

A Bayesian Model of Dose-Response for Cancer Drug Studies

Wesley Tansey^{*1}, Christopher Tosh², and David M. Blei^{2,3,4}

¹Department of Epidemiology & Biostatistics, Memorial Sloan Kettering Cancer Center, New York, NY, USA

²Data Science Institute, Columbia University, New York, NY, USA

³Department of Statistics, Columbia University, New York, NY, USA

⁴Department of Computer Science, Columbia University, New York, NY, USA

Abstract

Exploratory cancer drug studies test multiple tumor cell lines against multiple candidate drugs. The goal in each paired (cell line, drug) experiment is to map out the dose-response curve of the cell line as the dose level of the drug increases. We propose Bayesian Tensor Filtering (BTF), a hierarchical Bayesian model for dose-response modeling in multi-sample, multi-treatment cancer drug studies. BTF uses low-dimensional embeddings to share statistical strength between similar drugs and similar cell lines. Structured shrinkage priors in BTF encourage smoothness in the dose-response curves while remaining adaptive to sharp jumps when the data call for it. We focus on a pair of cancer drug studies exhibiting a particular pathology in their experimental design, leading us to a non-conjugate monotone mixture-of-Gammas likelihood. To perform posterior inference, we develop a variant of the elliptical slice sampling algorithm for sampling from linearly-constrained multivariate normal priors with non-conjugate likelihoods. In benchmarks, BTF outperforms state-of-the-art methods for covariance regression and dynamic Poisson matrix factorization. On the two cancer drug studies, BTF outperforms the current standard approach in biology and reveals potential new biomarkers of drug sensitivity in cancer. Code is available at <https://github.com/tansey/functionalmf>.

1 Introduction

To search for new therapeutics, biologists carry out exploratory studies of drugs. They test multiple drugs, at different doses, against multiple biological samples. The goal is to trace the dose-response curves, and to understand the efficacy of each drug.

^{*}tanseyw@mskcc.org (corresponding author)

This article concerns dose-response modeling in exploratory drug studies. In particular, we are concerned with studies where the experimental design makes it difficult to perform statistical inference on the resulting data. We consider two such studies, both involving anti-cancer drugs being tested *in vitro* on models of human tumors known as organoids [10]. A dose-response curve in the studies represents the expected cell survival rate (response) for a specific organoid as a function of the concentration (dose) of an anti-cancer drug. The studies differ primarily in their size. The first study is a small-scale pilot study conducted internally at Columbia University Medical Center with 35 drugs and 28 organoids; the second is a large-scale, “landscape” study conducted at Samsung Medical Center with 67 drugs and 284 organoids [31]. The experiments in each study are costly; each one can take weeks or months to conduct in the lab. Consequently, the exhaustive set of all (organoid, drug) combinations is not available. This leaves missing data, dose-response curves for which no data is available, that must be imputed.

Figure 1 shows data from the landscape study. Each panel illustrates the interaction of one type of drug with one organoid sample. The gray points are the results of a set of experiments, each set with 2 replicates measured at 7 different doses. The goal is to use the observations (gray points) to infer the true dose-response curves. The predictions—the orange lines and uncertainty bands—come from the Bayesian dose-response model we propose in this paper; each of the 9 panels in fig. 1 were held out from the model at fitting time.

Notice there is structure in the outcomes: each drug has similar effects on each organoid. Thus, we treat modeling of dose-response as a factorization problem. The structure in the data arises because organoids share latent molecular attributes, such as genomic mutations, and drugs share latent pharmaceutical attributes, such as chemical structures. In each experiment, organoid and drug attributes interact, creating the shared patterns of dose-response.

While traditional factorization considers a matrix of scalars, the entries of this matrix are latent dose-response curves subsampled at different doses. To model such curves, we model drug attributes as evolving with the dose level. While the effects usually vary smoothly between dose levels, there are occasional sharp jumps, such as between the final two dose levels of drug 1. Capturing latent structure in dose-response curves requires handling this type of non-stationarity.

The observation model for these experiments is also non-standard. The outcome measured is a positive, real-valued measurement of cell survival relative to a noisy baseline. The model we propose uses a non-conjugate mixture-of-Gamma-shapes likelihood with the latent dose-response entering through the scale parameter. This is reflected in fig. 1, where the uncertainty intervals shrink as the predicted survival rate drops.

As a final wrinkle, the drugs in these cancer studies are all cytotoxic, meaning that they will only kill cells, not facilitate growth. This biological prior knowledge implies the expected survival rate can only decrease as the dose increases. Cytotoxicity adds a shape constraint requiring the dose-response curves are all monotonic. Further, since these drugs will not facilitate growth, effects in each curve are also upper bounded.

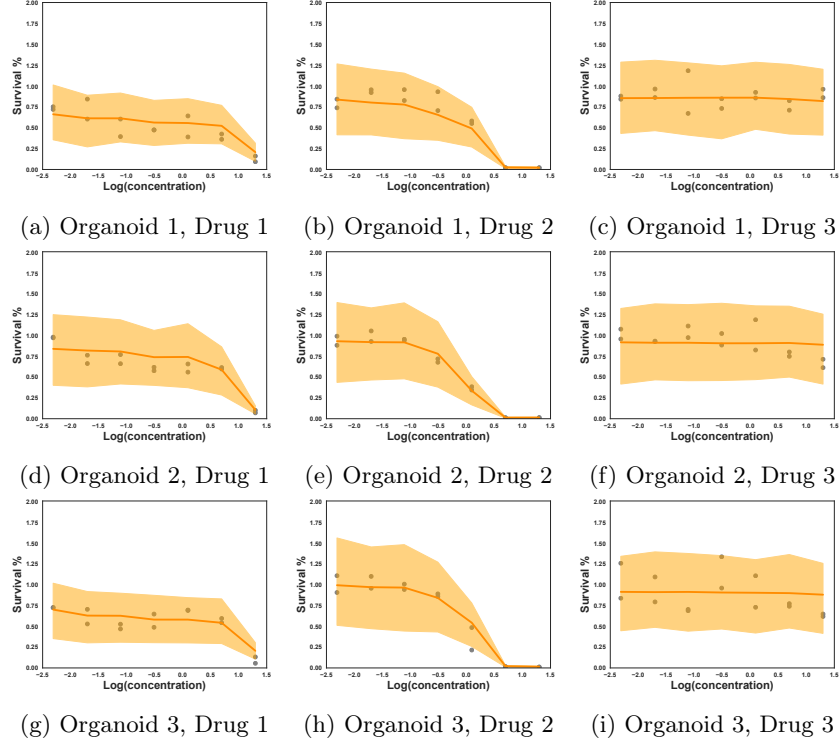


Figure 1: Sample of data from an organoid cancer drug experiment. Gray dots are observed outcomes; the orange line is the mean predicted response; bands represent 50% posterior predictive credible intervals for observations. All nine experiments were held out from the model at training time.

The dose-response model in this paper addresses these requirements. First, we propose Bayesian tensor filtering (BTF), a probabilistic method for smoothed tensor factorization. BTF uses structured shrinkage priors that encourage smoothness between successive dose levels, while simultaneously enabling sharp jumps when the data calls for it. Second, we develop generalized analytic slice sampling (GASS), a new MCMC inference algorithm that, when combined with BTF, enables the proposed dose-response model to support arbitrary likelihoods and linear inequality constraints on dose-response curves. The ability of GASS to handle linear constraints enables inference of monotone dose-response curves with upper and lower bounds.

The remainder of this paper is structured as follows. We first review related work on Bayesian dose-response modeling in Section 2. Section 3 then provides an overview of the two cancer drug studies and a motivation for the mixture-of-Gamma-shapes likelihood. Section 4.2 presents Bayesian tensor filtering, a flexible model for smoothed tensor factorization. In Section 5, we detail generalized analytic slice sampling, a procedure for sampling from posteriors with constrained

multivariate normal priors and non-conjugate likelihoods. Section 6 presents quantitative performance benchmarks for BTF, GASS, and the proposed dose-response model. Finally, in Section 7 the dose-response model is extended to handle side information in the form of molecular features about organoids. An analysis of the landscape study dataset with 115 features reveals potential new biomarkers of drug sensitivity in a subset of organoids.

2 Relevant literature

We survey a collection of the most relevant work to the proposed dose-response model. In each category of work, we focus on methods that enable uncertainty quantification, primarily through Bayesian inference.

Bayesian factor modeling Many models have been developed for Bayesian factor analysis with smooth structure. Zhang and Paisley [64] apply a group lasso penalty to the rows and columns of a matrix then derive a variational expectation maximization (EM) algorithm [5] for inference. Hahn et al. [20] use horseshoe priors [8] for sparse Bayesian factor analysis in causal inference scenarios with many instrumental variables. Kowal et al. [27] develop a time series factor model using a Bayesian trend filtering prior [12] on top of a linear dynamical system with Pólya–Gamma augmentation [48] for binomial observations. Schein et al. [50] develop Poisson-Gamma dynamical systems (PGDS), a dynamic matrix factorization model specifically for Poisson-distributed observations; we compare BTF with a tensor extension of PGDS in Section 6. Unlike the above models, BTF is likelihood-agnostic through GASS inference and enables modeling of independently-evolving columns rather than a common time dimension.

Independent dose-response curve estimation. A number of authors have investigated Bayesian methods for modeling monotone dose-response curves. These are typically done through a mixture of monotone functions. Perron and Mengersen [45] use a mixture of triangular distributions. Neelon and Dunson [43] use an autoregressive mixture prior of truncated normals in a piecewise linear spline model. Bornkamp and Ickstadt [6] propose a Bayesian nonparametric (BNP) model with a potentially-infinite mixture of two-sided power distributions. Shively et al. [51] also propose a BNP model which improves upon the model of Neelon and Dunson [43] with the key idea being to model the mean of the monotone curve as the integral of a positive function. Ghebretinsae et al. [18] present a Bayesian hierarchical model for nonnegative, real-valued comet assays with a Gamma outcome model on the shape. These methods all focus on the case of individual curve estimation. Lin and Dunson [34] propose a Gaussian process model with a posterior projection approach for shape-constrained curves. The datasets we consider here differ in that there are multiple samples and multiple drugs, with the goal to share statistical strength between samples and drugs to both denoise the existing curves and predict drug effects on samples without that specific (sample, drug) pair yet tested.

Joint dose-response curve estimation. A smaller body of work considers joint modeling of a set of dose-response curves. Both Vis et al. [58] and Abbas-Aghababazadeh et al. [1] use nonlinear mixed effect models for cancer drug response in a regression setting. Patel et al. [44] take a Bayesian B-spline approach to dose-response surface modeling. Fridley et al. [14] propose a hierarchical Bayesian log-linear model for dose-response in cytotoxicity studies. Fox and Dunson [13] consider the similar setting of covariance regression for influenza infection rates, imposing a sparse factorization on the covariance matrix that evolves over time in a Bayesian nonparametric setting. Wilson et al. [62] use monotone piecewise-linear splines in a hierarchical model of chemical toxicity assays, imposing a hierarchical Bayesian model that shrinks across similar molecules.

These models all have the common property that they shrink together samples via hierarchical priors, sharing statistical strength among rows to improve dose-response curve estimation. However, the hierarchical priors do not model any relational structure to shrink across samples and assays and do not provide any way to infer missing curves. Modeling the relational structure in multi-sample, multi-drug studies is crucial for cancer drug studies as often only a subset of samples have been tested for any given drug. Predictions about the missing curves can inform which experiments, among the many possible (sample, drug) combinations, show promise and should be carried out next. We discuss the deeper connections between BTF and both Fox and Dunson [13] and Wilson et al. [62] after presenting the details of BTF; we also compare against Fox and Dunson [13] with a Gaussian likelihood version of BTF in the benchmarks.

Predictive dose-response modeling. In many experiments, descriptive features are gathered representing useful side information about the samples or assays. A natural approach in these scenarios is to build predictive models that map from features to dose-response curves, enabling out-of-sample prediction for untested (sample, drug) pairs. Low-Kam et al. [36] proposed a Bayesian regression tree model with spline leaf nodes, enabling prediction of entire dose-response curves from chemical descriptors in nanoparticle experiments. Wheeler [61] modeled dose-response with molecular descriptors via additive Gaussian process tensor products over a real-valued feature space. In the case of binary or discrete data, Wheeler [61] first take a principal components decomposition to project features to a continuous space, making feature interpretation difficult. Both methods also assume a Gaussian noise model with no shape constraints; dealing with non-conjugate likelihoods and shape constraints would require a novel inference scheme similar to the proposed GASS algorithm. Tansey et al. [55] use deep neural networks to predict monotone dose-response curves from molecular features in an approximate Bayesian model, but require a large dataset of experiments and features to train the neural network. More generally, predictive methods generally assume features to be available and complete. In the pilot dataset, no features are present; in the landscape dataset, many samples are missing feature information. In section 7 we extend the dose-response model

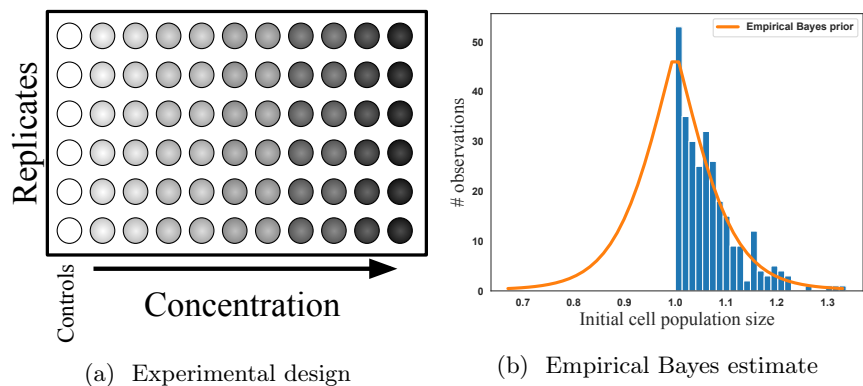


Figure 2: Left: The layout of each microwell plate experiment used to generate a single dose-response experiment. Cells are pipetted one column at a time, leading to correlated errors. Right: Estimate of the prior distribution of mean cell counts in each column, relative to the control column mean. The prior is estimated empirically assuming the lowest concentration had no effect if it had a higher mean.

to include potentially-missing features via a multi-view factorization approach.

3 Study design and dataset details

We detail the specific protocol used for the internal pilot study at Columbia University Medical Center. The landscape study we analyze uses a different number of plates, wells per plate, concentrations, drugs per plate, and replications. These distinctions are changes to the dimensions of the resulting tensor of observations, but the fundamental statistical inference challenges remain the same.

Each dose-response experiment is conducted on a microwell plate where a single drug is tested against a single biological sample. Each experiment measures a proxy for cell abundance 72 hours after applying the drug. Cell abundance is reported relative to a baseline control population where no drug was applied. For the control and each concentration level, 6 replicates are tested. Figure 2a shows the design of each 60-well plate experiment in the pilot study. The proxy measurement is a fluorescence assay; a fluorescent protein is added after 72 hours that binds to a molecule kept at near-constant levels in living cells. The degree of brightness of each microwell measures the relative abundance of cells alive but does not correspond to an exact cell count. All observations in the final dataset are normalized by dividing the brightness measurements in each microwell by the mean brightness for the control wells.

The details of how these organoid experiments were carried out in the wet lab are important, as they induce a particular form of correlated errors in the

observations. The first step in each experiment is to pipette an initial population of cells into each of the 60 microwells on the plate. This is a time consuming process for the biologist, often taking hours to pipette a single plate. To speed up the plating process, biologists use a multi-headed pipette that enables them to simultaneously fill an entire column of each plate. This reduces the burden on the biologist, but comes at a cost: correlated errors.

When a biologist fills a microwell, they first draw a pool of cells into the pipette. Given the small volumes involved in laboratory experiments, the actual number of cells drawn can vary substantially on a relative basis. Using a multi-headed pipette transforms this variation into a hierarchical model: first a pool of cells is drawn into the pipette, then it is split among all the heads. The majority of the variation comes in the initial sampling, with small noise added in the splitting process. This has the unintended side effect of creating correlated errors between all microwells in a single column. The drug concentrations are replicated along the same column, leading to a fundamentally unanswerable question: was the drug particularly effective at this dose or was the initial sample of cells particularly small by chance?

In other cancer drug experiments, such as high-throughput cancer cell line experiments [52, 15, 17], wells are filled using designs that avoid such correlated errors. In particular, these experiments typically use an orthogonal design, pipetting the drug replicates along rows and cells along columns. High-throughput screens exhibit their own correlated error pathologies, often referred to as batch effects. A number of techniques have been developed to preprocess high-throughput screening data and remove or limit batch effects [23, 32, 30, 38, 55]. These techniques account for the microenvironment similarities (e.g. temperature, humidity) that create spatial and temporal correlation between errors in wells nearby on the same plate, or run in the lab on the same day, in high-throughput screens.

Unfortunately, batch effect correction techniques are not applicable here. The correlation between wells is due to multi-headed pipetting, not similarities in the plate microenvironment between temporally- or spatially-related wells. Unlike smoothly-varying spatial batch effects, the experimental bias in the organoid datasets creates correlations between wells in the same column, but does not vary smoothly between columns. Drug concentrations are also varied across columns but constant across rows. This makes it impossible to disentangle the drug effect from the pipetting error without any assumptions.

4 Generative dose-response model for cancer organoid studies

To model the latent dose-response curves in the pilot and landscape studies we make two high-level design decisions. First, the unidentifiable experimental design necessitates making some assumption about the dose-response curves. We take an empirical Bayes approach, leveraging the fact that most drugs are ineffective at the smallest dose level. After specifying the likelihood, a hierarchical

structure is imposed to share statistical strength between similar organoids and similar drugs. We propose a Bayesian hierarchical model that imposes data-adaptive smoothness between successive doses and shares statistical strength via latent, low-dimensional embeddings.

4.1 Smallest-dose assumption and empirical Bayes likelihood estimation

The correlated errors in the columns render the exact effects unidentifiable. Each column has two latent variables affecting the final population size of cells: a dose-level effect from the drug and an initial population size from the pipetting. Since both of these variables affect all replicates in a column, disentangling them precisely is impossible.

We take an empirical Bayes approach to disentangling the variation in drug effects from the technical error in pipetting. In most experiments, the lowest concentration tested is too small to have any effect on cell survival. We therefore make the assumption that any experiment where the mean of the control replicates is lower than the mean of the replicates treated at the lowest concentration has effectively two sets of control columns. This enables estimation of the variation between means and an empirical Bayes prior for the pipetting error.

Specifically, we form a histogram of all lowest-concentration means greater than the control mean on the same plate. We then fit a Poisson GLM with 3 degrees of freedom to the histogram to estimate the prior probability that the mean of the initial population of cells was higher than the control mean. We assume the true distribution is symmetric and obtain an empirical Bayes prior on the means. Figure 2b shows the histogram and empirical Bayes prior estimate for the pilot study dataset. The within-column variance is identifiable and estimated using the control replicates.

We integrate out the uncertainty in the initial population mean. For each organoid sample $i = 1, \dots, N$ treated with drug $j = 1, \dots, M$ at dose level $t = 1, \dots, T$, with replicates $r = 1, \dots, R$, this yields a gamma mixture model likelihood,

$$P(y_{ijtr} \mid \mu_{ijt}) = \prod_{r=1}^R \left(\sum_{k=1}^K \hat{n}_k Ga(y_{ijtr}; \hat{a}_k, \hat{b}_k \mu_{ijt}) \right) \mathbb{1}[0 \leq \mu_{ijt} \leq 1], \quad (1)$$

where $(\hat{n}, \hat{a}, \hat{b})_k$ are derived from the empirical Bayes procedure. The latent drug effect $\mu_{ijt} \in [0, 1]$ is the probability of a cell in organoid i surviving treatment with drug j at dose level t . The drug effect enters in the Gamma scale as the initial cell population size sampled from $Ga(y_{ijtr}; \hat{a}_k, \hat{b}_k)$ is being multiplied by μ_{ijt} . The drug effect is constrained to be a proportion, as the drugs are known not to help any cells grow (i.e., the proportion can be at most 1) and a drug cannot kill more than all of the cells.

4.2 Bayesian tensor filtering for organoid dose-response modeling

To capture shared structure between dose-response curves of similar organoid samples and similar drugs, we place smoothed factor model priors on the dose-response curve in a hierarchical Bayesian model we term Bayesian tensor filtering (BTF),

$$\begin{aligned}
y_{ijtr} &\sim P(y_{ijtr} \mid \mu_{ijtr}) \mathbb{1}[\mu_{ijtr} \leq \mu_{ij(t-1)}] \\
\mu_{ijtr} &= w_i^\top v_{jtr} \\
w_i &\sim \mathcal{N}_D(\mathbf{0}, \sigma^2 I_D) \\
(\Delta^{(k)} V_j)_\ell &\sim \mathcal{N}_D(\mathbf{0}, \rho^2 \tau_{j\ell}^2 I_D) \\
\tau_{j\ell} &\sim \text{C}^+(0, \phi_{j\ell}) \\
\phi_{j\ell} &\sim \text{C}^+(0, 1) \\
\sigma^{-2} &\sim \text{Gamma}(0.1, 0.1) \\
\rho &\sim P(\rho).
\end{aligned} \tag{2}$$

In the above model, $w_i, v_{jt} \in \mathbb{R}^D$ are the latent factors and loadings for each organoid i and drug-dose (j, t) , respectively. The $T \times D$ drug loadings matrix V_j contains all (v_{j1}, \dots, v_{jT}) loading vectors for drug j . The choice of the number of latent factors, D , is a hyperparameter. We will occasionally refer to factors and loadings as embeddings or attributes, and D as the embedding dimension.

The generative model in eq. (2) incorporates a number of design decisions motivated by prior knowledge about biology and the nature of the experiments conducted. We explain the rest of the model in the context of these design decisions and the properties they induce in the resulting dose-response curves.

Monotonicity in the dose-response curve All drugs in both organoid studies are cytotoxic, only inducing a higher rate of cell death as the dose increases. This monotonic dose-response relationship is encoded by the hard constraint at the top of eq. (2). Each dose-response effect μ_{ijtr} is required to be at least as toxic as the previous $\mu_{ij(t-1)}$ effect; we assume $\mu_{ij0} = 1$.

Latent attributes for biological samples Organoid samples share molecular attributes. In cancer, different tumor samples contain similar patterns of genomic mutations, copy number alterations, and gene expression [60]. In mixed tissue experiments, cells that have differentiated into the same type will often respond similarly [e.g., 22]. These attributes are captured in BTF with a latent vector, $w_i \in \mathbb{R}^D$ for the i^{th} sample, as in standard matrix factorization. For identifiability [3], we assume a lower-triangular structure on the factors matrix $W = (w_1, \dots, w_N)$, though interpretable factors is not our primary goal here.

Independent dose-specific latent attributes for drugs For each drug in the dataset, BTF models each dose level $t = 1, \dots, T$ with its own embedding.

For a single drug j , there are T embeddings forming a drug embedding matrix $V_j \in \mathbb{R}^{T \times D}$. In BTF, columns (drug effects) are evolving independently, though potentially with similar latent attributes. This column independence distinguishes BTF from time-series tensor factorization models [63, 53, 16, 54] where all columns are progressing through time together. Independent column evolution captures the notion that two drugs treated at the same concentration may have totally different effects due to the molecular size of the drug, its targeting receptor, and its chemical structure. Different drugs may also be treated at entirely different concentration levels, as they have different molecular properties that require higher or lower concentration ranges to map out the dose-response relationship. For ease of notation, we assume all drugs are treated at T concentrations, however conceptually drugs could be treated at different numbers T_j of concentrations.

Non-stationary group smoothness priors on drug attributes We assume drug effects typically vary smoothly with dose, with occasional sharp jumps. To encode this, we place hierarchical smoothness priors on the differences between dose-specific drug embeddings. Specifically, we place a Bayesian group trend filtering prior on drug embeddings.

Trend filtering is an adaptive smoothing technique originally developed in the penalized regression case [24, 56]. The penalized regression formulation places ℓ_1 penalties on the k^{th} -order differences of neighboring points on a 1d grid. For instance, in the $k = 1$ case this is the total variation norm penalty,

$$\underset{\beta \in \mathbb{R}^n}{\text{minimize}} \ ||Y - \beta||_2^2 + \lambda \ ||\beta_{1:(n-1)} - \beta_{2:n}||_1 \quad (3)$$

The solution to the convex optimization problem in eq. (3) leads to piecewise constant plateaus in β due to the lasso penalty driving first differences to zero; higher order differences lead to piecewise-polynomial solutions [56].

Faulkner and Minin [12] extended trend filtering to the Bayesian context and considered a number of different priors on the k^{th} -order differences. They consider three different priors on the differences: normal priors, equivalent to ℓ_2 Laplacian smoothing in the regression case; Laplace priors, the direct Bayesian analog of ℓ_1 lasso priors; and horseshoe priors [8], a heavier tailed distribution that does not suffer from the non-diminishing bias of the lasso [57]. The horseshoe priors are shown to perform best across a range of problems. We adapt the horseshoe prior to the group trend filtering case in BTF.

We call $\Delta^{(k)} \in \mathbb{R}^{L \times T}$ the composite trend filtering matrix; it contains all linear operators needed to encode the $(1, \dots, k)^{\text{th}}$ -order differences. The ordinary trend filtering matrix encodes only the k^{th} -order differences, implicitly assuming all lower-order differences are not smooth. For example, the $k = 2$ case yields a

prior on the first and second order differences,

$$\Delta^{(2)} = \begin{bmatrix} 1 & 0 & 0 & 0 & \dots & 0 & 0 & 0 \\ 1 & -1 & 0 & 0 & \dots & 0 & 0 & 0 \\ 0 & 1 & -1 & 0 & \dots & 0 & 0 & 0 \\ & & & \dots & & & & \\ 0 & 0 & 0 & 0 & \dots & 0 & 1 & -1 \\ 1 & -2 & 1 & 0 & \dots & 0 & 0 & 0 \\ 0 & 1 & -2 & 1 & \dots & 0 & 0 & 0 \\ & & & \dots & & & & \\ 0 & 0 & 0 & 0 & \dots & 1 & -2 & 1 \end{bmatrix}. \quad (4)$$

The first line of eq. (4) places an independent prior on the embedding vector for the first dose level in each drug, v_{j1} , making the matrix $(\Delta^\top \mathcal{T} \Delta)$ non-singular, where $\mathcal{T} = \text{diag}(1/(\rho^2 \tau_j^2))$. This ensures the resulting prior on V_j is proper; see the supplementary material for details.

We impose a group trend filtering prior on the drug effects by placing a group horseshoe prior on the ℓ^{th} row of the $(\Delta^{(k)} V_j)$ differences matrix. To do this, in eq. (2) we adapt the Bayesian formulation of the group lasso [29] to the global-local shrinkage view of the horseshoe prior [47]. Each row $(\Delta^{(k)} V_j)_\ell$ in the differences matrix has both a local $\tau_{j\ell}^2$ variance and a global ρ^2 variance term. Small values of ρ^2 and $\tau_{j\ell}^2$ will shrink the ℓ^{th} difference vector to nearly zero, resulting in the curve being smoother; larger values enable the curve to jump in response to the data.

Following Bhadra et al. [4], we place a half-Cauchy prior on $\phi_{j\ell}$, the scale term in the local horseshoe shrinkage prior; this is referred to as the horseshoe+ prior. A full Bayesian specification could choose a reasonable prior for ρ , such as a standard Cauchy or Uniform(0, 1). If an estimate of the number of non-zero entries is available, Van Der Pas et al. [57] make an asymptotic argument for setting $\hat{\rho}$ to the expected number of non-zeros. We find BTF is robust to the choice of global shrinkage parameter; we default to a half-Cauchy prior on ρ in our implementation and also support performing a grid search over a range of discrete ρ values via deviance information criteria [9]. The value k in $\Delta^{(k)}$ is left as a hyperparameter; we suggest $k = 2$ as a reasonable default choice for most datasets.

4.3 Deeper connections to related work

The BTF model is closely related to two works: the Bayesian nonparametric covariance regression model of Fox and Dunson [13] and the zero-inflated piecewise log-logistic dose-response model of Wilson et al. [62]. Before proceeding to inference in BTF, we first provide a detailed discussion of the deeper connections to these methods.

Bayesian covariance regression. Fox and Dunson [13] introduced Bayesian nonparametric covariance regression (BNP-CovReg). The BNP-CovReg model

poses a Gaussian noise model for a set of p curves observed at n points,

$$\mathbf{y}_i = \boldsymbol{\mu}(\mathbf{x}_i) + \boldsymbol{\epsilon}_i, \quad \boldsymbol{\epsilon}_i \sim \mathcal{N}_p(0, \Sigma), \quad i = 1, \dots, n, \quad (5)$$

where $\mathbf{y}_i = \log \mathbf{r}_i$, the vector of observations in the p curves at point x_i , $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_p)$. A factor model is imposed on the latent mean,

$$\mathbf{y}_i = \mathbf{\Lambda}(\mathbf{x}_i)\boldsymbol{\eta}_i + \boldsymbol{\epsilon}_i, \quad \boldsymbol{\eta}_i \sim \mathcal{N}_k(\boldsymbol{\psi}(\mathbf{x}_i), I_k) \quad (6)$$

where $\mathbf{\Lambda}(\mathbf{x}_i)$ are the factor loadings at point \mathbf{x}_i and $\boldsymbol{\eta}_i$ the latent factors associated with observation \mathbf{y}_i . Here, $k \ll p$ imposes a low-rank assumption on the factor model and independent Gaussian process priors are placed on each ψ_h , for $h = 1, \dots, k$ with squared exponential kernel. Independent conjugate inverse-gamma priors are placed on each σ_i .

For computational feasibility, the factor loadings matrix $\mathbf{\Lambda}(\mathbf{x})$ is expressed as a weighted combination of a smaller set of L basis functions,

$$\mathbf{\Lambda}(\mathbf{x}) = \boldsymbol{\Theta}\boldsymbol{\xi}(\mathbf{x}), \quad (7)$$

where $\boldsymbol{\Theta}$ is a $p \times L$ matrix of coefficients and $\boldsymbol{\xi}(\mathbf{x})$ an $L \times k$ array of basis functions. A global-local shrinkage prior is placed on the elements of $\boldsymbol{\Theta}$ to effectively reduce the dimension of the basis to much smaller than L or k .

A number of follow-up works have investigated similar models. Kuniyama et al. [28] extend BNP-CovReg to longitudinal data with covariate information. Li et al. [33] use fixed factors, sacrificing the flexibility of the nonparametric approach of Fox and Dunson [13] for increased scalability. Heaukulani and van der Wilk [21] derive a variational inference approach for inverse Wishart processes that leads to a scalable approximate inference scheme for BNP-CovReg.

In the Gaussian likelihood case, BTF also uses a low-dimensional factor model for the response mean. However, rather than assuming independent sparsity and smoothness priors on a set of basis coefficients and latent factors, BTF imposes smoothness directly on the curves. This translates to a *group* smoothness assumption on the latent factors, enforcing that the k^{th} -order differences in successive latent means be shrunk to zero. As we show in section 6.1, BTF outperforms the BNP-CovReg model on the same dataset that motivated the design of BNP-CovReg. It also enjoys substantially faster computational times (around 10x faster on a 2018 MacBook Pro) and trivial parallelization in the Gibbs sampler if further scalability is needed. Furthermore, BTF is extensible to a number of other likelihoods, including non-conjugate models with linear constraints like monotonicity, through the GASS inference algorithm.

Piecewise log-logistic, monotone dose-response modeling Wilson et al. [62] introduced the zero-inflated piecewise log-logistic (ZIPLL) model for estimating dose-response curves in chemical toxicity assays. Similar to the cancer drug study scenario, the ZIPLL model is explicitly designed for multi-sample, multi-assay studies where observations form a tensor with concentration as a

third dimension. The ZIPLL model assumes a Gaussian noise distribution on observations y_{ijt} of sample i , assay j , concentration x_{ijt} , respectively,

$$y_{ijt} = f_{ij}(x_{ijt}) + \epsilon_{ijt}, \quad \epsilon_{ijt} \sim \mathcal{N}(0, \sigma^2). \quad (8)$$

In the toxicity study considered, the minimum and maximum values corresponding to no-effect and total toxicity are unknown and modeled as latent variables,

$$f(x_{ijt}) = \begin{cases} u_{ij} - (u_{ij} - l_{ij}) \times \text{Logistic}(g(x_{ijt}; a_{ij}, \mathbf{w}_{ij})) & \text{if } Z_{ij} = 1 \\ l_{ij} & \text{if } Z_{ij} = 0 \end{cases}, \quad (9)$$

where u_{ij} is the maximum response, l_{ij} is the minimum response, a_{ij} is the location parameter of the curve (the 50% survival point), and \mathbf{w}_{ij} are the shape parameters. The Z_{ij} variable captures the zero-inflated property of the chemicals in the toxicity dataset, where sparsity in effects is expected. The ZIPLL model uses a monotonic log-linear spline basis to model $g(x_{ijt}; a_{ij}, \mathbf{w}_{ij})$ using a fixed symmetric grid of internal knots. Multivariate normal priors are placed on $\boldsymbol{\theta}_{ij} = \log(l_{ij}, u_{ij}, a_{ij})$ and a Gaussian autoregressive prior is placed on $\log \mathbf{w}_{ij}$ to encourage smoothness in the dose-response curve. Shrinkage in ZIPLL is performed across samples in the same assay via a hierarchical prior,

$$\boldsymbol{\theta}_{ij} \mid \boldsymbol{\mu}_j, \Sigma_j \sim \mathcal{N}_3(\boldsymbol{\mu}_j, \Sigma_j), \quad \boldsymbol{\mu}_j \sim \mathcal{N}_3(\boldsymbol{\mu}, \Sigma), \quad (10)$$

where hard-coded hyperparameters are used for conjugate priors on $\boldsymbol{\mu}$, Σ_j , and Σ , as well as the other hyperparameters in the hierarchical model.

The ZIPLL model is closely related to the BTF dose-response model, in both modeling constraints and target application domain. The use of a monotone basis with Gaussian priors on log-transformed latent parameters ensures the posterior curves satisfy the monotonicity constraints in cancer dose-response modeling, and the logistic transform maps to the $[0, 1]$ interval as well. However, the Gaussian noise model is also misaligned with the cancer drug studies likelihood, and the stationary smoothness assumptions of the Gaussian autoregressive prior may not handle to sharp jumps found in the data.

Moreover, ZIPLL only pools statistical strength across samples (rows) in the dose-response tensor but not assays (columns); this prevents imputing assays that are missing entirely. Extending the comparatively simple scalar basis coefficients \mathbf{w}_{ijl} with autoregressive priors to a factor model with equivalent smoothness would be nontrivial, as posterior inference is non-conjugate and may require a novel approximate inference scheme. These are similar challenges to the inference in the BTF dose-response model, all of which led us to the development of the GASS algorithm in section 5. The ability to impute out of sample experiments is critical for the cancer organoid drug studies, where only a subset of samples are tested for each drug.

5 Posterior inference

Posterior inference in BTF is performed through Gibbs sampling. In its most basic form, Gibbs sampling requires us to sample from the conditional distributions for each of the parameters. The updates for the latent attributes W and V depend on the form of likelihood, $P(y_{ijtr}; w_i^\top v_{jt})$. The derivations for Gaussian and binomial likelihoods, as well as the horseshoe parameter updates, are provided in the supplementary material. Here we focus on the crux of the posterior inference challenge in BTF: the constrained, non-conjugate likelihood in eq. (1).

To sample from the complete conditional for each of the latent attributes, we run MCMC-within-Gibbs – running a separate Markov chain whenever a Gibbs step requires a sample from the conditional distributions of the latent attributes w_i and v_{jt} . The challenge to this step is that the likelihood imposes hard constraints on the values of entries. The inner product must be a probability, requiring $w_i^\top v_{jt} \in [0, 1]$. Each inner product must also be no greater than the previous inner product, requiring $w_i^\top (v_{jt} \leq v_{j(t-1)}) \leq 0$. Running a naive MCMC algorithm such as Metropolis Hastings within the Gibbs sampler is likely to have a high rejection rate and lead to poor mixing. Instead, we develop an MCMC algorithm that is capable of directly handling generic likelihoods and arbitrary linear constraints.

5.1 Generalized analytic slice sampling

Sampling from the conditional distributions of the latent attributes can be reduced to the problem of sampling from the posterior of a vector x with a multivariate normal prior constrained by a set of linear inequalities,

$$x \sim P(y; x) \text{MVN}(x; \mu, \Sigma) \mathbb{I}[Dx \geq \gamma]. \quad (11)$$

Note in eq. (11) we are describing a generic problem that is distinct from eq. (2). The variables $(y, x, \mu, \Sigma, D, \gamma)$ in eq. (11) are correspondingly generic variables that are distinct from those used in eq. (2).

A natural candidate for sampling from such a distribution is elliptical slice sampling [40], an empirically successful exact MCMC method for sampling from distributions with arbitrary likelihoods and multivariate normal priors. Elliptical slice sampling builds on the idea of slice sampling [42], an MCMC algorithm that, given a point x , samples a new point x' by first drawing a value u uniformly over the range 0 up to the likelihood of x , and then by drawing x' uniformly from the set of points whose likelihood is at least u . Computing the set of points whose likelihood is above a certain threshold is infeasible in general, and thus some form of rejection sampling is required. Unfortunately, in high dimensions, these rejection rates will be very large.

However, when we have a multivariate normal prior, we may utilize the fact that the contours of equal probability on the multivariate normal distribution are elliptical regions. Elliptical slice sampling exploits this observation by sampling a point v from the prior distribution and computing the ellipse $\{x \cos(\theta) + v \sin(\theta) : \theta \in [-\pi, \pi]\}$ containing both x and v . Then, as in slice sampling, it samples

Algorithm 1: Generalized analytic slice sampling (GASS) for constrained MVN priors

Data: Valid current point x , mean μ , covariance Σ , log-likelihood \mathcal{L} , constraints (D, γ)

Result: MCMC sample from

$$P(x') \propto \exp(\mathcal{L}(x')) \text{MVN}(x'; \mu, \Sigma) \mathbb{I}[Dx' \geq \gamma]$$

$t = \mathcal{L}(x) + \log \epsilon$, $\epsilon \sim U(0, 1)$;

Sample proposal $v \sim \text{MVN}(v; \mathbf{0}, \Sigma)$;

Grid approximation $\mathcal{G} = \text{grid}(-\pi, \pi)$;

foreach constraint $(d_i, \gamma_i) \in (D, \gamma)$ **do**

$a = d_i^\top (x - \mu)$, $b = d_i^\top v$, $c = \gamma_i - d_i^\top \mu$;

if $a^2 + b^2 - c^2 \geq 0$ and $a \neq -c$ **then**

 Get θ_1, θ_2 as in eq. (12);

if $a^2 > c^2$ **then**

$\mathcal{G} = \mathcal{G} \cap [\theta_1, \theta_2]$;

else

$\mathcal{G} = \mathcal{G} \cap ([-\pi, \theta_1] \cup [\theta_2, \pi])$;

end

end

end

Generate candidate samples $\mathcal{X} = \{x' : x \cos(\theta_g) + v \sin(\theta_g) + \mu, \theta_g \in \mathcal{G}\}$;

Select uniformly from sufficiently likely candidates

$\{x' : \mathcal{L}(x') \geq t, x' \in \mathcal{X}\}$.

a likelihood u and then performs a form of rejection sampling to sample x' uniformly from the set of points *on the ellipse* whose likelihood is at least u . Note that when the likelihood is reasonably smooth, there will always be a reasonably large interval of points around x on the interval above this likelihood, and we will not have too many rejections. However, when the likelihood has hard constraints, these regions can be very small and lead to high rejection rates.

To address this, we extend elliptical slice sampling to directly handle constrained multivariate normal priors. The approach, which we call generalized analytic slice sampling (GASS), is a natural extension of the analytic slice sampling procedure of Fagan et al. [11] for truncated multivariate normals. The key difference is that the original analytic slice sampler only considered centered truncated multivariate normals with no likelihood component. Generalizing this procedure to handle the more general case in eq. (11) requires handling several edge cases.

5.2 Algorithm

The full GASS procedure is presented in Algorithm 1. The idea of GASS is to note that the constraints can be pushed inside the proposal update. Namely, given an ellipse and a set of linear constraints, it is relatively straightforward to

exactly compute their intersection, and thus to restrict proposals to that region.

To see why, consider the case where we have a single linear constraint requiring that the output point satisfies $d^\top x' \geq \gamma$. Then a valid angle θ must satisfy $a \cos \theta + b \sin \theta - c \geq 0$, where $a = d^\top (x - \mu)$, $b = d^\top (v - \mu)$, and $c = \gamma - d^\top \mu$. Basic trigonometry implies that the feasible range of θ is a subset of $[-\pi, \pi]$ whose boundary points are given by

$$\theta_1, \theta_2 = 2 \arctan \left(\frac{b \pm \sqrt{a^2 + b^2 - c^2}}{a + c} \right). \quad (12)$$

There are two edge cases where the entire ellipse is valid: (i) $(a^2 + b^2 - c^2) < 0$ and (ii) $a = -c$. In the first case, we trivially have $a^2 + b^2 < c^2$ and therefore $a \cos \theta + b \sin \theta > c$ for all θ . In the second case, the only place the constraint boundary intersects the ellipse exactly at a single point of the ellipse and thus its selection has probability zero. For all other cases, the subset is determined based on the sign of $a^2 - c^2$. A positive sign indicates the quadratic in the inequality is concave and eq. (12) defines the boundaries of a contiguous region; a negative sign indicates convexity and thus the complement of the interval.

When there are many linear constraints, we can solve for the valid regions of each of the individual constraints separately and then take their intersection. Finally, after computing the region of valid θ 's, we approximate it with a fine-grained 1D grid. We draw a likelihood value u and filter out all the grid points with likelihood smaller than u . Finally, we draw a sample uniformly over the remaining grid points. An illustration of the algorithm is given in fig. 3.

It is not difficult to show that GASS is a valid Markov chain which will converge to the distribution given in eq. (11). For completeness, a proof of this convergence is provided in the supplementary material.

5.3 Conditioning heuristic

Elliptical slice sampling schemes like GASS can suffer from poor mixing when the likelihood overwhelms the multivariate normal prior. In such settings, the sampled ellipses, which are generated with respect the prior and not the likelihood, may only have a small region centered around the current sample with likelihood comparable with the current likelihood, causing the chain to take only very small steps. Motivated by this observation, Fagan et al. [11] suggest performing expectation propagation [39] to better align the prior with likelihood. Unfortunately, in the context of BTF this is impractical as the prior parameters for W are a function of V , and vice versa, which would require us to perform expectation propagation every iteration of the Gibbs sampler.

We instead approximate the entire tensor once at the start by a nonnegative, monotone tensor factorization. To do this, we find an approximate solution to the optimization problem,

$$\begin{aligned} \hat{W}, \hat{V} = & \underset{W, V}{\text{minimize}} && \sum_{ijtr} (y_{ijtr} - w_i^\top v_{jt})^2 \\ & \text{subject to} && 0 \leq w_i^\top v_{jt} \leq 1, \\ & && w_i^\top (v_{jt} - v_{j(t-1)}) \leq 0 \end{aligned}, \quad (13)$$

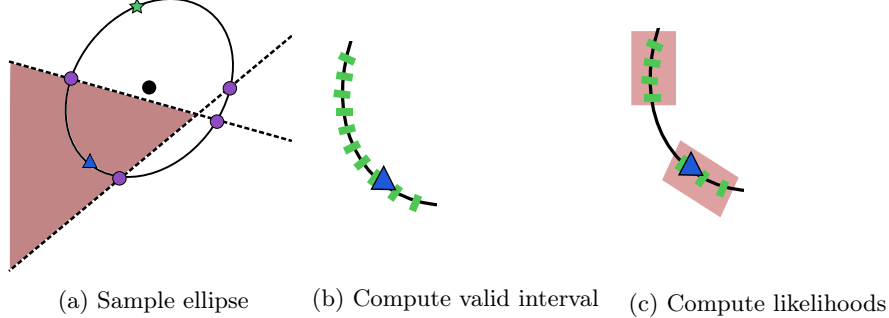


Figure 3: An illustration of one step of the GASS chain starting from the blue triangle. (a) First, the green star is sampled from the multivariate normal centered at the black circle. The green star and the blue triangle determine an ellipse around the black circle. Here the dashed lines denote the linear constraints, and the red region is the feasible region. (b) Then, the intersection of the feasible region with the ellipse is computed and gridded. (c) Finally, a likelihood u is sampled, and those points that have likelihood at least u are retained, denoted by the red region. The next state is randomly selected from the remaining grid points.

The approximation solution to eq. (13) is found via alternating constrained minimization for the rows and columns; we run the alternating minimization procedure until convergence.

After fitting the rows and columns, we calculate an over-estimate of the variance, analogous to an EP approximation, as a multiple of the empirical squared error in the estimate for each column and row,

$$\begin{aligned}\hat{\mu}_{ijt} &= \hat{w}_i^\top \hat{v}_{jt} \\ \hat{s} &= \frac{\sum_{ijtr} (y_{ijtr} - \hat{\mu}_{ijt})^2}{N \times M \times T \times R} \\ \hat{\Sigma} &= c\hat{s}I,\end{aligned}\tag{14}$$

where $c \geq 1$ is a hyperparameter. This over-estimates the empirical variance, accounting for a wider range of possible samples to correct for the error in the NMF procedure. The main sensitivity of BTF inference is initialization and conditioning with an accurate mean; the covariance is less important. A reasonably accurate $\hat{\mu}$ yields the majority of the gains, whereas the method is insensitive to reasonable choices of c . As we discuss in the benchmarks, the NMF model produces a reasonably accurate $\hat{\mu}$ in terms of RMSE, and serves as a good starting point for the BTF dose-response model; for all BTF benchmarks, we set $c = 3$.

The BTF dose-response model uses the conditioning heuristic in the W and V steps in the Gibbs sampler to calculate an adjusted prior. The log-likelihood used in the GASS procedure is then the original log-likelihood minus the log-

Google Flu Trends				
Model	In-sample		Out-sample	
	MAE	RMSE	MAE	RMSE
BNP-CovReg	0.31	0.41	0.31	0.42
Gaussian BTF (d=2)	0.15	0.22	0.15	0.21
Gaussian BTF (d=5)	0.10	0.14	0.12	0.16
Gaussian BTF (d=10)	0.07	0.10	0.13	0.21

Table 1: Posterior mean results on the Google Flu Trends dataset. The Gaussian BTF model outperforms the BNP-CovReg model of Fox and Dunson [13] with as few as 2 latent factors. MAE: mean absolute deviation from data; RMSE: root mean-squared error from data.

conditioning likelihood, leaving the resulting distribution equivalent but better aligning the prior and likelihood.

6 Benchmarks and performance comparisons

We study the performance of the proposed dose-response model and its components, BTF and GASS. We first benchmark GASS against different alternative methods for nonconjugate inference, where GASS mixes faster and has lower error. Then we study BTF on a dynamic matrix factorization problem with non-conjugate Poisson observations; BTF outperforms a recent Bayesian tensor decomposition approach designed for time-evolving count matrices. Finally, we apply the dose-response model to a real cancer drug study. We run 5 independent trials, holding out a different subset of entire dose-response curves and report averages over all trials; the BTF-based dose-response model outperforms all baselines in terms of log probability on held out data.

6.1 Gaussian BTF benchmarks

We benchmark the conjugate Gaussian-likelihood BTF model on the Google Flu Trends dataset¹ modeled in Fox and Dunson [13]. This dataset contains weekly influenza-like infection (ILI) counts from 183 different regions in the United States from 1996 to 2014. The set of regions contains nested information, including both major cities and entire states. For benchmarking purposes, we focus only on the 50 states to ensure that held out data is not leaked in through other nested regions; this makes the inference task strictly more difficult.

Fox and Dunson [13] model the weekly log-count of infections with the Gaussian noise model in eq. (5)–eq. (6), where $\mathbf{y}_i = \log \mathbf{r}_i$ is the vector of log Google-estimated ILI rates in the 50 states at time x_i .

¹<http://www.google.org/flutrends/>

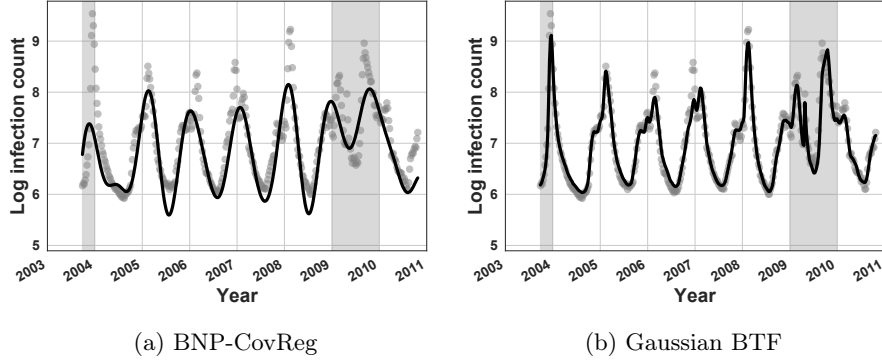


Figure 4: Example fits for BNP-CovReg and Gaussian BTF ($d = 5$) on a single state (Alabama) in the Google Flu Trends dataset. The gray shaded regions represent held out periods. BNP-CovReg over-smooths, failing to capture large peaks and short movements. BTF fits the data tightly both on in-sample examples and imputed held out periods, suggesting it has captured latent structure between states.

We compare against the performance of the Bayesian Nonparametric Covariance Regression (BNP-CovReg) model of Fox and Dunson [13], designed specifically for the Google Flu Trends data. For BNP-CovReg, we keep the same hyperparameter settings with truncation parameters $\bar{L} = 10$ and $\bar{k} = 20$; we use the reference implementation provided by the authors. For Gaussian BTF, we initialize the global shrinkage parameter ρ^2 to 0.1 and sample it with a HS+ prior; we also place an weakly-informative inverse-Gamma(0.1, 0.1) prior on the likelihood variance. To measure performance, we hold out 10% of all years, selected uniformly at random across all available state-years. Model performance is measured in both root mean squared error (RMSE) and mean absolute error (MAE) on held out data.

Table 1 shows the results for both the BNP-CovReg model and Gaussian BTF with $d = 2, 5$, and 10 latent factors. The BTF method outperforms BNP-CovReg in each case. The model performs best out of sample with $d = 5$ latent factors, suggesting it overfits with $d = 10$ factors. The BTF model also has good coverage, with 95% credible intervals having 95.83% coverage on in-sample data and 92.82% coverage for held out data with $d = 5$ factors.

Figure 4 shows an example of a single state (Alabama) comparing BNP-CovReg and $d = 5$ BTF. The BNP-CovReg model over-smooths, leading it to underestimate large peaks both in- and out-of-sample. The BTF model closely tracks the data, even in out-of-sample predicted weeks, suggesting it has learned latent structure between the states.

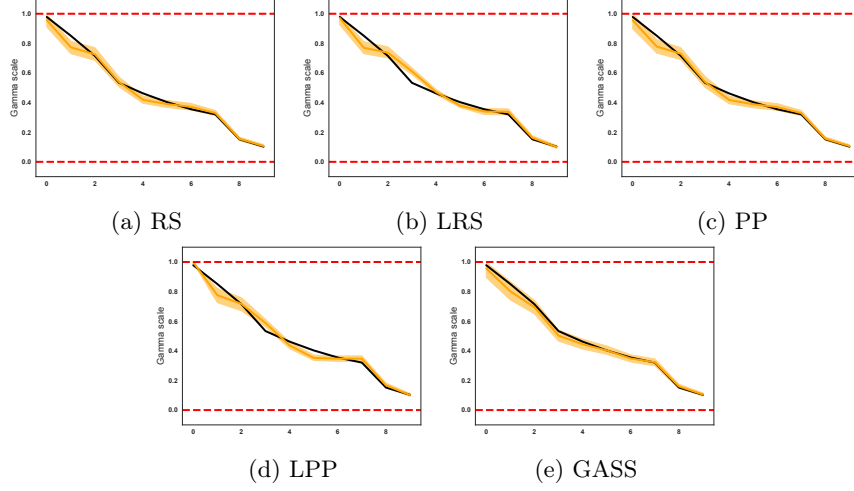


Figure 5: Sample fits for different methods on the gamma scale estimation benchmark; black is the true scale, orange is the estimated scale, bands are 90% credible intervals, dashed lines are constraint boundaries. GASS captures the shape of the curve and has good coverage after 10K Gibbs iterations.

6.2 GASS benchmarks

We benchmark GASS on a simulation study with a constrained multivariate normal prior with a non-conjugate gamma scale likelihood,

$$\begin{aligned}
 y_i^{(r)} &\sim \text{Gamma}(y_i; a, \theta_i) \\
 \boldsymbol{\theta} &\sim \text{MVN}(\boldsymbol{\theta}; \boldsymbol{\mu}, \Sigma) \mathbb{1}[0 \leq \boldsymbol{\theta} \leq 1] \prod_{i=1}^{n-1} \mathbb{1}[\theta_i \geq \theta_{i+1}] \\
 \boldsymbol{\mu} &= [0.95, 0.8, 0.75, 0.5, 0.29, 0.2, 0.17, 0.15, 0.01, 0.0001] \\
 \Sigma_{ij} &= \tau \exp\left(-\frac{1}{2b}(i-j)^2\right).
 \end{aligned} \tag{15}$$

The covariance matrix in the unconstrained prior corresponds to a squared exponential kernel. We set the hyperparameters $a = 100$, $\tau = 0.1$, and $b = 3$; all hyperparameters are assumed known. We use $R = 3$ replicates for \mathbf{y} . We compare GASS against four different variants of elliptical slice sampling (ESS),

- **Rejection sampling (RS).** Samples are drawn using the unconstrained ESS model, with the constraints pushed into the likelihood. Any violated constraint generates a zero probability and corresponds to a rejection sampler.
- **Logistic rejection sampling (LRS).** The ESS is used to model logits, which are then passed through the logistic transform to satisfy the $[0, 1]$ constraint; rejection sampling again handles the monotonicity constraint.

MSE ($\times 10^3$)					
Sampler	$m = 100$	$m = 500$	$m = 1000$	$m = 5000$	$m = 10000$
RS	1.54 ± 0.09	1.43 ± 0.09	1.43 ± 0.09	1.41 ± 0.09	1.41 ± 0.09
LRS	1.76 ± 0.14	1.75 ± 0.14	1.73 ± 0.14	1.66 ± 0.13	1.63 ± 0.13
PP	1.35 ± 0.09	1.31 ± 0.08	1.29 ± 0.09	1.29 ± 0.09	1.28 ± 0.08
LPP	1.66 ± 0.12	1.62 ± 0.12	1.56 ± 0.12	1.52 ± 0.13	1.50 ± 0.12
GASS	0.74 ± 0.05	0.66 ± 0.04	0.63 ± 0.04	0.52 ± 0.03	0.49 ± 0.03

90% Credible Interval Coverage					
Sampler	$m = 100$	$m = 500$	$m = 1000$	$m = 5000$	$m = 10000$
RS	0.37 ± 0.02	0.47 ± 0.02	0.49 ± 0.02	0.50 ± 0.02	0.52 ± 0.02
LRS	0.23 ± 0.01	0.34 ± 0.02	0.36 ± 0.02	0.46 ± 0.02	0.48 ± 0.02
PP	0.44 ± 0.02	0.54 ± 0.02	0.56 ± 0.02	0.57 ± 0.02	0.58 ± 0.02
LPP	0.30 ± 0.01	0.40 ± 0.02	0.44 ± 0.02	0.53 ± 0.01	0.53 ± 0.02
GASS	0.58 ± 0.02	0.73 ± 0.02	0.77 ± 0.02	0.86 ± 0.01	0.87 ± 0.01

Table 2: Benchmark performance for GASS versus alternative non-conjugate elliptical sampling approaches. Results are averages over 100 independent trials \pm standard error.

- **Posterior projections (PP).** No constraints are imposed on the model during posterior inference. Instead, we use the posterior projection approach of Lin and Dunson [34] to post-hoc enforce the constraints.
- **Logistic posterior projections (LPP).** A hybrid of the previous two approaches combined: modeling the logits for $[0, 1]$ constraints but projecting the posterior samples to the monotone surface.

All models use the true prior mean and covariance; for the logistic models, we empirically estimate the covariance of the logit-transformed θ . We compare performance with $2m$ MCMC steps, where the first m are a burn-in phase and the last m are used for posterior approximation; we consider $m = [100, 500, 1000, 5000, 10000]$. Performance is measured in terms of mean squared error (MSE) and coverage rate of the 90% credible intervals for every θ_i point. Results are averaged over 100 independent trials with (θ, \mathbf{y}) resampled from eq. (15) at the start of each trial.

Figure 5 shows examples of the fits for each method, with 90% credible intervals. GASS is the only procedure that results in good coverage of the true mean and captures the shape of the overall curve after 20K MCMC steps; the other methods tend to over-smooth the curve and underestimates the uncertainty.

Table 2 shows the aggregate results of the benchmarks. GASS outperforms all four comparison methods in terms of both error and coverage. After $m = 100$ samples the MSE for GASS is lower and coverage is higher than any of the other strategies after $m = 10000$ samples. Further, the model appears to have almost fully mixed after 5000 samples, with the coverage rate close to 90%.

d=2								
	Observations			True Rate		Coverage		
Model	MAE	RMSE	NLL	MAE	RMSE	50%	75%	95%
NMF	1.70	2.47	364.29	1.12	1.72	N/A	N/A	N/A
PGDS(0.25)	1.76	2.47	360.76	1.29	2.00	15.32	25.35	40.59
PGDS(0.5)	1.74	2.43	359.96	1.31	2.03	14.43	23.96	38.94
PGDS(1)	1.74	2.43	360.29	1.30	2.02	14.58	23.93	39.09
NBinom BTF	$> 10^2$	$> 10^3$	$> 10^5$	48.10	$> 10^2$	30.17	47.80	69.67
Poisson BTF	1.66	2.34	354.52	0.91	1.43	34.16	53.40	74.72

d=3								
	Observations			True Rate		Coverage		
Model	MAE	RMSE	NLL	MAE	RMSE	50%	75%	95%
NMF	2.26	3.48	552.10	1.15	1.88	N/A	N/A	N/A
PGDS(0.25)	1.72	2.42	365.26	1.06	1.68	21.58	35.83	56.24
PGDS(0.5)	1.71	2.42	364.39	1.06	1.68	21.51	35.58	55.69
PGDS(1)	1.71	2.41	363.86	1.05	1.68	21.27	35.00	54.81
NBinom BTF	$> 10^3$	$> 10^4$	$> 10^6$	$> 10^2$	$> 10^3$	37.14	59.20	82.60
Poisson BTF	1.78	2.47	358.95	0.79	1.29	40.45	62.11	82.87

d=5								
	Observations			True Rate		Coverage		
Model	MAE	RMSE	NLL	MAE	RMSE	50%	75%	95%
NMF	2.65	4.21	596.37	1.43	2.31	N/A	N/A	N/A
PGDS(0.25)	1.74	2.44	368.10	0.86	1.37	31.46	50.20	73.55
PGDS(0.5)	1.76	2.43	369.59	0.87	1.37	31.35	50.46	73.83
PGDS(1)	1.80	2.56	371.89	0.88	1.40	29.80	49.05	71.91
NBinom BTF	$> 10^2$	$> 10^2$	$> 10^4$	10.77	94.94	45.31	67.85	89.66
Poisson BTF	1.71	2.26	355.75	0.81	1.25	41.64	63.55	83.98

d=10								
	Observations			True Rate		Coverage		
Model	MAE	RMSE	NLL	MAE	RMSE	50%	75%	95%
NMF	4.53	7.45	$> 10^3$	1.89	3.12	N/A	N/A	N/A
PGDS(0.25)	1.86	2.74	391.60	0.82	1.33	38.25	60.21	83.84
PGDS(0.5)	1.89	2.80	402.63	0.79	1.29	39.20	61.48	84.40
PGDS(1)	1.79	2.52	376.62	0.79	1.26	39.21	61.52	84.61
NBinom BTF	31.87	42.16	$> 10^3$	3.13	11.14	43.33	66.17	87.55
Poisson BTF	1.72	2.43	359.12	0.83	1.32	44.05	67.03	86.59

Table 3: Mean results on the Poisson dynamical system benchmark; smaller is better for all metrics. NMF: nonnegative matrix factorization; PGDS(τ): Poisson-gamma dynamical system with hyperparameter τ ; NBinom-BTF: Bayesian tensor filtering with conditionally-conjugate negative binomial likelihood; Poisson-BTF: Bayesian tensor filtering with constrained, non-conjugate Poisson likelihood via GASS inference; NLL: negative log-likelihood; MAE: mean absolute deviation from truth; RMSE: root mean-squared error from truth. Bold indicates the best performance at each embedding dimension setting.

6.3 Non-stationary Poisson dynamical systems

We benchmark BTF on a synthetic Poisson tensor dataset where the observations are Poisson distributed with a latent rate curve for each function. The rate at every point in the curve is the inner product of two gamma random vectors,

$$\begin{aligned} h_{j\ell} &\sim \text{Bern}(0.2), & u_{j\ell d} &\sim (1 - h_{j\ell})\delta_0 + h_{j\ell}\text{Ga}(1, 1), & v_{jtd} &= \sum_{\ell=1}^t u_{j\ell d}, \\ w_{id} &\sim \text{Ga}(1, 1), & y_{ijt} &\sim \text{Pois}(\langle w_i, v_{jt} \rangle). \end{aligned}$$

The resulting true rates form a monotonic curve of constant plateaus with occasional jumps. As in the dose-response data, the columns evolve independently of each other, rather than through a common time parameter. We set the latent factor dimension to 3.

We compare a Poisson likelihood version of BTF with GASS inference (Poisson BTF) to nonnegative matrix factorization (NMF), the Poisson-Gamma dynamical system (PGDS) model of Schein et al. [50], and a negative binomial likelihood version of BTF (NBinom BTF) with Pólya-Gamma augmentation [48]. For PGDS, we contacted the authors who suggested we try three different values of the hyperparameter $\tau = (0.25, 0.5, 1)$. For NBinom BTF, we use MCMC-within-Gibbs and sample the latent rate parameter with 30 steps of random walk Metropolis-Hastings for every Gibbs step. For Poisson BTF, we initialize with NMF and did not use any conditioning heuristic.

We run all models for 5000 burn-in iterations and collect 5000 samples on an $11 \times 12 \times 20$ tensor with the upper left $3 \times 3 \times 20$ corner held out. We conduct 5 independent trials, regenerating new data each time and evaluating the models on the held out data and true latent rate. We measure performance in three categories of metrics: (i) mean absolute error (MAE), root mean squared error (RMSE), and negative log-likelihood (NLL) on held out observations; (ii) MAE and RMSE on the true latent rate; and (iii) posterior credible interval coverage of the true rate at 50%, 75%, and 90% targets. We evaluate all models at embedding dimensions $d = (2, 3, 5, 10)$ to compare sensitivity to the common hyperparameter.

Table 3 presents the results. The Poisson BTF model performs similarly across the range of dimension embeddings, whereas the other models are more sensitive. The NMF model generally performs better with a smaller embedding dimension while the PGDS model performs better with larger dimensions. In the case of NMF, this is due to overfitting without any smoothness prior built in. PGDS uses a canonical tensor decomposition, where the third tensor dimension is modeled with an embedding that is shared between rows and columns. This requires an inflation of the embedding dimension when columns are all evolving independently in order to sufficiently capture all of the latent structure in the data.

The negative binomial BTF model performs poorly on held out observations and true rate estimation. This is due to the instability of the model on held out data. The Pólya-Gamma approach to negative binomial likelihoods uses the

Cancer Drug Studies		
	Pilot Study	Landscape Study
Model	NLL	NLL
NMF	262.75 ± 308.12	25573.14
LMF	589.17 ± 582.29	$> 10^6$
BTF	-80.22 ± 9.67	-3268.11

Table 4: Left: mean results \pm standard error on held out data for the pilot cancer drug studies. Right: Results on a single test set of 1000 curves for the landscape study. NMF: nonnegative matrix factorization; LFM: logistic factor model; BTF: Bayesian tensor filtering; NLL: negative log-likelihood.

$NB(r, \sigma(\beta))$ parameterization, where r is an unknown dispersion parameter, σ is the logistic function, and β are log-odds. The mean of the distribution is

$$\frac{n\sigma(\beta)}{1 - \sigma(\beta)}.$$

When the failure rate $\sigma(\beta)$ is near zero or one, small changes in β lead to large changes in the mean. Thus, small errors in $\beta_{ijt} = w_i^\top v_{jt}$ on the held out data lead to large errors in the observational NLL and the MAE and RMSE metrics we consider.

For each embedding dimension, the Poisson BTF model performs competitively or better than the other models in each category. The negative log-likelihood on held out observations is always lowest for the Poisson BTF model. For other metrics, Poisson BTF is either the best performing model or within 10% of the best performing model. By contrast, the other models are off by more substantial amounts in certain categories at certain dimension embedding settings. The credible intervals for Poisson BTF also consistently have better coverage than PGDS for the same embedding dimension. This includes the true embedding dimension $d = 3$, where Poisson BTF shows the best performance in terms of MAE to the true rate, and competitive performance with other choices of d . By contrast, both NMF and PGDS perform substantially worse in the $d = 3$ regime than in other choices. Finally, unlike the negative binomial model, the Poisson BTF model maintains a stable prediction on held out entries.

6.4 Cancer drug study

We evaluate the proposed empirical Bayes dose-response model, built on top of BTF, on two cancer drug studies. First, we use a small internal pilot study conducted at Columbia University Medical Center. The pilot study tested 35 drugs against 28 tumor organoids, each at 9 different concentrations with 6 replicates. Second, we run on a large-scale, “landscape” study [31] that tested 67 drugs against 284 tumor organoids, each at 7 different concentrations with 2 replicates. For the pilot study, we run 5 independent trials, holding out 30 curves

at random, subject to the constraint that no column or row is left without any observations in the training set. For the landscape study, we hold out a single test set of 1000 curves ($\approx 5\%$ of the total entries). Since this is real data, MAE and RMSE from the truth are not available; we measure performance solely in terms of negative log-likelihood on the held out data.

The standard dose-response modeling approach in cancer datasets is a log-linear logistic model [58]. For a baseline, we extend that model to a logistic factor model (LFM), using the same preprocessing strategy. We also compare to NMF as a second baseline. To ensure the monotonicity, we project the NMF results to be monotone curves using the PAV algorithm as in Lin and Dunson [34]. We choose the factor size in both models by 5-fold cross-validation on the training set.

For BTF, we perform a grid search over hyperparameters: $\rho^2 = \{0.001, 0.01, 0.1\}$, factor size $D = \{1, 3, 5, 8\}$, and the order of the trend filtering matrix $k = \{0, 1\}$; we select the best model using the deviance information criterion [9]. We evaluate the BTF model using the average of the posterior draws, rather than the full Bayes estimate; this enables us to fairly compare with the NMF and LFM point estimates of the latent mean. We run 10000 Gibbs sampling steps in both studies, discarding the first 5000 as burn-in.

Table 4 present the results. The BTF dose-response model outperforms both baselines in terms of negative log-likelihood of the held out data in both studies. Results in terms of RMSE and MAE (not shown) on the raw observations were similar for all three models in both studies (e.g. RMSE 0.14 ± 0.01 , MAE 0.20 ± 0.01 in the pilot study). The BTF procedure is also more stable in the pilot study cross-validation, with a much lower reconstruction variance than either baseline. This suggests BTF not only forms a more accurate basis for a dose-response model, but is also more reliable.

Qualitative results on the held out predictions are in Figure 1 (orange). All 9 plots are for real data from the landscape study, with the gray observations held out. The orange line shows the posterior mean of the predicted curves. The curves have all of the desired properties: monotonicity with dose, bounded between zero and one, mostly smooth, locally adaptive to sharp jumps in the data, and highly predictive of the outcomes of the experiments. The orange bands show the 50% approximate posterior credible intervals using the empirical Bayes likelihood model. The credible intervals are conservative, estimating a larger variance than is actually observed in the outcomes. Even still, the NMF and LMF models far exceed these bands in certain points in the curve. This is due to the heteroskedastic nature of the likelihood and the misspecification of the NMF and LMF loss functions. Both competing models optimize for squared error, effectively making a heteroskedastic assumption on the model. In RMSE terms, all three models perform nearly identically, within ± 0.01 of each other on both datasets. Judging the models by RMSE would be misleading, since the high degree of noise in the first three dose levels dominates the overall loss and obscures the real fit of the model.

7 Landscape study analysis

In many cancer drug studies, molecular features such as gene expression, genomic mutations, or copy number alterations are gathered. These features represent useful side information that can help further denoise the dose-response data. Features also enable one to address the “cold-start” problem, enabling predictions for samples that have no dose-response data available. Features that are predictive of sensitivity or resistance in the dose-response experiments may represent “biomarkers”—diagnostic indicators of drug response. Biomarkers are candidate targets for future experimental investigation such as targeted drug development. The landscape study of Lee et al. [31] considered 115 molecular features, all binary, with a subset of the features gathered on a subset of organoid cell lines.

To incorporate potentially-missing features into the BTF dose-response model, we take a multi-view factorization approach. For each feature x_{im} , $m = 1, \dots, M$, for organoid cell line i , we assume a latent factor model between the cell line embeddings and feature embeddings,

$$\begin{aligned} x_{im} \mid w_i, u_m &\sim \text{Bern}(w_i^\top u_m) \mathbf{1}[0 \leq w_i^\top u_m \leq 1] \\ u_m &\sim \text{MVN}(\mathbf{0}, \sigma_u^2 I). \end{aligned} \quad (16)$$

Rather than a logistic link function with an unconstrained likelihood model, we use an identity link with a $[0, 1]$ constraint. A logistic link would put the feature likelihood and dose-response likelihoods on different scales. By using a constrained identity link function in eq. (16), a change in w_i has the same effect on the probability of dose-specific survival as in the probability of a feature being positive.

For the landscape study, we set $\sigma_u = 1$ and use a latent embedding dimension of 10. All other hyperparameters and settings are the same as in the benchmarks in section 6.4.

7.1 Site of origin structure captured by organoid embeddings

We first check whether the model has uncovered latent structure in the data. Tumor molecular profiles and drug response are strongly associated with the site of origin of the tumor. The BTF dose-response model was not supplied any direct site-of-origin information. Nonetheless, we expect that some clustering of organoids into site of origin should emerge.

Visualizing structure in the 10-dimensional organoid embeddings is challenging, as with any data of more than 4 dimensions. Figure 6 shows a 2-dimensional principal components analysis (PCA) of the posterior mean of the 10-dimensional organoid embeddings. The 2-d projection of the embeddings confirms that site of origin is associated with the first two principal components. Ovarian cancers are predominantly clustered in the bottom-center of the 2-d space, liver cancers skew left, and breast and gastric cancers skew right. Brain cancers, cover the entire space but also represent the largest and most diverse set of original tumors in the landscape dataset.

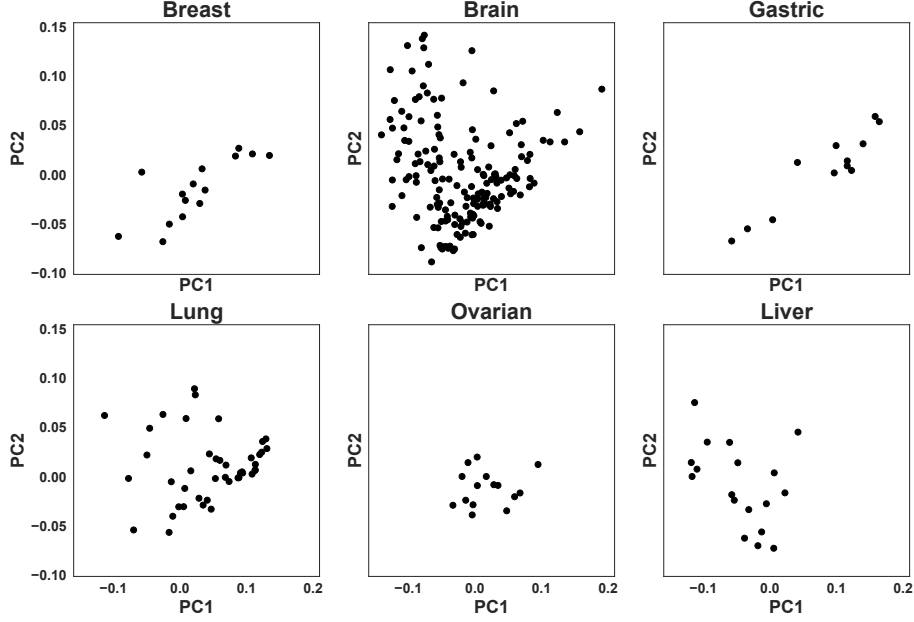


Figure 6: Two-dimensional PCA projection of the learned organoid embeddings, stratified by tumor site of origin. The embeddings reveal learned structure relating to site of origin, which was not an explicit input to the model.

7.2 Uncovering biomarkers associated with drug sensitivity

To discover potential biomarkers, we correlate the area under the dose-response curve (AUC) with the feature probability. To calculate the AUC, we approximate the curve with a piecewise-linear fit between successive dose intervals,

$$\text{AUC} = \frac{1}{T-1} \sum_{t=1}^{T-1} w_i^\top (v_{jt} - v_{j(t+1)}). \quad (17)$$

We fit independent linear models to predict the posterior mean AUC values as a function of the biomarker probability. Features are ranked by r^2 value and stratified into *sensitivity* or *resistance* based on the directionality of their slope. We filter out spurious results by removing features and drugs with standard deviation between samples is less than 0.05.

Figure 7 shows the top two results. Both features are flagged as potential biomarkers of drug sensitivity, indicating that as the probability of the feature increases, the AUC decreases. The top result (fig. 7a) associates the vIII rearrangement of the epidermal growth factor receptor (EGFR) gene with sensitivity to treatment with the drug Trametinib. EGFR is involved in many proliferation-inducing signaling pathways, including the important RAS/REF/MEK/ERK

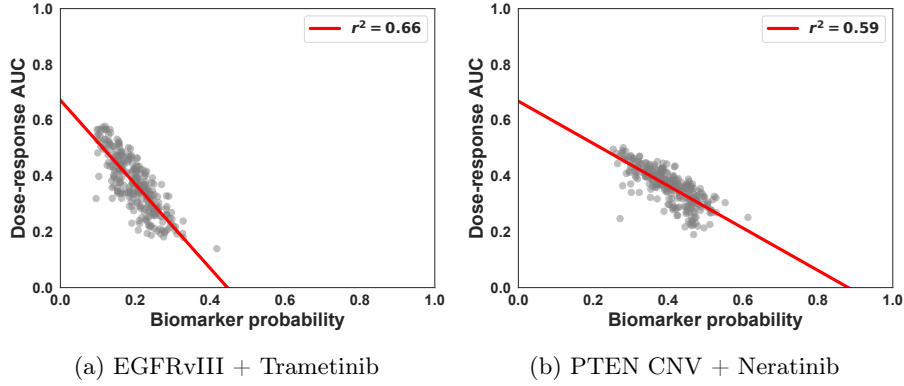


Figure 7: Top two biomarker results correlating drug sensitivity with biomarker presence. Left: The vIII rearrangement in the EGFR gene is associated with sensitivity to treatment with the drug Trametinib. Right: Copy number variation, typically in the form of a copy number loss, in the tumor suppressor gene PTEN is associated with sensitivity to treatment with the drug Neratinib.

pathway [25]. The vIII rearrangement of EGFR (EGFRvIII) leads to continual activation of EGFR, promoting oncogenesis [19]. Current approaches to targeting EGFRvIII have not seen major success [2]. The top result flags organoids with the EGFRvIII biomarker as associated with sensitivity to Trametinib, a MEK1/2 inhibitor. This may be due to Trametinib silencing the RAS/RAF/MEK/ERK pathway being constitutively activated by the EGFRvIII rearrangement, suggesting a subpopulation of patients with EGFRvIII would benefit from Trametinib.

The second top result (fig. 7b) is copy number alterations in the phosphatase and tensin homolog (PTEN) gene being associated with sensitivity to treatment with Neratinib. PTEN is a tumor suppressor gene frequently lost in cancer, in particular glioblastoma [26] (GBM). The landscape study includes a plurality of GBM organoids and correspondingly the vast majority of the PTEN copy number variations in the dataset are due to PTEN loss. PTEN loss predicts resistance in Trastuzumab [41] and Neratinib has been shown to overcome Trastuzumab resistance in a subset of breast cancer patients [7]. The association between PTEN loss and Neratinib sensitivity suggests a possible mechanistic explanation for these effects, as well as a biomarker for other treatment-resistant patients.

8 Discussion

Multi-sample, multi-drug cancer studies are time and resource intensive. The outcomes from these studies are noisy, often incomplete, observations of biological responses to candidate therapies. Denoising observations and imputing missing experiments is an important step in the scientific analysis and drug discovery pipeline. The Bayesian tensor filtering model we presented in this paper

enables scientists to flexibly model dose-response curves with consideration for measurement error and biological constraints on the shape of the curve. While the BTF model is an improvement over the state of the art, we believe there are several improvements that could be made.

The BTF model assumes a regularly-spaced grid. Extensions to irregular grids of dose levels may be possible via a Bayesian group trend filtering extension to an irregular grid approach for scalar trend filtering. The approach taken in Faulkner and Minin [12] was based on integrated Wiener processes. In the case of horseshoe priors, the method requires an approximation [35] that was only tractable up to second-order differences. We have found higher-order priors to not be necessary in our experience, suggesting irregular grid extensions here hold promise. Other methods from the trend filtering literature, such as the falling factorial basis [59], may also be adaptable to the Bayesian group trend filtering case.

The AUC values computed in section 7.2 are likely to be slightly biased due to the piecewise-linear approximation of the true curve. For more precise estimation of the AUC curves, one could adopt the Nadaraya–Watson kernel smoothing of Piegorsch et al. [46] by introducing pseudo-concentrations for each posterior sample. An alternative approach would be to introduce the pseudo-concentrations as missing data and have the model impute the values directly during posterior sampling with the composite trend filtering enforcing nonstationary smoothness. One could think of the latter approach as the more Bayesian approach whereas the post-hoc smoothing will be computationally faster. The modeling choice is analogous to how we enforce monotonicity directly in the posterior as opposed to post-hoc merging via the PAV approach of Lin and Dunson [34].

The current BTF model is computationally intensive. For small scale studies like the pilot study, the model runs in a few hours on a laptop. The landscape study required several days on a compute cluster to perform the hyperparameter search. Relative to the years required for the landscape experiments, the run time is negligible. Nevertheless, offering an alternative inference approach that can scale more efficiently, such as variational inference, may make the BTF model useful for a broader group of scientists.

Finally, BTF does not support adding chemical features about the drugs. In organoid studies, scientists generally know the class of approved and potentially-translatable drugs. Even in high-throughput screening studies for cancer cell lines (i.e. not organoids), only known drugs—mostly chemotherapy agents—are typically tested. The goal in these studies is to find biological markers of resistance or sensitivity to the set of available compounds. The molecular feature analysis in section 7 addresses this. However, if we were trying to discover entirely new drugs, such as might be done at a pharmaceutical company, extending the model to include drug features would be useful. This may be possible by including a second side-information matrix and adding a drug-specific embedding in a manner similar to canonical tensor decomposition.

References

- [1] Farnoosh Abbas-Aghababazadeh, Pengcheng Lu, and Brooke L Fridley. Nonlinear mixed-effects models for modeling *in vitro* drug response data to determine problematic cancer cell lines. *Scientific Reports*, 9(1):1–9, 2019.
- [2] Zhenyi An, Ozlem Aksoy, Tina Zheng, Qi-Wen Fan, and William A Weiss. Epidermal growth factor receptor and EGFRvIII in glioblastoma: Signaling pathways and targeted therapies. *Oncogene*, 37(12):1561–1575, 2018.
- [3] JM Bernardo, MJ Bayarri, JO Berger, AP Dawid, D Heckerman, A Smith, and M West. Bayesian factor regression models in the “large p, small n” paradigm. *Bayesian Statistics*, 7:733–742, 2003.
- [4] Anindya Bhadra, Jyotishka Datta, Nicholas G. Polson, and Brandon Willard. The horseshoe+ estimator of ultra-sparse signals. *Bayesian Analysis*, 12(4):1105–1131, 2017.
- [5] Christopher M Bishop. *Pattern recognition and machine learning*. Springer, 2006.
- [6] Björn Bornkamp and Katja Ickstadt. Bayesian nonparametric estimation of continuous monotone functions with applications to dose–response analysis. *Biometrics*, 65(1):198–205, 2009.
- [7] Alexandra Canonici, Merel Gijsen, Maeve Mullooly, Ruth Bennett, Noujoude Bouguern, Kasper Pedersen, Neil A O’Brien, Ioannis Roxanis, Ji-Liang Li, Esther Bridge, et al. Neratinib overcomes trastuzumab resistance in HER2 amplified breast cancer. *Oncotarget*, 4(10):1592, 2013.
- [8] Carlos M. Carvalho, Nicholas G. Polson, and James G. Scott. The horseshoe estimator for sparse signals. *Biometrika*, 97(2):465–480, 2010.
- [9] Gilles Celeux, Florence Forbes, Christian P. Robert, and D. Michael Titterton. Deviance information criteria for missing data models. *Bayesian Analysis*, 1(4):651–673, 2006.
- [10] Jarno Drost and Hans Clevers. Organoids in cancer research. *Nature Reviews Cancer*, 18(7):407–418, 2018.
- [11] Francois Fagan, Jalaj Bhandari, and John Cunningham. Elliptical slice sampling with expectation propagation. In *Uncertainty in Artificial Intelligence*, 2016.
- [12] James R. Faulkner and Vladimir N. Minin. Locally adaptive smoothing with Markov random fields and shrinkage priors. *Bayesian Analysis*, 13(1):225, 2018.
- [13] Emily B Fox and David B Dunson. Bayesian nonparametric covariance regression. *Journal of Machine Learning Research*, 16(1):2501–2542, 2015.
- [14] Brooke L Fridley, Gregory Jenkins, Daniel J Schaid, and Liewei Wang. A Bayesian hierarchical nonlinear model for assessing the association between genetic variation and drug cytotoxicity. *Statistics in Medicine*, 28(21):2709–2722, 2009.

- [15] Mathew J. Garnett, Elena J. Edelman, Sonja J. Heidorn, Chris D. Greenman, Anahita Dastur, King Wai Lau, Patricia Greninger, I. Richard Thompson, Xi Luo, Jorge Soares, et al. Systematic identification of genomic markers of drug sensitivity in cancer cells. *Nature*, 483(7391):570, 2012.
- [16] Laetitia Gauvin, André Panisson, and Ciro Cattuto. Detecting the community structure and activity patterns of temporal networks: A non-negative tensor factorization approach. *PLoS One*, 9(1), 2014.
- [17] Mahmoud Ghandi, Franklin W Huang, Judit Jané-Valbuena, Gregory V Kryukov, Christopher C Lo, E Robert McDonald, Jordi Barretina, Ellen T Gelfand, Craig M Bielski, Haoxin Li, et al. Next-generation characterization of the cancer cell line encyclopedia. *Nature*, 569(7757):503–508, 2019.
- [18] Aklilu Habteab Ghebretinsae, Christel Faes, Geert Molenberghs, Marlies De Boeck, and Helena Geys. A Bayesian, generalized frailty model for comet assays. *Journal of Biopharmaceutical Statistics*, 23(3):618–636, 2013.
- [19] Gao Guo, Ke Gong, Bryan Wohlfeld, Kimmo J Hatanpaa, Dawen Zhao, and Aamyn A Habib. Ligand-independent EGFR signaling. *Cancer Research*, 75(17):3436–3441, 2015.
- [20] P. Richard Hahn, Jingyu He, and Hedibert Lopes. Bayesian factor model shrinkage for linear IV regression with many instruments. *Journal of Business & Economic Statistics*, 36(2):278–287, 2018.
- [21] Creighton Heaukulani and Mark van der Wilk. Scalable Bayesian dynamic covariance modeling with variational Wishart and inverse Wishart processes. In *Advances in Neural Information Processing Systems*, pages 4582–4592, 2019.
- [22] Lei Huang, Shengnan Wu, and Da Xing. High fluence low-power laser irradiation induces apoptosis via inactivation of Akt/GSK3 β signaling pathway. *Journal of Cellular Physiology*, 226(3):588–601, 2011.
- [23] W Evan Johnson, Cheng Li, and Ariel Rabinovic. Adjusting batch effects in microarray expression data using empirical Bayes methods. *Biostatistics*, 8(1):118–127, 2007.
- [24] Seung-Jean Kim, Kwangmoo Koh, Stephen Boyd, and Dmitry Gorinevsky. ℓ_1 trend filtering. *SIAM review*, 51(2):339–360, 2009.
- [25] Walter Kolch, Melinda Halasz, Marina Granovskaya, and Boris N Kholodenko. The dynamic control of signal transduction networks in cancer cells. *Nature Reviews Cancer*, 15(9):515–527, 2015.
- [26] Dimpy Koul. PTEN signaling pathways in glioblastoma. *Cancer Biology & Therapy*, 7(9):1321–1325, 2008.
- [27] D. R. Kowal, D. S. Matteson, , and D. Ruppert. Dynamic shrinkage processes. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 2019.
- [28] Tsuyoshi Kuniyama, Carolyn T Halpern, and Amy H Herring. Non-parametric Bayes models for mixed scale longitudinal surveys. *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, 68(4):1091–1109, 2019.

- [29] Minjung Kyung, Jeff Gill, Malay Ghosh, George Casella, et al. Penalized regression, standard errors, and Bayesian lassos. *Bayesian Analysis*, 5(2):369–411, 2010.
- [30] Alexander Lachmann, Federico M Giorgi, Mariano J Alvarez, and Andrea Califano. Detection and removal of spatial bias in multiwell assays. *Bioinformatics*, 32(13):1959–1965, 2016.
- [31] Jin-Ku Lee, Zhaoqi Liu, Jason K. Sa, Sang Shin, Jiguang Wang, Mykola Boryd, Hee Jin Cho, Oliver Elliott, Timothy Chu, Seung Won Choi, Daniel I. S. Rosenbloom, In-Hee Lee, Yong Jae Shin, Hyun Ju Kang, Donggeon Kim, Sun Young Kim, Moon-Hee Sim, Jusun Kim, Taehyang Lee, Yun Jee Seo, Hyemi Shin, Mi-jeong Lee, Sung Heon Kim, Yong-Jun Kwon, Jeong-Woo Oh, Minsuk Song, Misuk Kim, Doo-Sik Kong, Jung Won Choi, Ho Jun Seol, Jung-Il Lee, Seung Tae Kim, Joon Oh Park, Kyoung-Mee Kim, Sang-Yong Song, Jeong-Won Lee, Hee-Cheol Kim, Jeong Eon Lee, Min Gew Choi, Sung Wook Seo, Young Mog Shim, Jae Ill Zo, Byong Chang Jeong, Yeup Yoon, Gyu Ha Ryu, Nayoung K. D. Kim, Joon Seol Bae, Woong-Yang Park, Jeongwu Lee, Roel G. W. Verhaak, Antonio Iavarone, Jeeyun Lee, Raul Rabadan, and Do-Hyun Nam. Pharmacogenomic landscape of patient-derived tumor cells informs precision oncology therapy. *Nature Genetics*, 50(10):1399–1411, 2018.
- [32] Jeffrey T Leek, Robert B Scharpf, Héctor Corrada Bravo, David Simcha, Benjamin Langmead, W Evan Johnson, Donald Geman, Keith Baggerly, and Rafael A Irizarry. Tackling the widespread and critical impact of batch effects in high-throughput data. *Nature Reviews Genetics*, 11(10):733–739, 2010.
- [33] Lingge Li, Dustin Pluta, Babak Shahbaba, Norbert Fortin, Hernando Ombao, and Pierre Baldi. Modeling dynamic functional connectivity with latent factor Gaussian processes. In *Advances in Neural Information Processing Systems*, pages 8263–8273, 2019.
- [34] Lizhen Lin and David B Dunson. Bayesian monotone regression using Gaussian process projection. *Biometrika*, 101(2):303–317, 2014.
- [35] Finn Lindgren and Håvard Rue. On the second-order random walk model for irregular locations. *Scandinavian Journal of Statistics*, 35(4):691–700, 2008.
- [36] Cecile Low-Kam, Donatello Telesca, Zhaoxia Ji, Haiyuan Zhang, Tian Xia, Jeffrey I Zink, and Andre E Nel. A Bayesian regression tree approach to identify the effect of nanoparticles’ properties on toxicity profiles. *Annals of Applied Statistics*, 9(1):383–401, 2015.
- [37] Enes Makalic and Daniel F. Schmidt. A simple sampler for the horseshoe estimator. *IEEE Signal Processing Letters*, 23(1):179–182, 2015.
- [38] Bogdan Mazouze, Robert Nadon, and Vladimir Makarenkov. Identification and correction of spatial bias are essential for obtaining quality data in high-throughput screening technologies. *Scientific Reports*, 7(1):11921, 2017.
- [39] Thomas P Minka. Expectation propagation for approximate Bayesian inference. In *Uncertainty in Artificial Intelligence*, pages 362–369, 2001.
- [40] Iain Murray, Ryan Adams, and David MacKay. Elliptical slice sampling. In *Artificial Intelligence and Statistics*, 2010.

- [41] Yoichi Nagata, Keng-Hsueh Lan, Xiaoyan Zhou, Ming Tan, Francisco J Esteva, Aysegul A Sahin, Kristine S Klos, Ping Li, Brett P Monia, Nina T Nguyen, et al. PTEN activation contributes to tumor inhibition by trastuzumab, and loss of PTEN predicts trastuzumab resistance in patients. *Cancer Cell*, 6(2):117–127, 2004.
- [42] Radford M Neal. Slice sampling. *Annals of Statistics*, 31(3):705–767, 2003.
- [43] Brian Neelon and David B Dunson. Bayesian isotonic regression and trend analysis. *Biometrics*, 60(2):398–406, 2004.
- [44] Trina Patel, Donatello Telesca, Saji George, and André E Nel. Toxicity profiling of engineered nanomaterials via multivariate dose–response surface modeling. *Annals of Applied statistics*, 6(4):1707, 2012.
- [45] Francois Perron and Kerrie Mengersen. Bayesian nonparametric modeling using mixtures of triangular distributions. *Biometrics*, 57(2):518–528, 2001.
- [46] Walter W Piegorsch, Hui Xiong, Rabi N Bhattacharya, and Lizhen Lin. Non-parametric estimation of benchmark doses in environmental risk assessment. *Environmetrics*, 23(8):717–728, 2012.
- [47] Nicholas G. Polson and James G. Scott. Shrink globally, act locally: Sparse Bayesian regularization and prediction. *Bayesian Statistics*, 9:501–538, 2010.
- [48] Nicholas G. Polson, James G. Scott, and Jesse Windle. Bayesian inference for logistic models using Pólya–Gamma latent variables. *Journal of the American Statistical Association*, 108(504):1339–1349, 2013.
- [49] G. O. Roberts and J. S. Rosenthal. General state space Markov chains and MCMC algorithms. *Probability surveys*, 1:20–71, 2004.
- [50] Aaron Schein, Hanna Wallach, and Mingyuan Zhou. Poisson-Gamma dynamical systems. In *Advances in Neural Information Processing Systems*, 2016.
- [51] Thomas S Shively, Thomas W Sager, and Stephen G Walker. A Bayesian approach to non-parametric monotone function estimation. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 71(1):159–175, 2009.
- [52] Robert H Shoemaker. The NCI60 human tumour cell line anticancer drug screen. *Nature Reviews Cancer*, 6(10):813–823, 2006.
- [53] Stephan Spiegel, Jan Clausen, Sahin Albayrak, and Jérôme Kunegis. Link prediction on evolving data using tensor factorization. In *Pacific-Asia Conference on Knowledge Discovery and Data Mining*, 2011.
- [54] Koh Takeuchi, Hisashi Kashima, and Naonori Ueda. Autoregressive tensor factorization for spatio-temporal predictions. In *International Conference on Data Mining*, 2017.
- [55] Wesley Tansey, Kathy Li, Haoran Zhang, Scott W Linderman, Raul Rabadan, David M Blei, and Chris H Wiggins. Dose–response modeling in high-throughput cancer drug screenings: An end-to-end approach. *arXiv preprint arXiv:1812.05691*, 2018.

- [56] Ryan J. Tibshirani. Adaptive piecewise polynomial estimation via trend filtering. *Annals of Statistics*, 42(1):285–323, 2014.
- [57] Stéphanie L. Van Der Pas, Bas J.K. Kleijn, and Aad W. Van Der Vaart. The horse-shoe estimator: Posterior concentration around nearly black vectors. *Electronic Journal of Statistics*, 8(2):2585–2618, 2014.
- [58] Daniel J. Vis, Lorenzo Bombardelli, Howard Lightfoot, Francesco Iorio, Mathew J. Garnett, and Lodewyk FA Wessels. Multilevel models improve precision and speed of IC50 estimates. *Pharmacogenomics*, 17(7):691–700, 2016.
- [59] Yu-Xiang Wang, Alex Smola, and Ryan Tibshirani. The falling factorial basis and its statistical applications. In *International Conference on Machine Learning*, pages 730–738, 2014.
- [60] John N. Weinstein, Eric A. Collisson, Gordon B. Mills, Kenna R. Mills Shaw, Brad A. Ozenberger, Kyle Ellrott, Ilya Shmulevich, Chris Sander, Joshua M. Stuart, and Cancer Genome Atlas Research Network. The cancer genome atlas pan-cancer analysis project. *Nature Genetics*, 45(10):1113, 2013.
- [61] Matthew W Wheeler. Bayesian additive adaptive basis tensor product models for modeling high dimensional surfaces: an application to high-throughput toxicity testing. *Biometrics*, 75(1):193–201, 2019.
- [62] Ander Wilson, David M Reif, and Brian J Reich. Hierarchical dose–response modeling for high-throughput toxicity screening of environmental chemicals. *Biometrics*, 70(1):237–246, 2014.
- [63] Liang Xiong, Xi Chen, Tzu-Kuo Huang, Jeff Schneider, and Jaime G. Carbonell. Temporal collaborative filtering with Bayesian probabilistic tensor factorization. In *International Conference on Data Mining*, 2010.
- [64] Aonan Zhang and John Paisley. Deep Bayesian nonparametric tracking. In *International Conference on Machine Learning*, pages 5828–5836, 2018.

A Local shrinkage updates

The local shrinkage parameters $\tau_{j\ell}$ can be updated through a double latent variable augmentation trick,

$$\begin{aligned}
(\tau_{j\ell} \mid -) &\sim \text{InvGamma}(D + 1, \left\| \Delta^{(k)} V_j \right\|_2^2 / 2 + 1/c_{j\ell}) \\
(c_{j\ell} \mid -) &\sim \text{InvGamma}(1, 1/\tau_{j\ell}^2 + 1/\phi_{j\ell}) \\
(\phi_{j\ell} \mid -) &\sim \text{InvGamma}(1, 1/c_{j\ell} + 1/\eta_{j\ell}) \\
(\eta_{j\ell} \mid -) &\sim \text{InvGamma}(1, 1/\phi_{j\ell} + 1).
\end{aligned} \tag{18}$$

The updates in eq. (18) come from the HS+ prior being a two-level horseshoe prior. The inverse-gamma latent variable augmentation for the horseshoe is fast and typically mixes quickly [37].

B Gaussian likelihood

When the likelihood is normal, $y_{ijt} \sim \mathcal{N}(w_i^\top v_{jt}, \nu^2)$, where ν^2 is a nuisance parameter, the factor and loading updates are conjugate. Let $\tilde{V} = (v_{1,1}, v_{1,2}, \dots, v_{1,T}, v_{2,1}, \dots, v_{M,T})$, and $\Omega^{-1} = \text{diag}\{1/\nu^2\}$, then the updates are multivariate normal,

$$\begin{aligned}
Q^{(i)} &= (\tilde{V}^\top \Omega^{-1} \tilde{V} + \text{diag}(\sigma^{-2}))^{-1} \\
(w_i \mid -) &\sim \text{MVN} \left(Q^{(i)} \tilde{V}^\top \Omega^{-1} \text{vec}(Y_i^\top), Q^{(i)} \right) \\
\mathcal{T}^{(j)} &= \text{diag}(1/(\rho^2 \tau_j^2)) \\
\Sigma^{(j)} &= (I_D \otimes \Delta^\top \mathcal{T} \Delta) + (W \otimes I_T)^\top \Omega^{-1} (W \otimes I_T) \\
(\text{vec}(V_j) \mid -) &\sim \text{MVN}(\Sigma^{(j)} (W \otimes I_T) \Omega^{-1} \text{vec}(Y_j^\top), \Sigma^{(j)}),
\end{aligned} \tag{19}$$

where **diag** diagonalizes the given vector, **vec** is the vectorization operator, and \otimes is the Kronecker product. In both the w_i and V_j updates the precision matrices will be sparse, making sampling from the conditionals computationally tractable.

C Binomial and related likelihoods via Pólya–Gamma augmentation

When the likelihood is binomial, $y_{ijt} \sim \text{Bin}(n_{ijt}, 1/\{1 + e^{w_i^\top v_{jt}}\})$, where n_{ijt} is a nuisance parameter, the updates are conditionally conjugate given a Pólya–Gamma (PG) latent variable sample [48],

$$(\psi_{ijt} \mid -) \sim \text{PG}(n_{ijt}, w_i^\top v_{jt}), \quad (w_i \mid -) \sim N(m_{\psi_i}, \Sigma_{\psi_i}), \quad (20)$$

where $\Sigma_{\psi_i} = (\tilde{V}^\top \Psi_i \tilde{V} + \sigma^{-2} I)^{-1}$, $m_{\psi_i} = \Sigma_{\psi_i} \tilde{V}^\top \kappa$, $\Psi_i = \text{diag}(\psi_{(i,1,1)}, \dots, \psi_{(i,M,T)})$, and $\kappa = (y_{(i,1,1)} - n_{(i,1,1)}/2, \dots, y_{(i,M,T)} - n_{i,M,T}/2)$. The updates for V follow analogously. PG augmentation can be applied to binomial, Bernoulli, negative binomial, and multinomial likelihoods, among others.

D GASS Convergence Proof

Recall the GASS sampling algorithm.

- Start in some state x_0 :
- For $t = 0, 1, 2, \dots$:
 - Sample $v \sim \mathcal{N}(0, \Sigma)$ and $u \sim \text{unif}([0, 1])$.
 - Compute the sub-ellipse

$$\mathcal{E}_{x_t, v} = \{x' = (x_t - \mu) \cos \theta + v \sin \theta + \mu : \theta \in [0, 2\pi] \text{ and } x_t \in S\}.$$
 - Compute the region $\mathcal{X}_{x_t, u} = \{x \in \mathbb{R}^d : \mathcal{L}(x) \geq \mathcal{L}(x_t) + \log(u)\}$.
 - Sample the new state $x_{t+1} \sim \text{unif}(\mathcal{E}_{x_t, v} \cap \mathcal{X}_{x_t, u})$.

where $\mathcal{L}(\cdot)$ is some log-likelihood function such that $L(x) > -\infty$ for all $x \in \mathbb{R}^d$, $\mathcal{N}(\cdot; \mu, \Sigma)$ is the multivariate normal density, $S \subset \mathbb{R}^d$ is some subset of \mathbb{R}^d with positive Lebesgue measure (in our case it will be an area defined by some linear inequalities).

In this section, we show that the iterates x_t of the GASS chain converge to the correct stationary distribution, whose density is given by

$$P(x) = \frac{1}{Z} \exp(\mathcal{L}(x)) \mathcal{N}(x; \mu, \Sigma) \mathbf{1}[x \in S] \quad (21)$$

where Z is the normalizing constant to make the density integrate to one.

To begin, we introduce some notation. We say $T(x, \cdot)$ is a Markov transition if $T(x, \cdot)$ is a probability density over \mathbb{R}^d for all $x \in \mathbb{R}^d$. Let T^n denote the transition operator applied n times. Note that T , along with a starting state $x_o \in \mathbb{R}^d$, induces a Markov chain $(X_n)_{n=0}^\infty$, where

$$\Pr(X_{t+n} \in A | X_t = x) = \int_A T^n(x, y) dy =: T^n(x, A)$$

for all measurable $A \subset \mathbb{R}^d$. The *total variation distance between* two probability distribution μ, ν is given by

$$\|\mu(\cdot) - \nu(\cdot)\|_{TV} := \sup_{A \subseteq \mathbb{R}^d} |\mu(A) - \nu(A)|.$$

We say a distribution π is a *stationary* distribution of T if for all measurable $x, y \in \mathbb{R}^d$, we have

$$\int_{x \in \mathbb{R}^d} \pi(dx) T(x, dy) = \pi(dy).$$

We also say T is *reversible* with respect to a distribution π if for $x, y \in \mathbb{R}^d$,

$$\pi(dx) T(x, dy) = \pi(dy) T(y, dx).$$

Let ϕ be a non-zero σ -finite measure, we say that T is ϕ -irreducible if for all subsets $A \subset \mathbb{R}^d$ with $\phi(A) > 0$ and all $x \in \mathbb{R}^d$, there exists a positive integer n such that $T^n(x, A) > 0$.

Finally, we say that a Markov chain with transition operator T and stationary distribution π is *aperiodic* if there does not exist a partition A_1, \dots, A_m of \mathbb{R}^k such that $\pi(A_1), \dots, \pi(A_m) > 0$ and

$$T(x, A_{(t \bmod m)+1}) = 1$$

for all $t = 1, \dots, m$ and $x \in A_t$.

The following result, which is a simplified version of Theorem 4 in [49], tells us that reversibility, irreducibility, and aperiodicity are enough to guarantee convergence to stationarity.

Theorem 1. *Suppose T is a transition operator and π is a probability distribution over \mathbb{R}^d satisfying*

- (i) T is reversible with respect to a distribution π ,
- (ii) T is π -irreducible, and
- (iii) T is aperiodic,

then for π -a.e. $x \in \mathbb{R}^d$, we have that $T^n(x, \cdot)$ converges to $\pi(\cdot)$ in total variation, i.e.

$$\lim_{n \rightarrow \infty} \|T^n(x, \cdot) - \pi(\cdot)\|_{TV} = 0.$$

Thus, to show that the GASS algorithm converges to the correct stationary distribution, it suffices to show that it meets the conditions of Theorem 1. We start by showing that the GASS transition operator is P -irreducible.

Lemma 1. *The GASS transition operator is P -irreducible, where P is given in equation (21).*

Proof. Let $A \subset \mathbb{R}^d$ such that $P(A) > 0$ and let $x \in \mathbb{R}^d$. Since $P(A) > 0$, we know there must exist $x_o \in A$ and $r_o > 0$ such that

- (i) $B(x_o, r_o) := \{z \in \mathbb{R}^d : \|z - x_o\| \leq r_o\} \subset A$ and
- (ii) $P(B(x_o, r)) > 0$ for each $r \in (0, r_o)$.

Furthermore, let $L_{\min} = \inf_{z \in B(x_o, r_o)} \mathcal{L}(z)$. Note that $L_{\min} > -\infty$. Then there is some positive probability that we sample $u \sim \text{unif}([0, 1])$ satisfying $\mathcal{L}(x) + \log(u) \leq L_{\min}$. Let us condition on this happening.

From (ii) above, we know that with some positive probability we choose some $v \in B(x_o, r_o/2)$ when sampling from $\mathcal{N}(0, \Sigma)$. Let us condition on this happening. Then there is some subset $B \subseteq [0, 2\pi]$ with positive Lebesgue measure such that

$$x \cos \theta + v \sin \theta \in B(v, r_o/2) \subseteq A$$

for all $\theta \in B$. Putting it all together, with positive probability, GASS transitions from x to A in a single step. \square

We now show that the GASS transition operator is reversible with respect to P . Note that the core ideas of this proof already appeared in the reversibility argument of elliptical slice sampling (ESS) [40].

Lemma 2. *The GASS transition operator is reversible with respect to P .*

Proof. Let $T(x, \cdot)$ denote the transition operator of GASS. Let $x, x' \in \mathbb{R}^d$. Note that if either x or x' is not in S , then we have

$$P(dx)T(x, dx') = 0 = P(dx')T(x', dx).$$

Thus, assume that $x, x' \in S$.

Now consider any $v \in \mathbb{R}^d$ such that there exists a $\theta \in [0, 2\pi]$ satisfying $x' = (x - \mu) \cos \theta + v \sin \theta + \mu$. Then note that if $v' = v \cos \theta - (x - \mu) \sin \theta$, we may rewrite

$$\begin{aligned} x &= (x' - \mu) \cos \theta + v' \sin \theta + \mu \\ v &= (x' - \mu) \sin \theta + v' \cos \theta. \end{aligned}$$

Note that the transformation $(x', v') \rightarrow ((x' - \mu) \cos \theta + v' \sin \theta + \mu, (x' - \mu) \sin \theta + v' \cos \theta) = (x, v)$ has Jacobian matrix with the following block structure:

$$J = \begin{pmatrix} I_d \cos \theta & I_d \sin \theta \\ I_d \sin \theta & I_d \cos \theta \end{pmatrix}$$

We can calculate the determinant of this matrix as

$$\begin{aligned} \det(J) &= \det(I_d \cos \theta + (I_d \sin \theta)(I_d \cos \theta)^{-1} I_d \sin \theta) \det(I_d \cos \theta) \\ &= (\cos^2 \theta + \sin^2 \theta)^d = 1. \end{aligned}$$

Moreover, we have

$$\begin{aligned} &\mathcal{N}(x; \mu, \Sigma) \mathcal{N}(v; 0, \Sigma) \\ &= \mathcal{N}((x' - \mu) \cos \theta + v' \sin \theta + \mu; \mu, \Sigma) \mathcal{N}((x' - \mu) \sin \theta + v' \cos \theta; \mu, \Sigma) \\ &= (2\pi)^{-d} \det(\Sigma)^{-1} \exp \left(-\frac{1}{2} ((x' - \mu) \cos \theta + v' \sin \theta)^\top \Sigma^{-1} ((x' - \mu) \cos \theta + v' \sin \theta) \right. \\ &\quad \left. -\frac{1}{2} ((x' - \mu) \sin \theta + v' \cos \theta)^\top \Sigma^{-1} ((x' - \mu) \sin \theta + v' \cos \theta) \right) \\ &= (2\pi)^{-d} \det(\Sigma)^{-1} \exp \left(-\frac{1}{2} (x' - \mu)^\top \Sigma^{-1} (x' - \mu) \cos^2 \theta - \frac{1}{2} (x' - \mu)^\top \Sigma^{-1} (x' - \mu) \sin^2 \theta \right. \\ &\quad \left. -\frac{1}{2} v'^\top \Sigma^{-1} v' \cos^2 \theta - \frac{1}{2} v'^\top \Sigma^{-1} v' \sin^2 \theta \right) \\ &= \mathcal{N}(x'; \mu, \Sigma) \mathcal{N}(v'; 0, \Sigma) \end{aligned}$$

Putting the above together, we have

$$\mathcal{N}(x; \mu, \Sigma) \mathcal{N}(v; 0, \Sigma) dx dv = \mathcal{N}(x'; \mu, \Sigma) \mathcal{N}(v'; 0, \Sigma) dx' dv'.$$

Moreover, the elliptical regions satisfy $\mathcal{E}_{x,v} = \mathcal{E}_{x',v'}$.

Finally, define $f(z|x, v)$ as the transition probability density function from x conditioned on having chosen v to form the elliptical region $\mathcal{E}_{x,v}$. Then it is not hard to see that

$$\exp(\mathcal{L}(x))f(x'|x, v) = \exp(\mathcal{L}(x'))f(x|x', v').$$

Putting everything together, we have the desired result:

$$\begin{aligned} P(dx)T(x, dx') &= \frac{1}{Z} \exp(\mathcal{L}(x))\mathcal{N}(x; \mu, \Sigma) \int_{\mathbb{R}^d} \mathcal{N}(v; 0, \Sigma) f(x'|x, v) dx' dv dx \\ &= \frac{1}{Z} \exp(\mathcal{L}(x'))\mathcal{N}(x'; \mu, \Sigma) \int_{\mathbb{R}^d} \mathcal{N}(v'; 0, \Sigma) f(x|x', v') dx' dv dx \\ &= P(dx')T(x', dx). \end{aligned} \quad \square$$

Finally, we show that the GASS transition operator is aperiodic.

Lemma 3. *The GASS transition operator is aperiodic.*

Proof. From Lemma 2, we know P is a stationary distribution of P . So take A_1, \dots, A_m to be any partition of \mathbb{R}^k satisfying $P(A_1), \dots, P(A_m) > 0$. Then from the proof of Lemma 3, we know that for any $x \in A_1$

$$T(x, A_1) > 0$$

where T is the GASS transition operator. Thus, for any $x \in A_1$, we have $T(x, A_2) < 1$, which implies that GASS is aperiodic. \square

Putting it all together, we have that the GASS chain converges to the desired distribution.

Theorem 2. *For P -a.e. $x \in \mathbb{R}^d$, the GASS chain starting at x converges to P in total variation distance.*