse priors. *The Annals of Statistics*, 43(5):1986 – 2018, 2015.
- S. Chen and D. Zhang. Robust image segmentation using FCM with spatial constraints based on new kernel-induced distance measure. *IEEE Transactions on Systems, Man and Cybernetics, Part B*, 34(4): 1907–1916, 2004.
- K. Cutajar, E. V. Bonilla, P. Michiardi, and M. Filippone. Random feature expansions for deep Gaussian Processes. In *ICML*, 2017.
- A. C. Damianou and N. D. Lawrence. Deep Gaussian Processes. In *AISTATS*, 2013.
- D. Duvenaud, H. Nickisch, and C. Rasmussen. Additive Gaussian processes. In *NeurIPS*, 2011.
- D. Eriksson and M. Jankowiak. High-dimensional Bayesian optimization with sparse axis-aligned subspaces. In *Proceedings of the Thirty-Seventh Conference on Uncertainty in Artificial Intelligence*, 2021.
- M. Filippone and M. Girolami. Pseudo-marginal Bayesian inference for Gaussian processes. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 36(11):2214–2226, 2014.
- J. Friedman, T. Hastie, and R. Tibshirani. Sparse inverse covariance estimation with the graphical lasso. *Biostatistics*, 2008.
- L. Gan, N. N. Narisetty, and F. Liang. Bayesian regularization for graphical models with unequal shrinkage. *Journal of the American Statistical Association*, 114:1218 – 1231, 2018.
- P. Hegde, Ç. Yıldız, H. Lähdesmäki, S. Kaski, and M. Heinonen. Variational multiple shooting for Bayesian ODEs with Gaussian processes. In *UAI*, 2022.
- M. Heinonen, C. Yildiz, H. Mannerström, J. Intosalmi, and H. Lähdesmäki. Learning unknown ODE models with Gaussian processes. In *ICML*, 2018.
- J. Hensman, N. Fusi, and N. D. Lawrence. Gaussian Processes for Big Data. In *UAI*, 2013.
- A. Kalaitzis. Image inpainting with Gaussian Processes. Master’s thesis, University of Edinburgh, 2009.
- K. Khare and B. Rajaratnam. Wishart distributions for decomposable covariance graph models. *The Annals of Statistics*, 39(1):514 – 555, 2011.
- K. Krauth, E. V. Bonilla, K. Cutajar, and M. Filippone. AutoGP: Exploring the capabilities and limitations of Gaussian process models. In *UAI*, 2017.
- M. Lázaro-Gredilla, J. Quinero-Candela, C. E. Rasmussen, and A. R. Figueiras-Vidal. Sparse Spectrum Gaussian Process Regression. *Journal of Machine Learning Research*, 11:1865–1881, 2010.
- K. Lee and J. Lee. Estimating large precision matrices via modified Cholesky decomposition. *Statistica Sinica*, (31):173–196, 2021.
- Y. Li, B. A. Craig, and A. Bhadra. The graphical horseshoe estimator for inverse covariance matrices. *Journal of Computational and Graphical Statistics*, 28:747 – 757, 2017.
- D. J. C. MacKay. Probable networks and plausible predictions – a review of practical bayesian methods for supervised neural networks. *Network: Computation in Neural Systems*, 6(3):469, 1995.
- J. R. Magnus and H. Neudecker. The elimination matrix: Some lemmas and applications. *SIAM Journal on Matrix Analysis and Applications*, 1(4):422–449, 1980.
- B. Matérn. *Spatial Variation*. Springer, 1960.

- A. G. Matthews, M. van der Wilk, T. Nickson, K. Fujii, A. Boukouvalas, P. León-Villagrà, Z. Ghahramani, and J. Hensman. GPflow: A Gaussian Process library using TensorFlow. *Journal of Machine Learning Research*, 18(40):1–6, April 2017.
- K. Murphy. *Probabilistic Machine learning: Advanced topic*. The MIT Press, 2023.
- I. Murray and R. P. Adams. Slice sampling covariance hyperparameters of latent Gaussian models. In J. D. Lafferty, C. K. I. Williams, J. Shawe-Taylor, R. S. Zemel, and A. Culotta (eds.), *Advances in Neural Information Processing Systems 23: 24th Annual Conference on Neural Information Processing Systems 2010. Proceedings of a meeting held 6-9 December 2010, Vancouver, British Columbia, Canada.*, pp. 1732–1740. Curran Associates, Inc., 2010.
- R. M. Neal. *Bayesian Learning for Neural Networks (Lecture Notes in Statistics)*. Springer, 1 edition, August 1996.
- M. Noack, G. Doerk, R. Li, J. Streit, R. Vaia, K. Yager, and M. Fukuto. Autonomous materials discovery driven by gaussian process regression with inhomogeneous measurement noise and anisotropic kernels. *Scientific reports*, 2020.
- C. Oh, J. M. Tomczak, E. Gavves, and M. Welling. *Combinatorial Bayesian optimization using the graph cartesian product*. Curran Associates Inc., Red Hook, NY, USA, 2019.
- M. Opper and O. Winther. Gaussian processes for classification: Mean-field algorithms. *Neural Computation*, 12(11):2655–2684, 2000.
- C. Paciorek and M. Schervish. Nonstationary covariance functions for Gaussian Process Regression. In *NIPS*, 2003.
- T. Poggio and F. Girosi. Networks for approximation and learning. In *Proceedings of the IEEE*, 1990.
- A. Rahimi and B. Recht. Random Features for Large-Scale Kernel Machines. In *NIPS*, 2008.
- C. E. Rasmussen and C. Williams. *Gaussian Processes for Machine Learning*. MIT Press, 2006.
- S. Rossi, M. Heinonen, E. Bonilla, Z. Shen, and M. Filippone. Sparse Gaussian Processes revisited: Bayesian approaches to inducing-variable approximations. In *AISTATS*, 2021.
- A. Roverato. Cholesky decomposition of a hyper inverse Wishart matrix. *Biometrika*, 87(1):99–112, 2000.
- A. Roverato. Hyper inverse Wishart distribution for non-decomposable graphs and its application to bayesian inference for gaussian graphical models. *Scandinavian Journal of Statistics*, 29(3):391–411, 2002.
- K. Sagar, S. Banerjee, J. Datta, and A. Bhadra. Precision matrix estimation under the horseshoe-like prior–penalty dual. *Electronic Journal of Statistics*, 18(1):1 – 46, 2024.
- H. Salimbeni and M. Deisenroth. Doubly Stochastic Variational Inference for Deep Gaussian Processes. In *NeurIPS*, 2017.
- J. Shawe-Taylor and N. Cristianini. *Kernel Methods for Pattern Analysis*. Cambridge University Press, New York, NY, USA, 2004.
- R. Silva and Z. Ghahramani. The hidden life of latent variables: Bayesian learning with mixed graph models. *Journal of Machine Learning Research*, 10(41):1187–1238, 2009.
- E. Snelson and Z. Ghahramani. Sparse Gaussian Processes using Pseudo-inputs. In *NIPS*, 2005.
- M. Titsias and M. Lazaro-Gredilla. Variational Inference for Mahalanobis Distance Metrics in Gaussian Process Regression. In *NIPS*, 2013.
- M. K. Titsias. Variational Learning of Inducing Variables in Sparse Gaussian Processes. In *AISTATS*, 2009.

- W. van den Boom, A. Beskos, and M. D. Iorio. The g-Wishart weighted proposal algorithm: Efficient posterior computation for Gaussian graphical models. *Journal of Computational and Graphical Statistics*, 31(4):1215–1224, 2022.
- F. Vivarelli and C. Williams. Discovering hidden features with Gaussian processes regression. In *NIPS*, 1998.
- H. Wang. Bayesian Graphical Lasso Models and Efficient Posterior Computation. *Bayesian Analysis*, 7(4): 867 – 886, 2012.
- H. Wang. Scaling It Up: Stochastic Search Structure Learning in Graphical Models. *Bayesian Analysis*, 10 (2):351 – 377, 2015.
- C. K. I. Williams and D. Barber. Bayesian classification with Gaussian Processes. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 20:1342–1351, 1998.
- C. K. I. Williams and M. Seeger. Using the Nyström method to speed up kernel machines. In *NIPS*, 2000.
- A. Wilson and H. Nickisch. Kernel Interpolation for Scalable Structured Gaussian Processes (KISS-GP). In *ICML*, 2015.
- A. G. Wilson, Z. Hu, R. Salakhutdinov, and E. P. Xing. Deep Kernel Learning. In *AISTATS*, 2016.
- R. Xiang, K. Khare, and M. Ghosh. High dimensional posterior convergence rates for decomposable graphical models. *Electronic Journal of Statistics*, 9(2):2828 – 2854, 2015.

A Experimental details

In this section, we present details to reproduce our experimental campaign. All the experiments were conducted on Google Colab. The code to reproduce the results on the UCI and MoCap datasets can be found at https://github.com/mattyred/BayesianSGP_Automatic_Coupling_Determination and <https://github.com/mattyred/gaussian-process-odes-acd/> respectively.

BSGP model We use $M = 500$ inducing points initialized by a k-means algorithm as commonly used in practice and we place a Normal prior $p_{\xi}(\mathbf{Z})$ over the inducing locations \mathbf{Z} . For inference, we use an adaptive version of Stochastic Gradient Hamiltonian Monte Carlo (SGHMC) in which the hyperparameters are automatically tuned during a burn-in phase. We set the default hyperparameter of the number of SGHMC steps to $K = 10$. Exclusively for regression datasets with Gaussian likelihood, we employ an Adam optimizer with a learning rate set at 0.01 for optimizing the variance of the likelihood.

ARD kernel We use the Radial Basis Function (RBF) kernel with Automatic Relevance Determination (ARD) placing a LogNormal prior with unit variance and means equal to 1 and 0.05 for the lengthscales and variance, respectively.

ACD kernel We place a LogNormal prior with unit variance and mean 0.05 over the kernel variance σ_f^2 while over the precision matrix $\mathbf{\Lambda}$ we explore a wide range of priors. Note that for the Horseshoe prior we use the exponential integral approximation of Barry et al. (2000).

Table 3: Parameter settings for the UCI experiments.

| parameter | value |
|-------------------------|-------|
| num. of inducing points | 500 |
| mini-batch size | 1000 |
| num. iterations | 10500 |
| step size | 0.01 |
| momentum | 0.05 |
| num. of burn-in steps | 1500 |
| num. of samples | 50 |
| thinning interval | 180 |

B Simulated dataset

In this section we carry out an experiment using simulated datasets with known underlying precision matrices. In particular we assess the ability of the BSGP model using a ACD kernel to recover the true precision $\mathbf{\Lambda}$ while fitting simple regression problems. Here it’s described how the simulated regression datasets are constructed and some experiments conducted to show the behaviour of the model. We consider N input-label pairs $\{\mathbf{X}, \mathbf{y}\} = \{(\mathbf{x}_n, y_n)\}_{n=1}^N$ with $\mathbf{x}_n \in \mathbb{R}^D$ and $y_n \in \mathbb{R}$ defined as follows:

$$\begin{aligned}
 \mathbf{x}_n &\sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \\
 \mathbf{K}_{\mathbf{xx}} : \mathbf{K}_{\mathbf{xx}}[i, j] &= \sigma_f^2 ((\mathbf{x}_i - \mathbf{x}_j)^\top \mathbf{\Lambda} (\mathbf{x}_i - \mathbf{x}_j)) \\
 \mathbf{y} &\sim \mathcal{N}(\mathbf{0}, \mathbf{K}_{\mathbf{xx}} + \sigma_n \mathbf{I})
 \end{aligned} \tag{27}$$

Once that the underlying precision $\mathbf{\Lambda}$ has been constructed, specifying a value for the kernel variance σ_f^2 and another one for the Gaussian noise in observations via σ_n is sufficient. The regression dataset $\{\mathbf{X}, \mathbf{y}\}$ can be used to train a BSGP model by means of the ACD kernel. Through acquiring samples of the precision matrix $\mathbf{\Lambda}$, we aim to recover the original underlying precision used to generate the data. A visual insight into this experiment is given in Fig. 8. Note also that the covariance/precision of the input covariates is the identity; therefore, analyzing the covariates by themselves would not reveal any covariate couplings, while our model is capable to retrieve the couplings that were used to generate the labels.

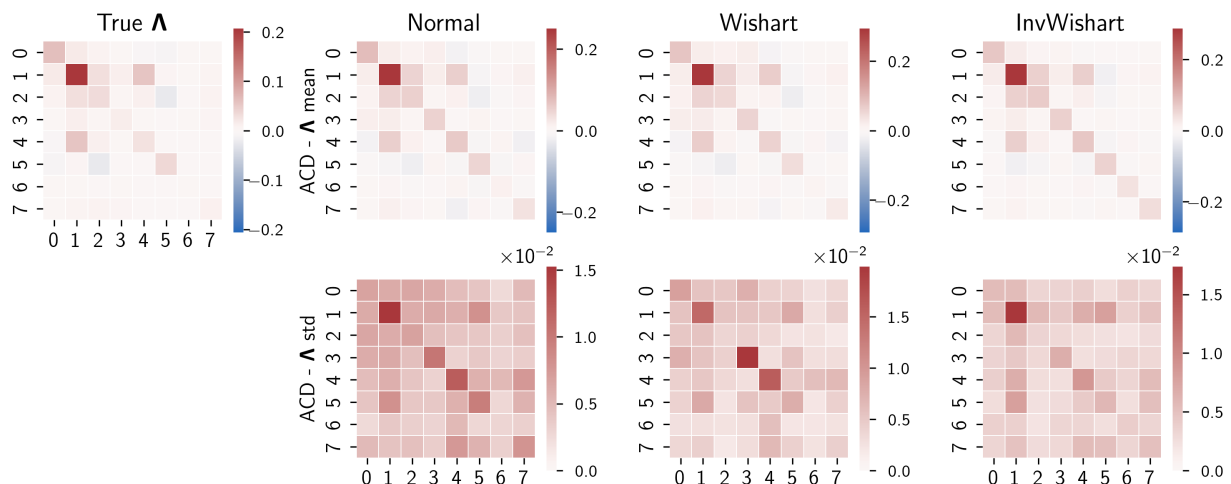


Figure 8: Underlying sparse precision (top left) compared with mean and standard deviation of the Λ samples obtained with different priors. The dataset is made of $N = 1000$ samples and the labels are obtained according to Eq. 27 setting $\sigma_f^2 = 1$ and $\sigma_n = 0.1$

C Additional results

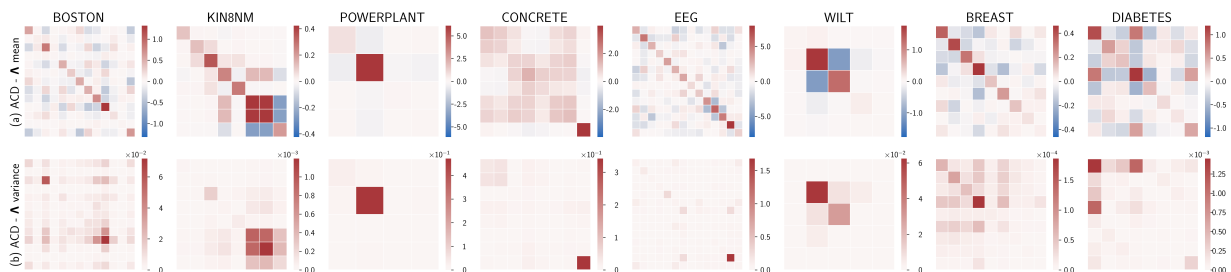


Figure 9: Posterior precision matrix mean (a) and variance (b) with Wishart prior.

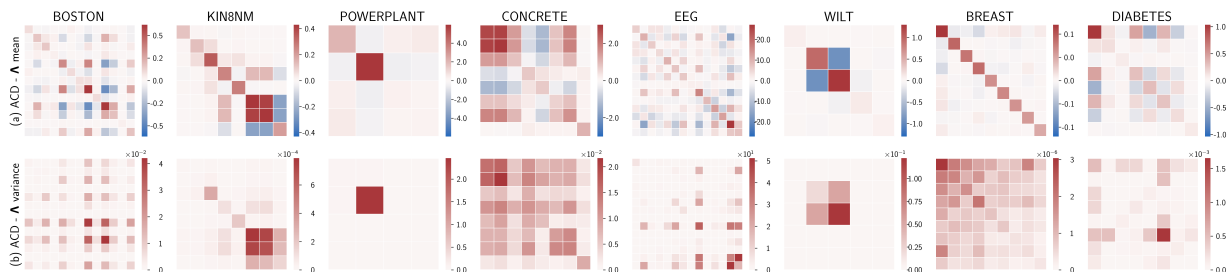


Figure 10: Posterior precision matrix mean (a) and variance (b) with Inverse Wishart prior.

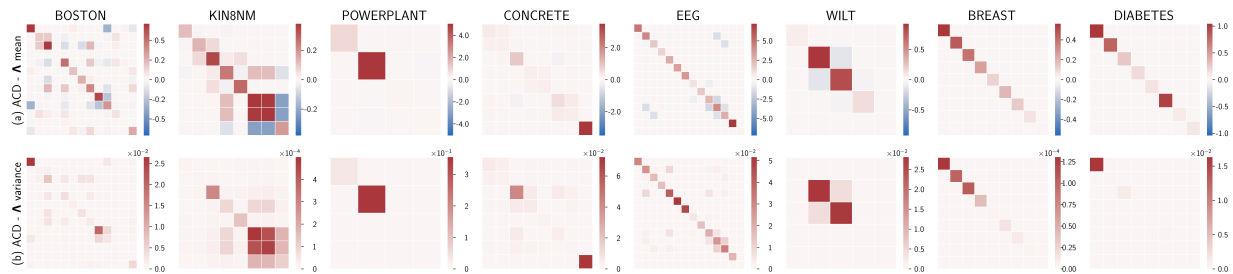


Figure 11: Posterior precision matrix mean (a) and variance (b) with Laplace prior $b = 0.1$.

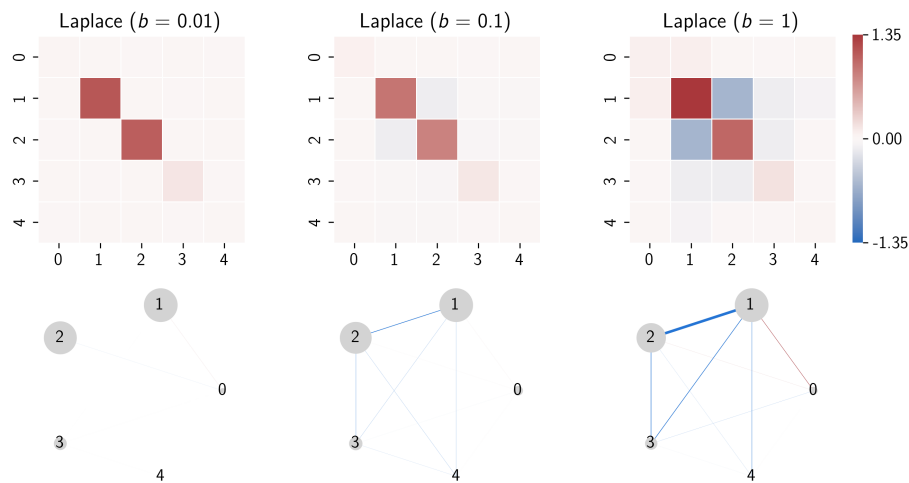


Figure 12: The precision matrices of the `wilt` dataset using Laplace prior show a progressive level sparsity.

Table 5: Runtime comparison of BSGP training and inference on the `boston` dataset using ARD kernel and ACD (with Wishart prior) kernel. Values are reported in seconds as mean \pm standard deviation across three different folds.

| | ARD | ACD |
|----------------|---------------------|--------------------|
| Training time | 2082.57 \pm 16.61 | 2513.56 \pm 3.14 |
| Inference time | 0.12 \pm 0.00 | 0.13 \pm 0.00 |

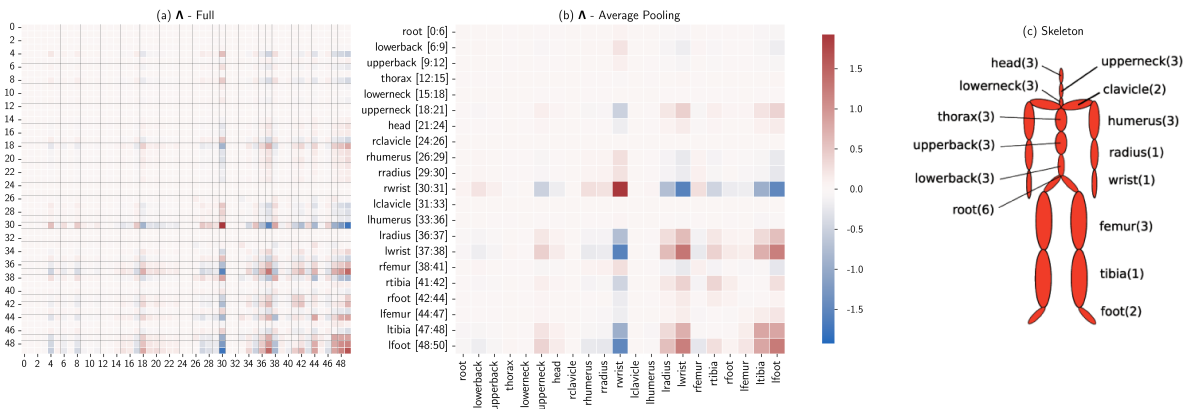


Figure 13: The full motion capture precision matrix (a), a pooled part-wise version (b) and the reference skeleton connectivity (c).

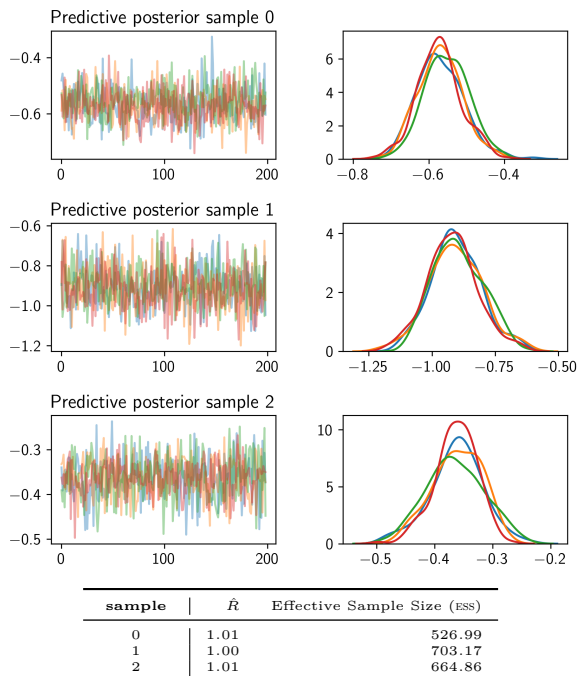


Table 4: UCI datasets used, including number of data-points and dimensionalities.

| Dataset | N | D |
|------------|--------|-----|
| boston | 506 | 13 |
| breast | 683 | 9 |
| diabetes | 783 | 8 |
| concrete | 1,030 | 8 |
| wilt | 4,839 | 5 |
| kin8nm | 8,192 | 8 |
| powerplant | 9,568 | 4 |
| eeg | 45,730 | 14 |

Figure 14: Traces of the mean of the predictive distribution for three test points on **boston** dataset with Inverse Wishart prior (4 chains, 200 samples represented); the table reports \hat{R} and Effective Sample Size (ESS) statistics for each set of 4 chains.

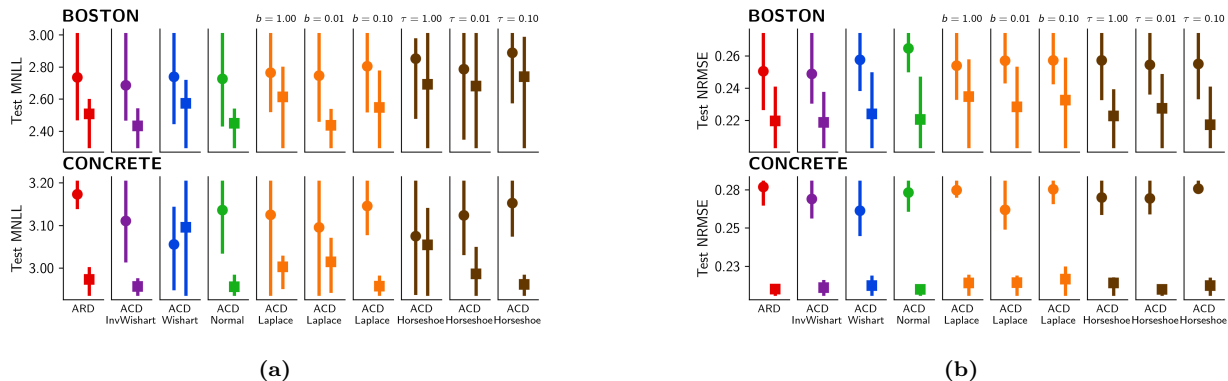


Figure 15: Comparison of full GPs (\square) vs BSGPs (\circ) with 200 inducing points on two UCI regression data sets. The metrics are MNLL in (a) and normalized RMSE in (b).

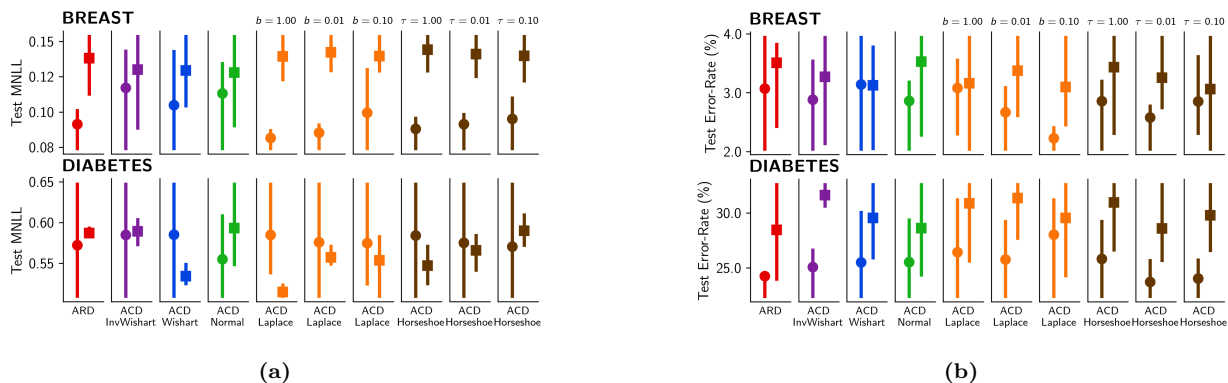


Figure 16: Comparison of full GPs (\square) vs BSGPs (\circ) with 500 inducing points on two UCI classification data sets. The metrics are MNLL in (a) and Error-Rate in (b).

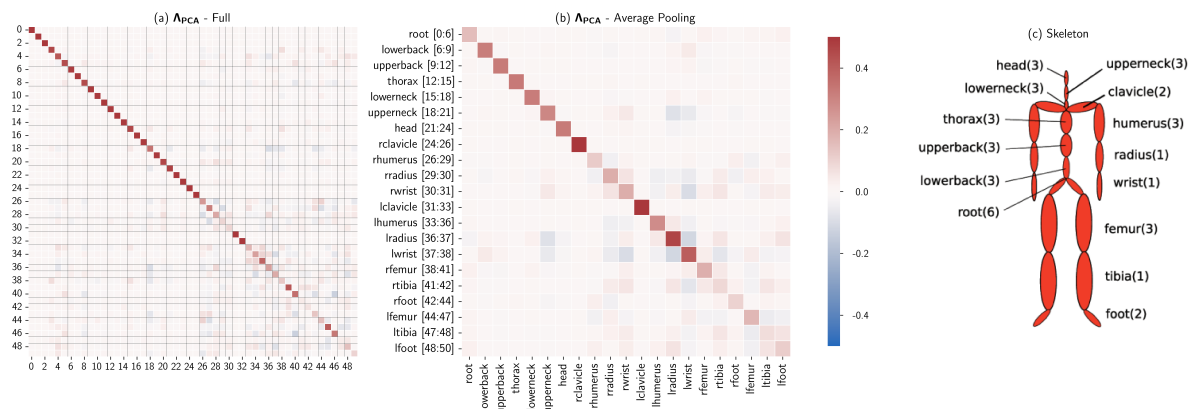


Figure 17: (a): The full estimated precision matrix $\Lambda_{PCA} = (\mathbf{P}_d^T \mathbf{S} \mathbf{P}_d)^{-1}$ obtained applying PCA with rank 15 on the training dataset; (b): The pooled version of the same matrix; (c): The reference skeleton connectivity.

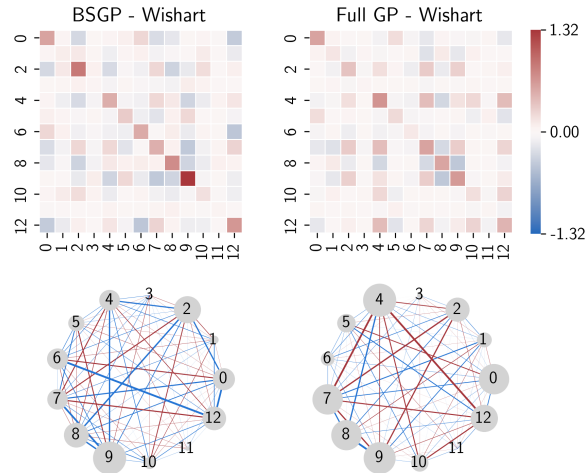


Figure 18: Comparison of posterior mean of the precision matrix \mathbf{A} on the `boston` dataset with Wishart prior for full GP vs BSGP with 500 inducing points.

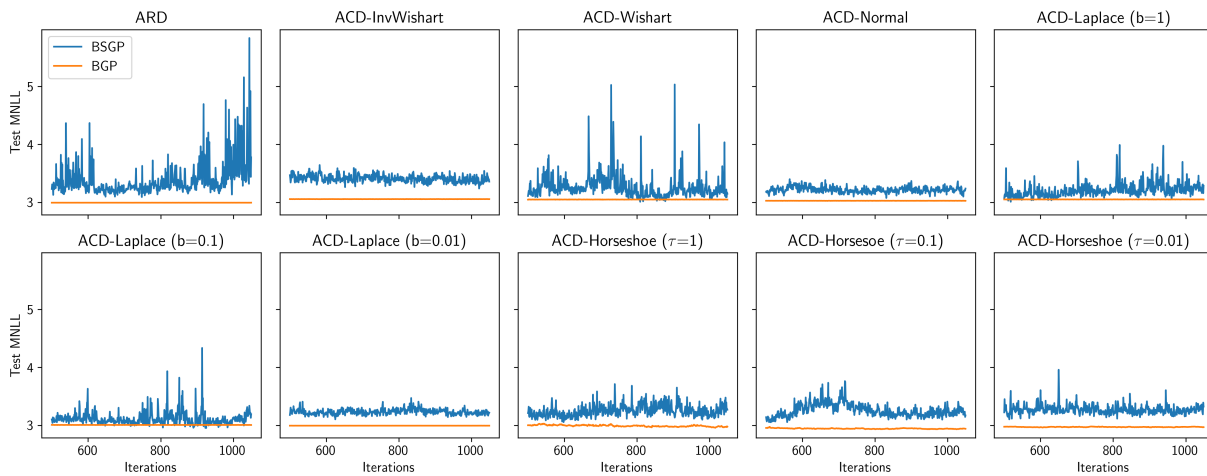


Figure 19: MNLL vs iterations for BSGP with 500 inducing points and for full GP on `concrete` dataset. The plots show one value every 10 of the 10,000 iterations.

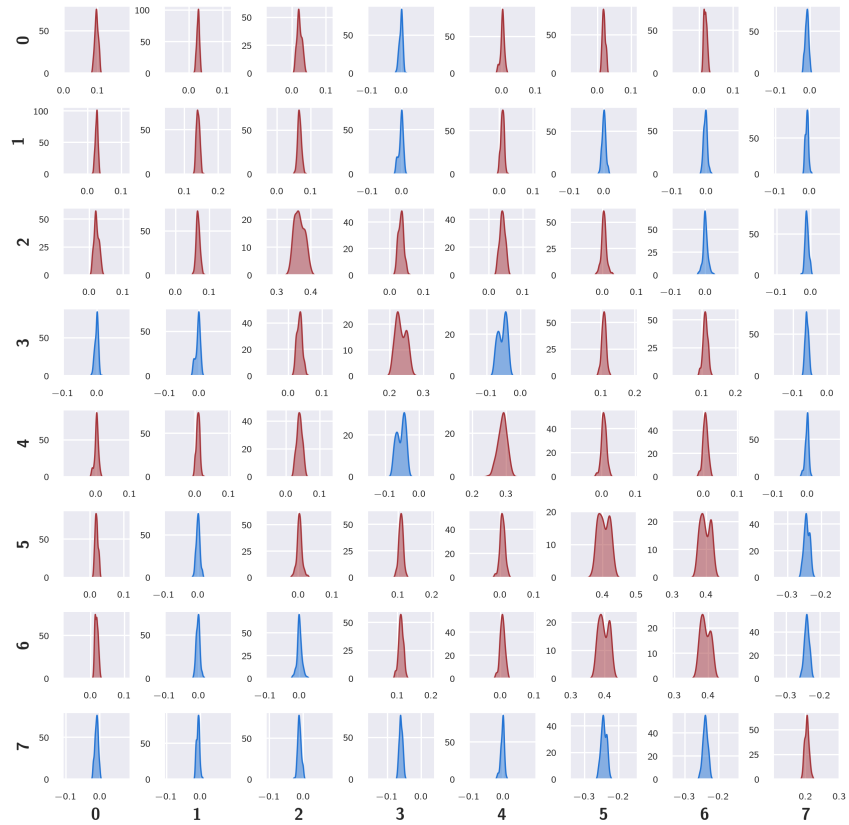


Figure 20: Posterior samples distribution of precision matrix entries for kin8nm dataset with Horseshoe ($\tau = 0.1$) prior.

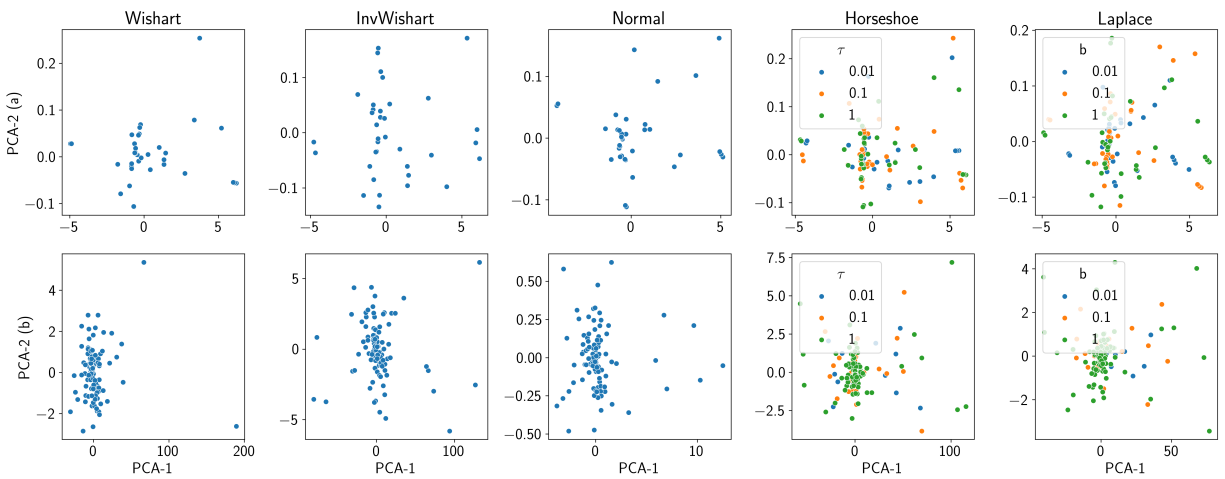


Figure 21: PCA representation of vectorized posterior precision matrices. Each point in the 2D space represents a posterior sample (precision matrix). (a) kin8nm dataset, (b) eeg dataset.