

An ANOVA Model for Dependent Random Measures

Maria De Iorio¹, Peter Müller², Gary L. Rosner²,
Steven N. MacEachern³

Abstract

We consider dependent non-parametric models for related random probability distributions. For example, the random distributions might be indexed by a categorical covariate indicating the treatment levels in a clinical trial, and might represent random effects distributions under the respective treatment combinations. We propose a model that describes dependence across random distributions in an ANOVA type fashion. We define a probability model in such a way that marginally each random measure follows a Dirichlet process (DP) and use the dependent Dirichlet process (MacEachern, 1999) to define the desired dependence across the related random measures. The resulting probability model can alternatively be described as a mixture of ANOVA models, with a DP prior on the unknown mixing measure. The main features of the proposed approach are ease of interpretation and computational simplicity. Since the model follows standard ANOVA structure, interpretation and inference parallels conventions for ANOVA models. This includes the notion of main effects, interactions, contrasts, etc. Of course, the analogies are limited to structure and interpretation. The actual objects of the inference are random distributions instead of the unknown normal means in standard ANOVA models. Besides interpretation and model structure, another important feature of the proposed approach is ease of posterior simulation. Since the model can be rewritten as a DP mixture of ANOVA models it inherits all computational advantages of standard DP mixture models. This includes availability of efficient Gibbs sampling schemes for posterior simulation and ease of implementation of even high dimensional applications. Complexity of implementing posterior simulation is – at least conceptually – dimension independent.

¹Imperial College, London, U.K.

²Department of Biostatistics, The University of Texas M. D. Anderson Cancer Center, Houston, TX.

³Department of Statistics, The Ohio State University, Columbus, OH.

Research supported by NIH/NCI grant 2 R01 CA75981-04A1 and NSF grant DMS-0072526

1 Introduction

We consider dependent non-parametric models for related random probability distributions or functions. A typical application arises in modeling random effects distributions for related submodels in a hierarchical model. We propose a model that describes dependence across random distributions in an ANOVA type fashion. Specifically, assume that random distributions F_x are indexed by a p -dimensional vector $x = (x_1, \dots, x_p)$ of categorical covariates. For example, in a clinical trial F_{x_1, x_2} could be the random effects distribution for patients treated at levels x_1 and x_2 of two drugs. We define a non-parametric probability model for F_x in such a way that marginally, for each x , the random measure F_x follows a Dirichlet Process (DP), $DP(M, F_x^o)$, with total mass parameter M and base measure F_x^o (Ferguson, 1973). But we introduce dependence across x , i.e, dependence for $(F_x, x \in X)$, using the dependent Dirichlet process (DDP) as defined by MacEachern (1999, 2001). The random measures F_x are almost surely discrete with the point masses generated marginally from the base measure F_x^o . The DDP introduces dependence across x by imposing dependence in the distribution of these point masses. We use the DDP to define ANOVA type dependence across related random measures by assuming ANOVA models for these point masses. The resulting probability model defines an overall average effect and offsets for each level of the categorical covariates. If desired this can be generalized to include interaction effects. We propose a Markov chain Monte Carlo scheme to implement full posterior inference in the proposed model.

Our model is based on DP prior distribution (Ferguson, 1973; Antoniak, 1974). The DP is a probability model for random probability distributions. It plays a central role in non-parametric Bayesian inference and it has been successfully applied in many problems. One of the critical properties is the a.s. discreteness of a random measure $F \sim DP(M, F_0)$. Letting $\delta(x)$ denote a point mass at x we can write $F = \sum_{h=0}^{\infty} w_h \delta(\theta_h)$. Here w_h are the weights of point masses at locations θ_h . Sethuraman (1994) gives a constructive definition of the DP. The weights are generated from rescaled Beta distributions, $w_h / \prod_{i=1}^{h-1} (1 - w_i) \sim Be(1, M)$, and the locations θ_h are i.i.d. samples from the base measure F_0 .

Another property that will feature importantly in the following discussion is the Polya Urn representation for the marginal distribution of a sample from a random DP distribution. Assume $y_i \sim F, i = 1, \dots, n$, is sampled from an unknown distribution F which in turn, is generated by a DP, $F \sim DP(M, F_0)$. The marginal distribution of $y = (y_1, \dots, y_n)$ is described by the following Polya Urn scheme (Blackwell and MacQueen, 1973): $y_1 \sim F_0$ and

$$p(y_m | y_1, \dots, y_{m-1}) = \begin{cases} \delta(y_i) & \text{with probability } 1/(M + m - 1), i = 1, \dots, m - 1 \\ F_0 & \text{with probability } M/(M + m - 1). \end{cases} \quad (1)$$

The m -th sample point is either a tie with a previous sample y_i , or it is a new draw from the base measure. The positive probability of ties in (1) is due to the discrete nature of the random distribution $F \sim DP(M, F_0)$.

In many data analysis applications this discreteness is inappropriate. DP mixture models (DPM) avoid this discreteness in the sampling distribution by adding an additional convolution with a continuous kernel. The typical DPM model assumes

$$y_i \stackrel{iid}{\sim} H \text{ with } H(y) = \int f(y|\mu) dF(\mu), \quad F \sim DP(M, F_0), \quad (2)$$

i.e., a mixture with a DP prior on the random mixing measure F . Many applications use a normal kernel $f(y|\mu) = N(\mu, S)$ with a common covariance matrix S , leading to a discrete mixture of normals $H(y) = \sum_{h=1}^{\infty} w_h N(\mu_h, S)$. One of the main attractions of DPM models like (2) is computational simplicity. Also, posterior simulation algorithms are dimension independent. See, for example, Escobar and West (1998) or MacEachern and Müller (2000) for a review of models based on (2).

Several papers have considered extension of DP and DPM models to hierarchical models over related random distributions. In the context of parametric models, i.e., models with finite dimensional parameter vector, such hierarchies with submodels for related experiments are standard modeling tools. In non-parametric models such extensions are complicated by the infinite dimensional nature of the random distribution. Some of the first developments of dependent DP models appear in Cifarelli and Regazzini (1978) who define dependence across related random measures by introducing a regression for the base measure of marginally DP distributed random measures. The model is used, for example, in Muliere and Petrone (1993) who define dependent non-parametric models $F_x \sim DP(M, F_x^o)$ by assuming a regression in the base measure, $F_x^o = N(\beta x, \sigma^2)$. Similar models are discussed in Mira and Petrone (1996) and Giudici et al. (2002). A similar strategy of linking dependent random DP measures F_x at the level of the base measure is used in Tomlinson and Escobar (1999). They achieve increased flexibility by assuming a DPM hyperprior on the common base measure. Gelfand and Kottas (2001) define dependent non-parametric models by considering a representation of random measures as products of DP distributed factors. This allows them to enforce stochastic ordering. Tomlinson and Escobar (1999) and Gelfand and Kottas (2001) are appropriate to model dependence across several related random measures. However, they are not naturally extended to include regression on covariates.

MacEachern (1999) defines the dependent DP (DDP) to allow a regression on a covariate x . Consider a family of random measures $\mathcal{F} = (F_x, x \in X)$ indexed by a covariate x . MacEachern (1999) defines a probabilty model for \mathcal{F} such that marginally, for each x , $F_x = \sum w_h \delta(\theta_{xh})$ follows a DP. We use an additional subindex x for the point masses θ_{xh} to indicate the point masses in the random measure F_x . In the basic DDP model the weights

w_h are common to all F_x . The DDP model induces dependence across x by assuming that $\theta_h = (\theta_{xh}, x \in X)$ are i.i.d. realizations of a stochastic process (in x). For example, θ_{xh} might be assumed a Gaussian process. Independence across h , together with the stick breaking prior for the weights w_h , guarantees that F_x marginally follows a DP. Dependence in the sample path of the stochastic process θ_h introduces the desired dependence across x . We use this DDP structure to develop an ANOVA like probability model over an array of random distributions. The DDP model provides a convenient starting point for the discussion. But the proposed model is more general. It can be rewritten as a DP mixture model. With minimal changes in the computational algorithms the DP can be replaced by any non-parametric model that allows a constructive definition by a stick breaking algorithm as in Sethuraman (1994). See, for example, Ishwaran and James (2001) for discussions of such probability distributions.

In Section 2 we develop the basic model as a dependent DP model. In Section 3.1 we rewrite the model as a DP mixture of ANOVA models. Building on this representation we discuss computational implementation issues. Section 3.3 discusses the use of the ANOVA DDP model to define random effects distributions in hierarchical models and other applications. Section 4 illustrates the proposed models with three examples. Section 5 concludes with a final discussion.

2 The ANOVA DDP

Assume $\mathcal{F} = (F_x, x \in X)$ is an array of random distributions, indexed by a categorical covariate x . For simplicity of explanation, assume for the moment that $x = (v, w)$ is bivariate with $v \in \{1, \dots, V\}$ and $w \in \{1, \dots, W\}$. The covariates (v, w) could be, for example, the levels of two treatments in a clinical trial, and the distributions F_x might be sampling distributions for recorded measurements on each patient, or distributions for random effects. In the latter case, an additional layer in the model hierarchy defines a sampling distribution for the observed outcomes conditional on the random effects.

In this context we wish to develop a probability model for the random distributions F_x that will enable us to build an ANOVA type dependence structure. For example, we want the random distributions F_x and $F_{x'}$ for $x = (v_1, w_1)$ and $x' = (v_1, w_2)$ to share a common main effect due to the common factor v_1 . The model should allow us to incorporate prior information about the presence of interaction between the covariates. If interactions are present, the effect of $v = v_1$ should be allowed to depend on the level of the other covariate w . The model described below gives a formal definition to notions like “main effect” and “interaction”. Briefly, instead of a non-zero additive effect on the mean of the response variable in an ANOVA model, an effect is recast as a difference in distribution of some

quantity that has, in turn, an impact on the distribution of the final response. Thus, the models we create allow us to transfer both the interpretation and structure used for unknown normal means in the traditional ANOVA model to unknown random distribution functions.

Like standard ANOVA models the proposed model can be justified by a judgement of partial exchangeability for observed data. Assume y_{xi} , $i = 1, 2, 3, \dots$ are observed data, with y_{xi} denoting the i -th observation under condition x , for example, the response of the i -th patient assigned treatment combination x in a clinical trial. If for each x , the subsequence $y_x = (y_{x1}, y_{x2}, \dots)$ is judged exchangeable, then by de Finetti's representation theorem y_{xi} can be assumed an i.i.d. sample from some distribution F_x . Bernardo and Smith (1994, chapter 4) show how additional assumptions, including in particular that sample mean and variance are sufficient predictive statistics, leads to an ANOVA model. Stopping short of the assumption about predictive sufficiency naturally leads to assuming a non-parametric prior for F_x

We achieve the desired model structure by using the DDP framework. Specifically, let $F_x = \sum w_h \delta(\theta_{xh})$ for $x = (v, w)$. We assume Sethuraman's (1994) stick breaking prior for the common weights, $w_h / \prod_{i=1}^{h-1} (1 - w_i) \sim Be(1, M)$. On the locations θ_{xh} we impose additional structure,

$$\theta_{xh} = m_h + A_{vh} + B_{wh}. \quad (3)$$

As in standard ANOVA models we need to introduce an identifiability constraint for interpretability. We may impose any of the standard constraints, for example, $A_{1h} = B_{1h} \equiv 0$. For the remaining parameters we assume $m_h \stackrel{iid}{\sim} p_m^o(m_h)$, $A_{vh} \stackrel{iid}{\sim} p_{Av}^o(A_{vh})$, and $B_{wh} \stackrel{iid}{\sim} p_{Bw}^o(B_{wh})$, with independence being across h , v and w . We refer to the joint probability model on \mathcal{F} as $(F_x, x \in X) \sim \text{ANOVA DDP}(M, p^o)$. The model is parametrized by the total mass parameter M and the base measure p^o on the ANOVA effects in (3). Marginally, for each $x = (v, w)$, the random distribution F_x follows a DP with mass M and base measure F_x^o given by the convolution of p_m^o , p_{Av}^o and p_{Bw}^o . Model (3) defines dependence across x by defining the covariance structure of the point masses θ_{xh} across x . As in standard ANOVA the structural relationships are defined by the additive structure (3) and the level of the dependence is determined by the variances in p_m^o , p_{Av}^o and p_{Bw}^o . For example, consider two treatment combinations $x = (v, 1)$ and $x' = (v, 2)$ and random samples $y_x \sim F_x$ and $y_{x'} \sim F_{x'}$. Assuming normal priors $p_m^o = N(\mu_m, \sigma_m^2)$, $p_{Av}^o = N(\mu_{Av}, \sigma_A^2)$, and $p_{Bw}^o = N(\mu_{Bw}, \sigma_B^2)$ we find the marginal covariance $\text{cov}(y_x, y_{x'}) = (\sigma_m^2 + \sigma_A^2)/(M + 1)$. Compare this with the covariance $\text{cov}(y_x, y_{x'}) = (\sigma_m^2 + \sigma_A^2)$ that would arise in a standard ANOVA with the base measure p^o as prior for the ANOVA effects. The ANOVA DDP model introduces an additional level of uncertainty by defining the random measures $(F_x, F_{x'})$. The resulting covariance structure remains unchanged except for the attenuation factor $1/(M + 1)$ corresponding to the additional uncertainty about F_x . The same result remains true for arbitrary ANOVA structure,

including more factors and possibly interactions. In general, the marginal covariance for the observable responses y_x is the covariance under the corresponding fully parametric ANOVA model, reduced by a factor $1/(M+1)$. The marginalization is over the random distributions F_x with respect to the ANOVA DDP prior.

Model (3) is not constrained to univariate distributions F_x . The point masses θ_{xh} and the ANOVA effects m_h, A_{vh}, B_{wh} can be q -dimensional vectors. This is important, for example, if the random distributions F_x are used as random effects models in a hierarchical model. In the example discussed in Section 4.2 we use 7-dimensional random effects vectors. It is a critical advantage of the ANOVA DDP model that model specification and computation are dimension independent.

Another important generalization of model (3) is to more complex ANOVA structure. The model is easily generalized to a p -dimensional categorical covariate $x = (x_1, \dots, x_p)$:

$$\theta_{xh} = m_h + \sum_{i=1}^p A_i(x_i), \quad (4)$$

where $A_i(x_i)$ is the main effect due to treatment x_i . Further extensions to include interactions A_{ij} , etc., are equally straightforward. For example, in (3) we could introduce additional terms $C_{vw,h}$. See the examples in Section 4 for an illustration. However, the same caveat as with corresponding parametric models applies. Meaningful inference for interactions can only be achieved with sufficiently many data observed under different covariate levels x .

Like standard Bayesian ANOVA, the model allows us to incorporate differential prior information for the various levels of the covariate. This is accomplished through choice of different prior distributions $p_{A_v}^o$ for the different levels of v . In the context where $v = 1$ indicates a control and $v = 2, \dots, V$ are exchangeable treatments, we might take $p_{A_v}^o$ to be degenerate at 0 for the control and to be an identical distribution with a larger spread for each of the treatments. As an analog of a linear contrast in standard ANOVA, we might take the distributions $p_{A_v}^o$ to have non-zero means falling along a line; including further structure on the means of these distributions lets us expand our models in a fashion similar to the classical expansion through orthogonal polynomials, though the realizations will not exactly follow the possibly lower dimensional model.

One can also place constraints on the estimated effects. Enforcing a dependence above that forces the A_{vh} to lie on a line (to do so, we need to violate the condition of independence of A_{vh} across levels of v) produces a lower-dimensional component in the model. Alternatively, a constraint such as monotonicity of the effect A_{vh} in v can be enforced. Such a constraint ensures that the random distributions F_x are stochastically ordered with respect to v . This type of constraint is meaningful, for example, if v is the toxicity level of an anti cancer agent in a chemotherapy treatment.

3 Mixture of ANOVA Models

3.1 A DP Mixture of ANOVA Models

Most applications of DP models in data analysis add an additional layer in the model to convolve the discrete measure generated from a DP with a continuous, typically normal, kernel, to construct DPM models as in (2). For the same reasons we propose to add an additional normal mixture to the ANOVA DDP model. This leads to models of the form

$$\begin{aligned} (y_i | x_i = x) &\sim H_x(y_i), \text{ with } H_x(y) = \int N(y | \mu, S) dF_x(\mu), \\ (F_x, x \in X) &\sim \text{ANOVA DDP}(M, p^\circ), \end{aligned} \quad (5)$$

with appropriate hyperpriors for the common normal variance S and the ANOVA DDP parameters.

Implementation of posterior inference in the ANOVA DDP model (5) is most easily developed on the basis of an equivalent reformulation of the model as a mixture of ANOVA models. Consider, for example, model (3) for q -dimensional random measures F_x . Let N denote the number of ANOVA effects in (3) and let $\alpha_h = [m_h, A_{2h}, \dots, A_{Vh}, B_{2h}, \dots, B_{Wh}]$ denote the $(q \times N)$ matrix of ANOVA parameters corresponding to the h -th point mass in the random measures. Let d_i denote a design vector to select the appropriate ANOVA effects corresponding to x_i , i.e., $\theta_{xh} = \alpha_h d_i$ for $x = x_i$. Using this notation, model (5) with base measure $(p_m^\circ, p_{Av}^\circ, p_{Bw}^\circ)$ can be rewritten as:

$$(y_i | x) \sim H_x(y_i), H_x(y) = \int N(y | \alpha d_i, S) dF(\alpha), \text{ and } F \sim DP(M, p^\circ). \quad (6)$$

In words, data y_i are sampled from a mixture of ANOVA models, with a DP prior on the unknown mixing measure. As usual in mixture models posterior simulation is based on breaking the mixture in (6) by introducing latent variables α_i ,

$$y_i = \alpha_i d_i + \epsilon_i, \quad \alpha_i \sim F \text{ and } F \sim DP(M, p^\circ), \quad (7)$$

with $\epsilon_i \sim N(0, S)$. It follows from this equivalence that any Markov chain Monte Carlo (MCMC) scheme for DP mixture models can be used for posterior simulation in DDP ANOVA models of the type (5). The conjugate nature of the base measure p° and the kernel in the error distribution $p(\epsilon_i)$ in (7) greatly simplify posterior simulation. See, for example, MacEachern and Müller (1998) for details of posterior MCMC simulation for DPM models. A description of the relevant modifications needed for the ANOVA DDP is given in Appendix A. Since the MCMC simulation proceeds by marginalizing with respect to the unknown measures F_x , inference about the unknown distribution F_x itself is not a standard

part of MCMC in DP mixture models and can be difficult. See Gelfand and Kottas (2002) for a discussion and for inference about F_x in general. However, some important simplifications are possible. For any fixed y , the value of $H_x(y)$ in (6) takes the form of a linear functional of a DP random measure for which Regazzini et al. (2002) give results for exact inference. Guglielmi and Tweedie (2001) discuss an MCMC implementation to evaluate linear functionals. Further, in many applications the desired inference on F_x is constrained to reporting the posterior mean measures $E(H_x | Y)$, and sampling from $p(H_x | Y)$. Here Y denotes the observed data. But the posterior mean $E(H_x | Y)$ is equal to the posterior predictive distribution $p(y_{n+1} | x_{n+1} = x, Y)$ for a future observation in the DPM model (6). We can exploit this to easily evaluate $E(H_x | Y)$. Proceeding with Rao-Blackwellization as usual in MCMC implementations (Casella and Robert, 1996) it is possible to further simplify computations and estimate $E(H_x | Y)$ as an average of complete conditional posterior distributions that are already evaluated in the course of the MCMC simulation. Details of this estimate as well as a computationally efficient algorithm for approximate sampling from $p(H_x | Y)$ are described in Appendix A.

While the DP mixture representation of the ANOVA DDP model is helpful to derive posterior simulation methods, the proposed model goes beyond standard DP mixture models. One crucial difference between the DDP model that we present and corresponding DPM models is in where the models are defined. The DPM models are defined at a fixed, finite set of covariates while the DDP model may be defined over a much broader space. The study in Section 4.2 concerns a factorial design with two factors (doses of two drugs) studied at 3 and 4 levels, respectively. Potential doses for each drug range along a continuum. The DDP model specifies a prior distribution which, when updated with the data, allows us to make posterior inference at an *arbitrary* combination of doses. Such a specification is consistent with the DDP model of reduced dimension that we use for computational purposes. Of course such inference would require the prior specification for the ANOVA effects to extend over a continuum of covariates. For example, in (3) we could use a Gaussian process prior on A_{vh} as a function of a continuous dose v .

Inference for a DDP model with a typical prior specification has a particularly attractive feature, namely appropriate asymptotic degeneracy and non-degeneracy. Recall that $\alpha = [m, A_2, \dots, A_V, B_2, \dots, B_W]$, $\alpha \sim F$, is the matrix of ANOVA effects in (6). Suppose that an augmented set of treatment combinations includes a dose combination x^* that was not recorded in the observed data, and let d^* denote the design vector corresponding to x^* . If the base measure p^o is a non-singular multivariate normal distribution, then the conditional distribution of the mean effect αd^* at the new dose combination, conditional on the mean effect at the studied dose combinations, is a non-degenerate normal distribution. This conditional non-degeneracy is typical of a DDP prior specification. The consequence of

non-degeneracy is that, even if αd_x at the studied dose combinations x were known exactly, the posterior distribution at the new dose combination would be non-degenerate. Thus, even if the posterior distribution on the mean effect at observed treatment combinations concentrates asymptotically, the posterior distribution at other dose combinations does not concentrate. This behavior is in keeping with our prior beliefs in most settings.

The parameter M induces a distribution on the number of clusters into which the observations fall. Clusters are defined by ties in the latent α_i introduced in (7). In DPM models, summaries of this distribution, such as its mean and variance are often used to judge whether a particular prior distribution provides a match to prior beliefs. See, for example, Escobar (1994), who also uses these summaries to motivate a distribution over M . The DDP model which we use in this paper relies on a single mass parameter, M . For this DDP model, clusters of observations occur both within a treatment and across treatments. The number of clusters is stochastically increasing in the number of observations. See Appendix B for a formal statement and proof. Consequently, varying sample sizes across treatments imply different distributions for the number of within-treatment clusters as well as a different distribution of the number of clusters across all treatments.

For illustration we computed summaries of the prior distribution on the number of clusters in a setup as in the case study analyzed in Section 4.2, i.e., we assume the same number of dose combinations x and the same sample sizes within each dose combination x . Let k_x denote the number of clusters of observations with covariate x , and let k denote the number of clusters among all observations. Note that $k \leq \sum_x k_x$ since clusters can extend over multiple values of x . Considering four alternative prior assumptions for the total mass parameter M we evaluated prior mean and standard deviation for the number of clusters. Specifically, we computed prior summaries for k_x and k under a Gamma prior, $M \sim Ga(5, 0.5)$, and assuming fixed values $M = 1, 10$ and 25 . Let n_x denote the number of observations with covariate level x . For a treatment with $n_x = 6$ observations, and assuming $M \sim Ga(5, 0.5)$, $M = 1, 10$ and 25 , we find the prior mean $E(k_x)$ (standard deviation $SD(k_x)$) to be 4.8 (1.0), 2.5 (0.98), 4.9 (0.91) and 5.5 (0.68), respectively. For a treatment with $n_x = 10$ observations we find prior moments of 6.9 (1.6), 2.9 (1.2), 7.2 (1.3) and 8.6 (1.1), respectively. For a treatment with a relatively large number of observations, $n_x = 12$, we find prior moments of 7.9 (1.9), 3.1 (1.2), 8.2 (1.5) and 10.0 (1.1), respectively. Finally, for the entire study, with $n = 52$ patients, we find prior moments for k of 18.0 (5.2), 4.5 (1.7), 18.7 (3.1), and 28.5 (3.3), respectively. In summary, the impact of the different numbers of observations (6, 10, and 12) at the different treatments is relatively modest. There is, however, a substantial difference when moving from the individual treatments with $n_x = 6, 10$ and 12 , to the entire experiment, with $n = 52$.

3.2 Other Mixture of ANOVA Models

Rewriting the ANOVA DDP as (7) highlights the generality of the underlying model structure. The use of a DP prior for the discrete mixing measure is motivated by technical convenience, and because the parsimonious parametrization of the DP avoids difficult prior elicitation problems. On the other hand, the fact that the DP is parametrized by a base measure and one scalar precision parameter M only could sometimes be a limitation. Green and Richardson (2001) argue for the use of more general mixture models and show appropriate posterior simulation schemes. Muliere and Tardella (1998), Gelfand and Kottas (2002) and Ishwaran and James (2001) discuss finite truncations of DP priors. Ishwaran and James (2001) propose alternative non-parametric priors based on similar stick-breaking representations. Any of these non-parametric priors can be substituted in (7) without changing the model structure, and requiring only minimal changes in the posterior simulation schemes.

3.3 Hierarchical Models and Other Extensions

Consider a generic hierarchical model of the form

$$y_i \sim p(y_i|\theta_i), \quad (\theta_i|x_i = x) \sim H_x(\theta|\phi). \quad (8)$$

In words, data y_i for the i -th sampling unit, e.g., a patient, is sampled from a probability model parametrized by a random effects vector θ_i . For example, this could take the form of a non-linear regression

$$y_{ij} = f(t_{ij}; \theta_i) + \epsilon_{ij} \quad (9)$$

with a mean function $f(\cdot; \theta)$ parametrized by θ_i and evaluated at known times t_{ij} , $j = 1, \dots, n_i$. The θ_i are generated from a random effects distribution H_x . The random effects distribution depends on a covariate specific to the sampling unit and possibly additional hyperparameters ϕ .

If the covariates are, for example, treatment indicators for the i -th patient, then ANOVA DDP models as in (5) are appropriate prior probability models for $(H_x, x \in X)$. The random effects vector θ_i takes the place of y_i in (5), and $H_x(\cdot) = \int N(\cdot|\mu, S) dF_x(\mu)$. In general, the ANOVA DDP model can be used whenever the random effects distributions H_x are indexed by some categorical covariates x_i specific to the i -th unit with a notion of ANOVA type dependence across the random distributions. Section 4.2 discusses a typical example. Posterior inference is implemented as in (5), with an additional step to update the random effects vectors θ_i which now replace the data y_i in the ANOVA DDP model (5). The details of this step are problem specific. For example, if $p(y_i|\theta_i)$ is a non-linear regression as in (9) then updating θ_i amounts to a posterior draw from the parameters in a non-linear regression with

data y_{ij} , $j = 1, \dots, n_i$ and prior $\theta_i \sim N(\alpha_i d_i, S)$. Here α_i is the latent variable introduced to break the mixture model in (7). The latent variable α_i is imputed in the course of updating the ANOVA DDP model. See, for example, MacEachern and Müller (1998) for details.

The hierarchical model (8) is only one of many diverse modeling contexts and applications that gives rise to structure as in (5). The ANOVA DDP model (5) allows a convenient formalization of regression for survival data. Consider observations y_i , censoring indicators δ_i , with $\delta_i = 1$ if y_i is a censoring time and $\delta_i = 0$ if y_i is an observed event. Assume different subpopulations are indexed by a categorical covariate x_i . Using (2) to represent $p(y_i | x_i, \delta_i = 0)$ one can represent the sampling distribution of event times under x_i as a DP mixture model. Such models are used, for example, in Doss (1994), Gelfand and Kuo (1991) and Kuo and Mallick (1997). Kuo and Mallick (1997) include a regression in an accelerated failure time model. The DDP-ANOVA model allows to introduce dependence across subpopulations indexed by x_i when stronger assumptions like proportional hazards, additive hazards, or accelerated failure time are too stringent.

The basic ANOVA-DDP model (5) is easily extended to categorical or ordinal data. Let $z_i = (z_{i1}, \dots, z_{ip})$ denote a vector of observed categorical outcomes for experimental unit i . A multivariate probit model for z_i represents the sampling distribution for z_i by introducing a latent multivariate normal vector y_i . Specifically, assume $z_{ik} \in \{1, \dots, R\}$ is ordinal with R possible outcomes. We introduce cutoffs $-\infty = \theta_0 < \dots < \theta_r < \dots < \theta_R = \infty$ and assume $z_{ik} = r$ if and only if $\theta_{r-1} < y_{ik} \leq \theta_r$. See, for example, Albert and Chib (1993), Cowles et al. (1996), and Chen and Dey (2000) for inference in such models. Replacing the normal distribution of the latent variable y_i by an unknown distribution with a non-parametric prior provides a natural generalization of the probit model. Kottas et al. (2002) show that without loss of generality the cutoffs θ_j can be fixed when using a non-parametric prior model for the latent variables y_i . Consider now data recorded for different subpopulations indexed by x_i . For example x_i could be different hospitals, different treatments or different raters. The DDP-ANOVA model (5) formalizes a non-parametric regression of z_i on the covariates x_i . The categorical cell probabilities $\pi_{zx} = Pr(z | x)$ are modeled as corresponding quantiles of H_x , i.e., $\pi_{zx} = \int dH_x(y)$ with the integral extending over the range of y corresponding to categorical outcomes z .

Another interesting application is to classification. Assume we have data (x_i, y_i) , $i = 1, \dots, n$, under sampling model (5) for $p(y_i | x_i)$. Data could be repeated measurement data with the basic model being extended as in (8). Assume a new observation y_{n+1} with unknown label x_{n+1} is recorded, and interest focuses on inference about x_{n+1} . Augmented by a prior $p(x_i)$, the DDP-ANOVA model allows a model based flexible framework to implement the desired classification as $p(x_{n+1} | y_{n+1}, Y)$.

4 Examples

4.1 Simulation Example

Consider a two way bivariate ANOVA model with two factors v and w . Assume the number of levels are $V = 2$ and $W = 2$ for v and w , respectively. Let $\alpha = [m, A_1, A_2, B_1, B_2]$ denote

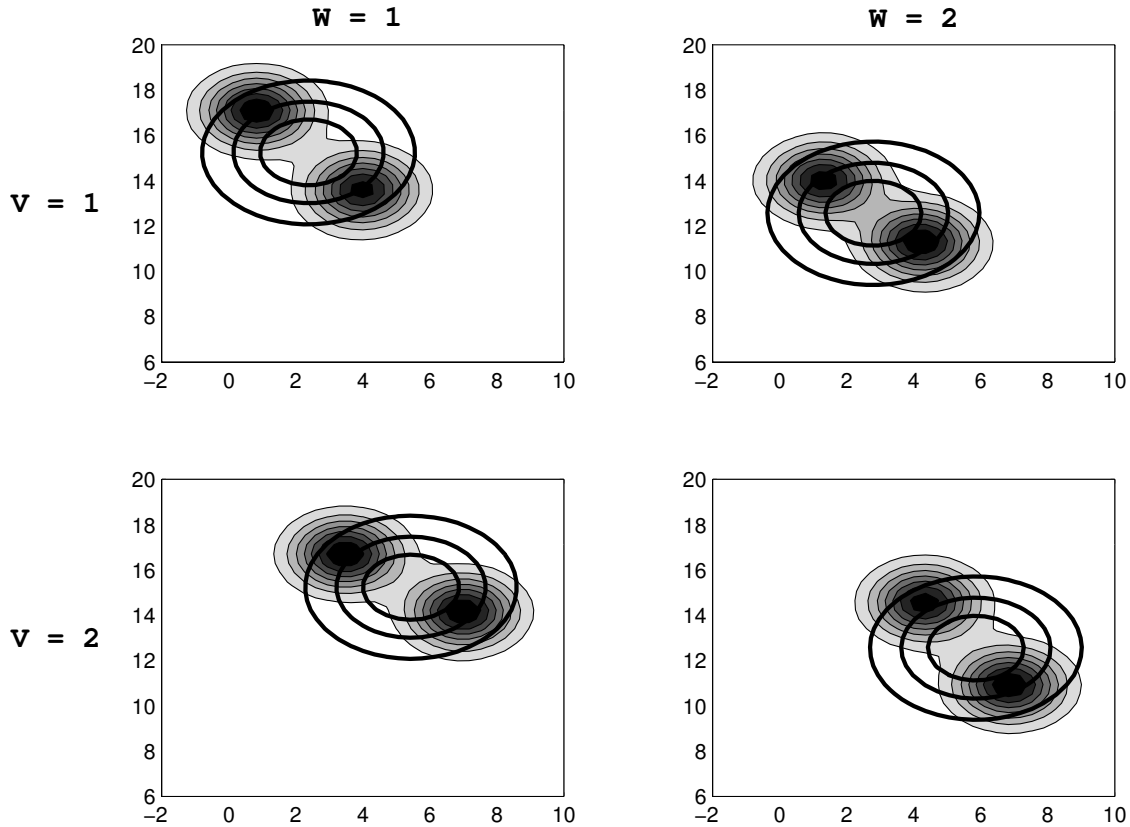


Figure 1: Posterior estimated distributions. The grey shades show the posterior mean distributions $E[H_{vw}(\cdot) | Y]$. For comparison the overlaid contours show the estimated distributions under a bivariate two-way ANOVA model, assuming a single bivariate normal sampling distribution for each combination of (v, w) .

overall mean and main effects for $v = 1, 2$ and $w = 1, 2$, respectively. Such data might arise, for example, when investigating differential expression patterns of proteins across major subtypes of ovarian cancer. Assume expression is measured by immunohistochemical staining of tissues. The bivariate outcome could be the level of expression of the two proteins of interest, and histologically different subtypes might be defined by type v (serous, endometrioid and mucinous) and degree of differentiation w (low and high).

We simulated 100 observations $y_i \sim N(\alpha_i d_i, I)$ with randomly selected designs d_i and generating

$$\alpha_i = \begin{cases} \begin{bmatrix} 3 & 2 & 5 & -1 & -1 \\ 1 & 6 & 6 & 7 & 4 \end{bmatrix} & \text{with probability 0.5} \\ \begin{bmatrix} 3 & 0.5 & 3.5 & -2.5 & -2.5 \\ 1 & 7.5 & 7.5 & 8.5 & 5.5 \end{bmatrix} & \text{with probability 0.5,} \end{cases}$$

i.e., we generate from model (6) with a two point discrete mixing measure $F(\alpha)$. We proceeded to fit the ANOVA DDP model (5) with an ANOVA structure as in (3). Let $H_{vw}(\cdot) = \int N(\cdot | \mu, S) dF_{vw}(\mu)$ denote the unknown distribution for factor levels (v, w) and let $Y = (y_1, \dots, y_n)$ denote the observed data. Figure 1 plots contours of the estimated distributions $E(H_{vw} | Y)$. Posterior inference correctly recovers the true mixture. For comparison we estimated a bivariate normal regression on (v, w) , i.e., a standard bivariate ANOVA model. For a fair comparison the prior probability model for the bivariate ANOVA model and the hyperprior for the base measure p^o in the DDP model are matched. In both models the residual covariance matrix S is assumed to be diagonal, $S = \sigma^2 I$, with a conjugate Gamma prior on $1/\sigma^2$. Thus both models have almost the same number of hyperparameters, with the only additional hyperparameters in the ANOVA DDP model being the prior parameters for the total mass M . The posterior mean bivariate ANOVA model is shown in Figure 1, as contours overlaid over the grey shades for the estimated DDP model. Forced to fit a single normal the bivariate ANOVA model approximates the mixture by a bivariate normal distribution centered between the locations of the two mixture terms. Failing to model the heterogeneity in the data the bivariate ANOVA model reports the mode of the distribution in a low density area between the two modes. The estimated residual variance is inflated to allow for the approximation. The marginal posterior means $E(\sigma^2 | Y)$ are 1.01 and 3.57 under the ANOVA DDP model and the bivariate ANOVA model, respectively.

4.2 Hierarchical Models

Müller and Rosner (1997) describe an analysis of hematologic data arising from a dose-escalation study (Lichtman et al. 1993). The data are white blood cell counts over time for each of n patients receiving relatively high doses of cancer chemotherapy. The treatment included a commonly used drug, cyclophosphamide, that is known to lower a person's white blood cell counts as the dose increases. Having very low white blood cell counts can be life-threatening, since these cells include the ones that are part of the human immune system. Thus, the investigators are keenly interested in knowing about the effect of dose on the white blood cell counts in order to guard their patients against severe and life-threatening toxicity.

Unfortunately, there is no standard measure in use for characterizing this form of toxicity. Is it the smallest WBC or is it the length of time a person's WBC is below a threshold value? Analyzing the entire profile of blood counts over time allows one to characterize depressed counts in any way one would like. The modeling approach of Müller and Rosner (1997) allows such full modeling. However, a critical limitation of the approach discussed in Müller and Rosner (1997) is the way how it implements regression on patient-specific covariates, in particular the assigned dose of chemotherapy. We re-analyze the data, using the DDP ANOVA model instead. The ANOVA structure allows us to characterize dose-specific differences in WBC profiles between the doses.

The data record white blood cell counts over time for each of $n = 52$ chemotherapy patients. Denote with y_{it} the measured response on day t for patient i . The profiles of white blood cell counts over time look similar for most patients. Figure 2 shows some typical patients. There is an initial base line, followed by a sudden decline when chemotherapy starts, and a slow, S-shaped recovery back to approximately base line after the end of the treatment. Profiles can be reasonably well approximated with a regression function with two change points corresponding to the beginning of the decline and the nadir, a horizontal line for the base line to the left of the first change point, a straight line for the rapid decline between the two change points, and a shifted and scaled logistic for the final S-shaped recovery to the right of the second change point. We parametrize this piecewise linear-linear-logistic regression with a 7-dimensional parameter vector θ (Müller and Rosner, 1997). But the non-linear regression parameters differ significantly across patients. Thus we introduce a patient specific random effects vector θ_i . Conditional on θ_i we assume a non-linear regression using the piecewise linear-linear-logistic regression model

$$y_{it} = f(t; \theta_i) + \epsilon_{it}. \quad (10)$$

The model is completed with a random effects model H_x . The random effects distribution H_x depends on the treatment levels for patient i . There are two treatments, the actual anti-cancer agent cyclophosphamide (CTX), and a second drug (GM-CSF) given to mitigate some adverse side effects of the chemotherapy. We impose an ANOVA structure on H_x with rows and columns in the two-way ANOVA indicating levels of CTX and GM-CSF. We code the levels of CTX as $v = 1, \dots, 4$ (corresponding to dose levels 1.5, 3.0, 4.5 and 6.0 gm/m^2) and the levels for GM as $w = 1, 2, 3$ (dose levels 2.5, 5 and 10 $\mu g/kg$). Let $x_i = (v_i, w_i)$ denote the treatment for patient i . The number of patients observed at each treatment combination x varies between $n_x = 6$ for $x \in \{(1, 3), (2, 1), (2, 2), (2, 3), (4, 2)\}$, $n_x = 10$ for $x = (3, 3)$, and $n_x = 12$ for $x = (3, 2)$. We assume

$$(\theta_i | x_i = x) \sim H_x(\theta), \quad (H_x, x \in X) \sim \text{ANOVA DDP}(M, p^o) \quad (11)$$

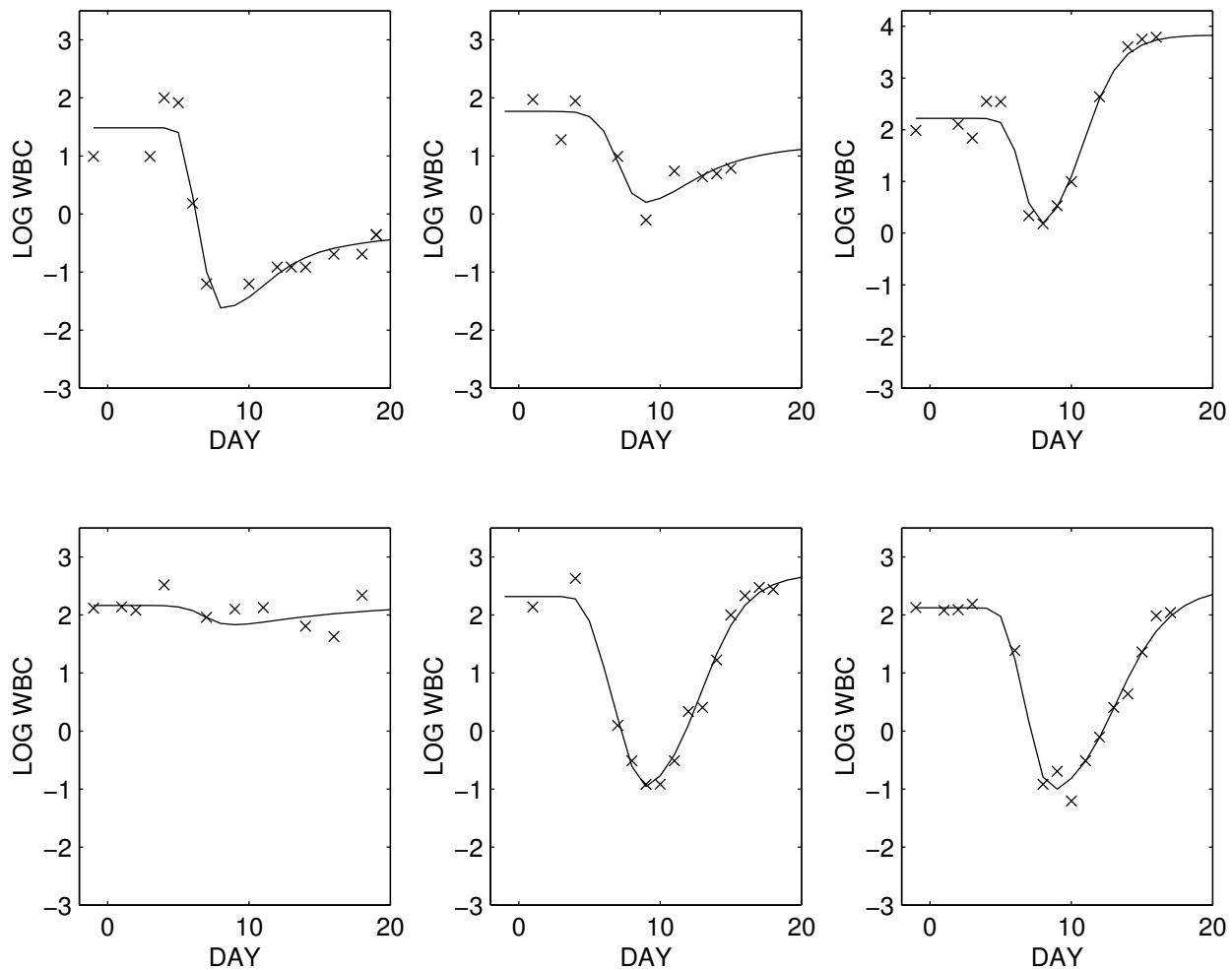


Figure 2: Some typical patients. The crosses plot observed log white blood cell count y_{it} for six typical patients. The curves show the posterior fitted mean response curves $f_{x_i}(t) = E \left\{ \int f(t; \theta) dH_{x_i}(\theta) | y \right\}$ for these patients.

Posterior predictive inference for future patients depends on the observed historical data only indirectly through learning about the random effects distribution (11). Conditional on the random effects distributions H_x , observed and future data are independent. Thus a structured, flexible hyperprior for H_x is critical to achieve the desired learning.

Figure 3 summarizes some critical aspects of the analysis. In this example, the random effects distributions $H_x(\cdot)$ are 7-dimensional. We summarize inference on $H_x(\cdot)$ by showing implied profiles. Let Y denote the observed data. For example, for $x = (1, 1)$ we show $f_x(t) = E \left\{ \int f(t; \theta) dH_x(\theta) | Y \right\}$, the posterior mean hematologic profile over time t for a patient treated with doses $x = (1, 1)$. Imposing identifiability constraints $A_{1h} = B_{1h} \equiv 0$ the corresponding design vector in (6) is $d = (1, 0, \dots, 0)$, including a main effect only. The figure shows the estimated mean curve $f_x(t)$ in the right lower panel.

For patients treated at other doses, we display corresponding offsets in a familiar ANOVA fashion. For example, for $x = (2, 1)$ the figure plots the posterior expected mean curve corresponding to design vector $d = (0, 1, \dots, 0)$ with an offset for $v = 2$. The second panel in the top row shows $f_x(t) = E \left\{ \int f(t; \theta = \alpha d) dF(\alpha) | Y \right\}$. The other panels in the same figure have analogous interpretations. Interpretation of Figure 3 needs to take into account that the additive ANOVA structure applies only for the ANOVA effects in (3). The plotted mean profiles are only convenient summaries of the 7-dimensional distributions, and are not additive. Since inference is based on posterior simulation, any other desired posterior summary can be derived. For example, Figure 4 shows posterior predictive distributions for the likely responses of a future patient treated at any of the 12 possible treatment combinations.

4.3 Categorical Outcome

Fluorescence activated cell-sorter (FACS) analysis is used to measure properties of cells in flow (Melamed et al., 1990). Freedman et al. (2002) use FACS to measure the proportion of monocytes in a blood sample that express certain surface markers. Using multi-color FACS analysis up to four different surface markers can be recorded simultaneously. The instrument reports the number of cells that express each of the identified surface markers beyond a certain preset threshold. If data is reported for blood samples collected under different treatment conditions $x_i \in X$ the resulting data format fits the categorical data model described in Section 3.3, with $z_{ik} \in \{0, 1\}$ reporting presence and absence of the k -th surface marker.

We set up a simulation study mimicking the setup in a currently ongoing clinical trial. The trial proposes chemoimmunotherapy for ovarian cancer. Two agents, denoted here as G and I, are used for the immunotherapy. FACS analysis is used to measure the expression

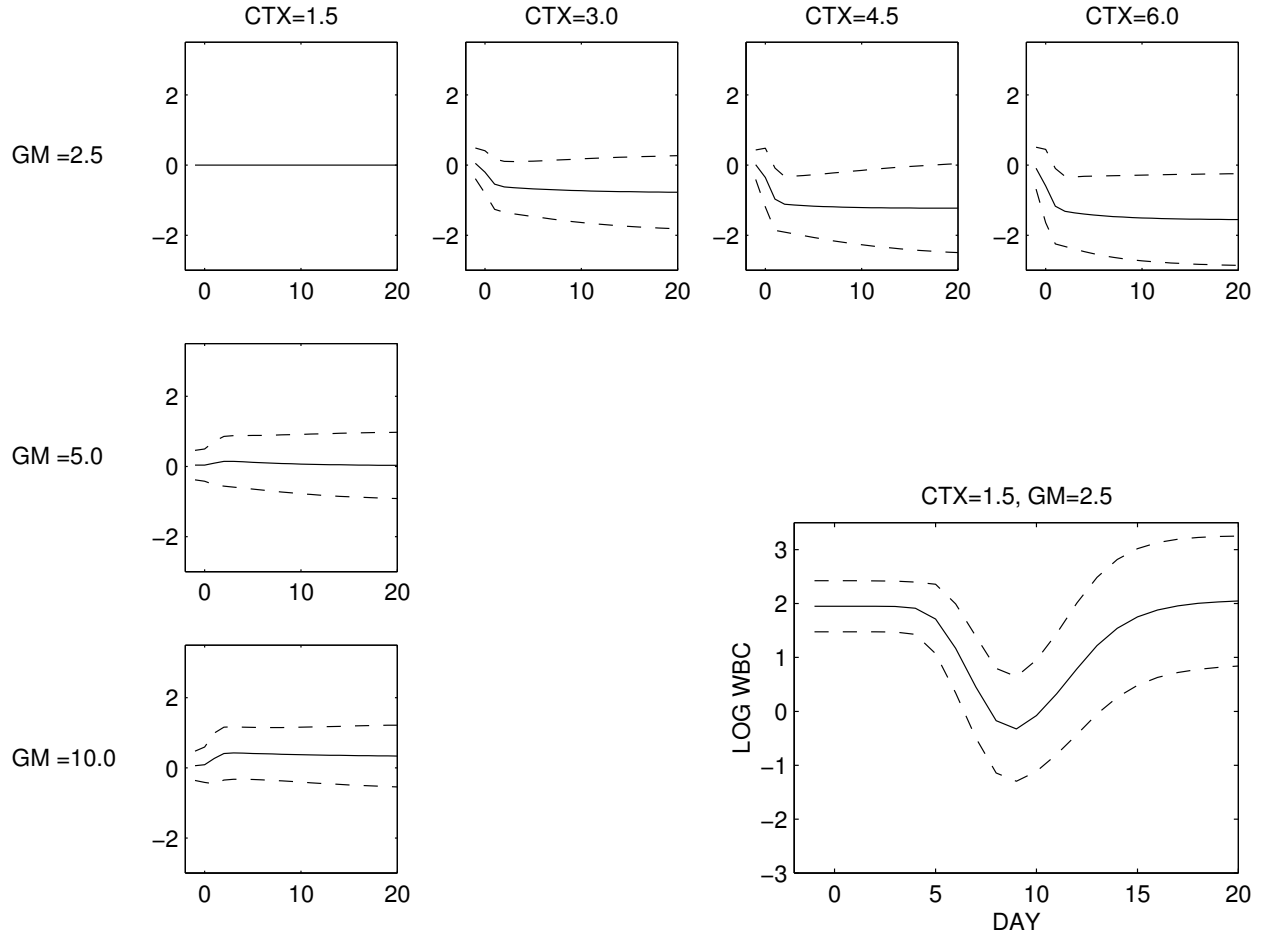


Figure 3: Estimated main effects A_v , B_w and overall mean effect m . For identifiability, we constrained $A_{1h} = B_{1h} \equiv 0$. The right lower panel shows the estimated hematologic profile for a patient treated with dose levels $(CTX = 1.5, GM = 2.5)$, i.e., corresponding to the overall mean m . The other figures summarize the offset corresponding to the respective main effect. For each ANOVA effect, the plot shows $E \left\{ \int f(t; \theta = \alpha d) dF(\alpha) \mid Y \right\}$, where d is the design vector corresponding to the respective ANOVA effect. The dashed curves show corresponding pointwise one posterior standard deviation margins. See the text for a detailed explanation.

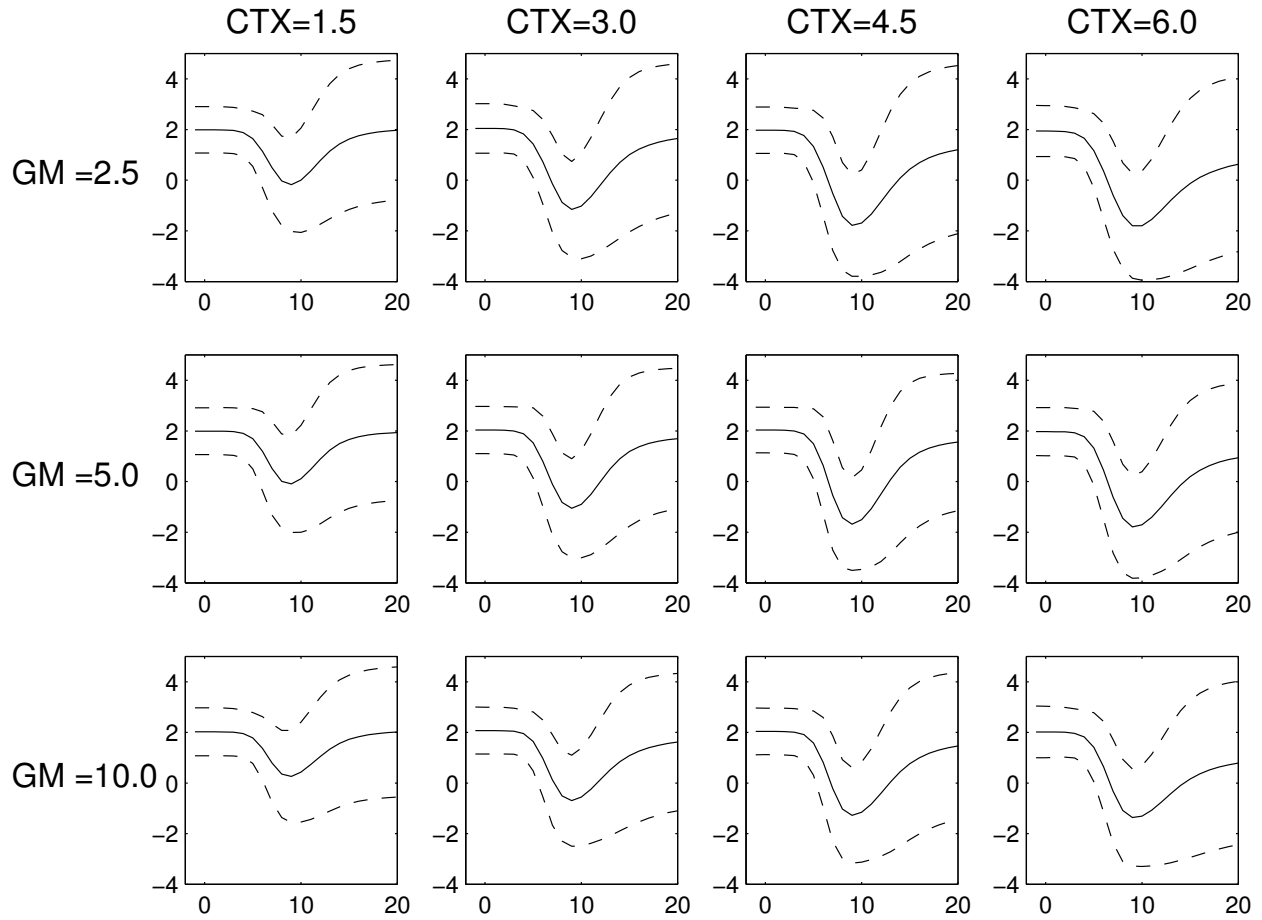


Figure 4: Posterior predictive distributions $p(y_{it}|x_i, y)$ for a future patient $i = n + 1$ assuming treatment combination $x_i = (v, w)$, arranged by v and w . The solid line shows the pointwise posterior predictive mean $E(y_{it}|x_i, y)$, plotted against t . The dashed lines show pointwise one posterior standard deviation margins.

Table 1: Observed (a) and fitted (b) cell probabilities $\pi_{zx} = Pr(z | x)$. Table (b) reports in parentheses the marginal posterior standard deviations. For comparison, table (c) reports the posterior mean cell probabilities under a bivariate probit model.

(a) Observed frequencies (in %)					(b) Estimated probabilities (in %)				
(G,I)	$z = (z_1, z_2)$				(G,I)	$z = (z_1, z_2)$			
	(1,1)	(1,0)	(0,0)	(0,1)		(1,1)	(1,0)	(0,0)	(0,1)
(0,0)	29	7	59	5	(0,0)	27 (5)	8 (3)	59 (6)	6 (3)
(0,1)	26	8	64	2	(0,1)	27 (5)	8 (4)	62 (6)	3 (2)
(1,0)	10	26	4	60	(1,0)	10 (4)	25 (6)	5 (3)	60 (6)
(1,1)	7	25	39	29	(1,1)	7 (4)	25 (5)	39 (6)	29 (6)

(c) Estimated probabilities (in %) (bivariate probit)				
(G,I)	(1,1)	(1,0)	(0,0)	(0,1)
(0,0)	16	21	45	18
(0,1)	12	22	51	15
(1,0)	28	8	24	40
(1,1)	15	17	46	21

of a variety of surface markers in blood samples collected under different treatment conditions. The first agent, G, has a growth factor stimulating effect. It can expand and prime monocytes, a type of white blood cell. The hypothesized effect of the second agent, I, is to serve as an activator to mature and fully activate these cells. As part of the protocol, blood samples will be analyzed for the number and activation status of monocytes. FACS analysis will be used to record expression of selected markers that are modified on activated monocytes. For this simulation study we consider two markers and samples under four different treatment conditions x_i , $x_i = (G, I) \in \{(0, 0), (0, 1), (1, 0), (1, 1)\}$. A pair of binary outcomes, $z_i = (z_{i1}, z_{i2})$ reports whether markers 1 and 2, respectively, are upregulated.

The conjectured role of I as an activator of cell populations expanded by G makes it important to include interactions in the model. We proceed with the ANOVA DDP model for ordinal data proposed in Section 3.3. We include main effects for G and I , plus an additional interaction $I * G$. We generated hypothetical data for four blood samples, one under each treatment combination. We assume 100 cells are measured in each sample, with observed frequencies reported in Table 1a. Let $\pi_{zx} = Pr(z | x)$ denote the probability of binary out-

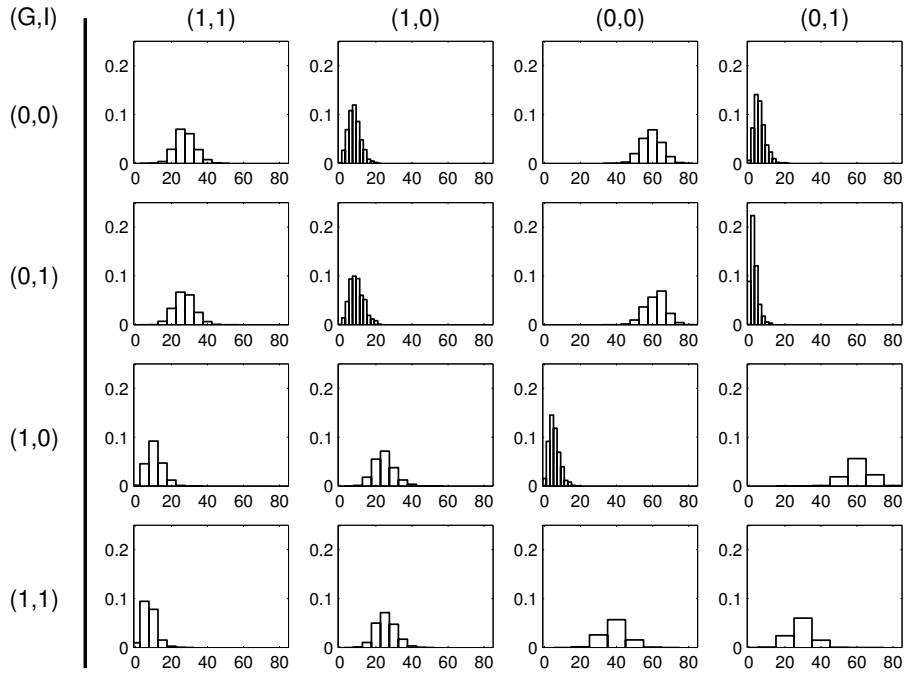


Figure 5: Posterior distribution for the unknown cell probabilities π_{zx} arranged by covariate $x = (G, I)$ (rows) and by outcome z (columns).

comes z under subpopulation x . Using the DDP ANOVA model the cell probabilities π_{zx} are expressed as the corresponding quantiles under H_x . For example, $\pi_{00} = \int_{-\infty}^0 \int_{-\infty}^0 dH_x(y_1, y_2)$, using the cutoff $\theta_1 = 0$.

Table 1b shows the posterior estimated cell probabilities, together with marginal posterior standard deviations. Marginal posterior distributions for the categorical cell probabilities are shown in Figure 5. For comparison Table 1c shows inference under a bivariate probit regression. To allow a fair comparison we use the base measure p^o as prior for the parametric probit model. Although the two models have a comparable number of hyperparameters, the non-parametric mixture in the ANOVA DDP model allows significantly more flexibility in fitting the observed frequencies.

5 Discussion

We have introduced a probability model for random distributions arranged in an ANOVA like array. The main features of the proposed model are ease of interpretation, facility to impose structure in the usual ANOVA like fashion, and efficient computation.

Limitations of the ANOVA DDP model are the need of MCMC simulation for posterior

inference, and the practical limitation to stick-breaking priors for the non-parametric mixing measure. Also, the model inherits limitations inherent in the DP prior. For example, the model includes only one scalar precision parameter for the random mixing measure. This makes it impossible to express different levels of prior precision across the sample space of the base measure. On the other hand, the parsimonious prior parametrization facilitates prior elicitation.

The ANOVA model that we have developed and illustrated in the examples relies on a particularly simple version of the DDP model. In this simple model, the weights of the point masses, w_h , do not depend on the level of the treatment. This model has the advantage of allowing us to rely on computational strategies developed for DPM models. More complex DDP models allow the w_h to vary across the treatments, but would necessitate more complex computational strategies that include additional parameters (the treatment-specific w_h or equivalent) in the Gibbs sampling steps. Inclusion of these additional parameters can, in other words, be described as forsaking a possible marginalization of them. The result of not marginalizing parameters when they can be marginalized is a poorer mixing Markov chain (Liu, 1994; MacEachern, 1994). An alternative approach to fitting the more complex DDP models is to use the model we have fit as an importance sampler. The two complex and simple models may be matched so as to provide the same marginal prior distributions at each level of the treatment. They would differ only in how the structure is connected across treatments, with this difference appearing in the prior distribution on the configuration vector, s (see the appendix for a description of s). As a consequence, the importance sampling weights would be determined by the two distributions on s .

Appendix

A. Posterior MCMC Simulation

We briefly describe the implementation of posterior simulation in the ANOVA DDP model (6). Since F is almost surely discrete (see Ferguson, 1973), there is a positive probability for ties among the α_i . Write $\{\alpha_1^*, \dots, \alpha_k^*\}$ for the $k \leq n$ distinct elements in $\{\alpha_1, \dots, \alpha_n\}$. Set $s_i = j$ iff $\alpha_i = \alpha_j^*$. Let n_j be the number of s_i equal to j , i.e. n_j is the size of the j th cluster and let $\Gamma_j = \{i : s_i = j\}$. Let η denote possibly unknown hyperparameters, including hyperparameters in the base measure p^o , the total mass parameter M and the covariance matrix S in (5). See, for example, MacEachern and Müller (1998) or Bush and MacEachern (1996) for a full description of Gibbs sampling scheme to estimate DPM models.

1. *Resampling s_i given all the other parameters:*

Marginalize over α_i and sample s_i from

$$\Pr(s_i = j \mid s_{-i}, y, \eta, S) \propto \begin{cases} n_j^- p(y_i \mid s_i = j, s_{-i}, S, \eta, y_{-i}) & j = 1, \dots, k^- \\ M \int \mathbf{N}(y_i; \alpha d_i, S) dp^o(\alpha) & j = k^- + 1 \end{cases}$$

where $y_{-i} = \{y_1, \dots, y_{i-1}, y_{i+1}, \dots, y_n\}$, $s_{-i} = \{s_1, \dots, s_{i-1}, s_{i+1}, \dots, s_I\}$,

$$n_j^- = \begin{cases} n_j - 1 & \text{if } j = s_i \\ n_j & \text{otherwise} \end{cases}$$

and $k^- = \#$ of clusters with α_i removed. If $n_{s_i}^- = 0$, we relabel the remaining clusters $j = 1, \dots, k^- = k - 1$. After sampling s_i , set

$$k = \begin{cases} k^- & \text{if } s_i \leq k^- \\ k^- + 1 & \text{if } s_i = k^- + 1 \end{cases}$$

Note that by conditional independence, in $p(y_i \mid s_i = j, s_{-i}, S, \eta, y_{-i})$ the data y_{-i} can be replaced by $\{y_l, l \in \Gamma_j, l \neq i\}$.

2. *Resampling α_j^* :*

The posterior distribution $p(\alpha_j^* \mid s, y, \eta, S)$ is

$$p(\alpha_j^* \mid s, y, \eta, S) \propto \left[\prod_{i \in \Gamma_j} \mathbf{N}(y_i; \alpha_j^* d_i, S) \right] p^o(\alpha_j^* \mid \eta).$$

With a normal base measure p^o the posterior distribution of α_j^* is normal with mean and covariance matrix that can be found by standard calculations.

3. Resampling η conditional on the currently imputed values s, α^*, k follows standard posterior simulation for DP mixture models.

To report posterior inference on H_x the MCMC simulation can be augmented as follows. Let θ denote the vector of all model parameters, and let $\theta^{(i)}$ denote the parameters imputed after i iterations of the MCMC simulation. We evaluate $E(H_x \mid Y) = p(y_{n+1} \mid x_{n+1} = x, Y)$ as

$$p(y_{n+1} \mid x_{n+1}, Y) = E[p(y_{n+1} \mid x_{n+1}, Y, \theta) \mid Y] \approx \frac{1}{T} \sum_{i=1}^T p(y_{n+1} \mid x_{n+1}, \theta^{(i)}, Y) = \frac{1}{T} \sum_{i=1}^T p(y_{n+1} \mid x_{n+1}, \theta^{(i)}).$$

The terms in the last average are easily computed. We use a superindex $^{(i)}$ to identify the imputed parameter values after i iterations of the MCMC simulation. Let d_x denote the design vector for $x = x_{n+1}$. We find

$$p(y_{n+1} | x_{n+1}, \theta^{(i)}) \propto \sum_{j=1}^{k^{(i)}} n_j^{(i)} \text{N}(y_{n+1}; \alpha_j^{*(i)} d_x, S^{(i)}) + M^{(i)} \int \text{N}(y_i; \alpha d_x, S^{(i)}) dp^o(\alpha) \quad (12)$$

Uncertainty in H_x is illustrated through posterior draws of H_x . For the following argument we consider augmenting the imputed parameter vector $\theta^{(i)}$ with the random distribution F defined in (6). Given $\theta^{(i)}$ the conditional posterior for F is a DP with updated parameters,

$$(F | \theta^{(i)}, Y) \sim \text{DP}(q, M^{(i)} + n) \text{ with } q \propto M^{(i)} p^o + \sum_{j=1}^{k^{(i)}} n_j^{(i)} \delta_{\alpha_j^{*(i)}}. \quad (13)$$

The large total mass parameter $M^{(i)} + n$ implies that the random measure F is close to the conditional expectation q , the DP base measure in (13). We exploit this to approximate a posterior draw $F \sim p(F | \theta^{(i)}, Y)$ as $F \approx q$, and thus a posterior draw for H_x as $\int \text{N}(\alpha d_x, S^{(i)}) dq(\alpha)$, i.e.,

$$H_x(y) \propto M^{(i)} \int \text{N}(y; \alpha d_x, S^{(i)}) dp^o(\alpha) + \sum_{j=1}^{k^{(i)}} n_j^{(i)} \text{N}(y; \alpha_j^{*(i)} d_x, S^{(i)})$$

The latter is simply the predictive distribution conditional on $\theta^{(i)}$ in (12).

B. Stochastic Monotonicity of k

The two results below are stated for the ANOVA DDP model that we develop in this paper. The results apply more generally, with substitution of essentially any DP or DDP model for the model in (5). The first result describes the impact of the mass parameter on the distribution of the number of clusters in prior and posterior.

Result 1. Consider a family of DDP models of the form given in (5), indexed by M . That is, the models have identical prior specifications except for the mass parameter of the DDP. Then both the prior distribution and the posterior distribution on the number of clusters is stochastically increasing in M .

Proof. The probability of obtaining k clusters from a sample of size n when the mass parameter of the Dirichlet process is M is given by the expression

$$P_M(k) = \frac{M^k \Gamma(M)}{\Gamma(M+n)} \sum \prod_{i=1}^k \Gamma(n_i)$$

where the sum runs over vectors of length k having components n_i , with $n = \sum_{i=1}^k n_i$. This expression follows from the Polya urn scheme and holds for $k = 1, \dots, n$.

Consider two values of the mass parameter, $M_1 < M_2$. Then

$$\frac{P_{M_1}(k)}{P_{M_2}(k)} = \left(\frac{M_1}{M_2}\right)^k \frac{\Gamma(M_1)}{\Gamma(M_1+n)} \frac{\Gamma(M_2+n)}{\Gamma(M_2)}$$

for $k = 1, \dots, n$. This ratio of probabilities is monotone decreasing in k , and so the prior distribution of the number of clusters for mass M_2 is stochastically greater than that for mass M_1 .

Conditional on the number of clusters, k , the distribution on the configuration vector, s , does not depend on M . Consequently, with matched priors on hyperparameters and the shape of the base measure, the likelihood for the data depends on M only through k . Thus we may write

$$\frac{P_{M_1}(k|Y)}{P_{M_2}(k|Y)} \propto \frac{P_{M_1}(Y|k)P_{M_1}(k)}{P_{M_2}(Y|k)P_{M_2}(k)} \propto \left(\frac{M_1}{M_2}\right)^k.$$

As before, monotonicity of these ratios in k leads to the conclusion that the *posterior* distribution for the number of clusters is stochastically increasing in M .

This result describes the impact of the sample size on the number of clusters. While the qualitative result follows from the Polya urn scheme (a new draw from the DDP may either join an already existing cluster, or it may begin a new cluster—it cannot remove an existing cluster), the recursion in the proof below is useful for quick computation of the prior distribution on the number of clusters. This distribution is useful for marginalizing the mass parameter through preintegration, as described in MacEachern (1998).

Result 2. Consider the model in (5). The number of clusters is stochastically increasing in n .

Proof. Let $p^{(n,M)}$ represent the vector of length n with components $P_M(k)$. $p^{(1,M)} = (1)$ for all M . To find the corresponding distribution for a sample of size $n+1$, given the distribution for a sample of size n , note that

$$p^{(n+1,M)} = \frac{n}{M+n} (p^{(n,M)}, 0) + \frac{M}{M+n} (0, p^{(n,M)}).$$

This yields the result.

References

- Albert, J. H. and Chib, S. (1993), “Bayesian Analysis of Binary and Polychotomous Response Data,” *Journal of the American Statistical Association*, 88, 669–679.
- Antoniak, C. E. (1974), “Mixtures of Dirichlet Processes with Applications to Bayesian Nonparametric Problems,” *Annals of Statistics*, 2, 1152–1174.
- Bernardo, J. M. and Smith, A. F. M. (1994), *Bayesian Theory*, New York: John Wiley & Sons.
- Blackwell, D. and MacQueen, J. B. (1973), “Ferguson distributions via Pólya urn schemes,” *The Annals of Statistics*, 1, 353–355.
- Bush, C. A. and MacEachern, S. N. (1996), “A semiparametric Bayesian model for randomised block designs,” *Biometrika*, 83, 275–285.
- Casella, G. and Robert, C. P. (1996), “Rao-Blackwellisation of sampling schemes,” *Biometrika*, 83, 81–94.
- Chen, M.-H. and Dey, D. K. (2000), “Bayesian Analysis for Correlated Ordinal Data Models,” in *Generalized Linear Models: A Bayesian Perspective*, eds. D. K. Dey, W. S. Ghosh, and B. Mallick, 135–162, New York: Marcel Dekker.
- Cifarelli, D. and Regazzini, E. (1978), “Problemi statistici non parametrici in condizioni di scambiabilità parziale e impiego di medie associative,” Technical report, Quaderni Istituto Matematica Finanziaria, Torino.
- Cowles, M. K., Carlin, B. P., and Connett, J. E. (1996), “Bayesian Tobit Modeling of Longitudinal Ordinal Clinical Trial Compliance Data With Nonignorable Missingness,” *Journal of the American Statistical Association*, 91, 86–98.
- Doss, H. (1994), “Bayesian nonparametric estimation for incomplete data via successive substitution sampling,” *The Annals of Statistics*, 22, 1763–1786.
- Escobar, M. D. (1994), “Estimating normal means with a Dirichlet process prior,” *Journal of the American Statistical Association*, 89, 268–277.
- Escobar, M. D. and West, M. (1998), “Computing nonparametric hierarchical models,” in *Practical Nonparametric and Semiparametric Bayesian Statistics*, eds. D. Dey, P. Müller, and D. Sinha, 1–22, New York, NY, USA.
- Ferguson, T. S. (1973), “A Bayesian analysis of some nonparametric problems,” *Annals of Statistics*, 1, 209–230.
- Freedman, R., et al. (2002), “Pilot Study of FLT3 Ligand in Patients with Peritoneal Carcinomatosis Associated with Müllerian and Gastrointestinal Carcinomas,” Technical report, University of Texas, M.D. Anderson Cancer Center.
- Gelfand, A. E. and Kottas, A. (2001), “Nonparametric Bayesian Modeling for Stochastic Order,” *Annals of the Institute of Statistical Mathematics*, 53, 865–876.

- (2002), “A Computational Approach for Full Nonparametric Bayesian Inference under Dirichlet Process Mixture Models,” *Journal of Computational and Graphical Statistics*, 11, 289–305.
- Gelfand, A. E. and Kuo, L. (1991), “Nonparametric Bayesian bioassay including ordered polytomous response,” *Biometrika*, 78, 657–666.
- Giudici, P., Mezzetti, M., and Muliere, P. (2002), “Mixtures of Dirichlet process priors for variable selection in survival analysis,” *Journal of statistical planning and inference*, to appear.
- Green, P. and Richardson, S. (2001), “Modelling heterogeneity with and without the Dirichlet process,” *Scandinavian Journal of Statistics*, 28, 355–375.
- Guglielmi, A. and Tweedie, R. L. (2001), “Markov chain Monte Carlo estimation of the law of the mean of a Dirichlet process,” *Bernoulli*, 7, 573–592.
- Ishwaran, H. and James, L. (2001), “Gibbs sampling methods for stick-breaking priors,” *Journal of the American Statistical Association*, 96.
- Kottas, A., Müller, P., and Quintana, F. (2002), “Nonparametric Bayesian modeling for multivariate ordinal data,” Technical report, University of California Santa Cruz.
- Kuo, L. and Mallick, B. (1997), “Bayesian semiparametric inference for the accelerated failure-time model,” *The Canadian Journal of Statistics*, 25, 457–472.
- Liu, J. S. (1994), “The collapsed Gibbs sampler in Bayesian computations with applications to a gene regulation problem,” *Journal of the American Statistical Association*, 89, 958–966.
- MacEachern, S. N. (1994), “Estimating normal means with a conjugate style Dirichlet process prior,” *Communications in Statistics, Part B – Simulation and Computation*, 23, 727–741.
- MacEachern, S. N. (1998), “Computational Methods for Mixture of Dirichlet Process Models,” in *Practical Nonparametric and Semiparametric Bayesian Statistics*, eds. D. Dey, P. Müller, and D. Sinha, 23–44, New York, NY, USA.
- MacEachern, S. N. (1999), “Dependent nonparametric processes,” in *ASA Proceedings of the Section on Bayesian Statistical Science*, American Statistical Association, Alexandria, VA.
- (2001), “Decision theoretic aspects of dependent nonparametric processes,” in *Bayesian Methods with Applications to Science, Policy and Official Statistics*, ed. E. George, 551–560, ISBA.
- MacEachern, S. N. and Müller, P. (1998), “Estimating mixture of Dirichlet process models,” *Journal of Computational and Graphical Statistics*, 7, 223–239.
- (2000), “Efficient MCMC Schemes for Robust Model Extensions using Encompassing Dirichlet Process Mixture Models,” in *Robust Bayesian Analysis*, eds. F. Ruggeri and D. R. Insua, New York, NY, USA.

- Melamed, M. R., Lindmo, T., and Mendelsohn, M. L., eds. (1990), *Flow cytometry and sorting*, Wiley-Liss.
- Mira, A. and Petrone, S. (1996), “Bayesian Hierarchical Nonparametric Inference for Change-Point Problems,” in *Bayesian Statistics 5*, eds. J. M. Bernardo, J. O. Berger, A. P. Dawid, and A. F. M. Smith, Oxford: Oxford University Press.
- Muliere, P. and Petrone, S. (1993), “A Bayesian Predictive Approach to Sequential Search for An Optimal Dose: Parametric and Nonparametric Models,” *Journal of the Italian Statistical Society*, 2, 349–364.
- Muliere, P. and Tardella, L. (1998), “Approximating distributions of random functionals of Ferguson-Dirichlet priors,” *Canadian Journal of Statistics*, 26, 283–297.
- Müller, P. and Rosner, G. (1997), “A Bayesian population model with hierarchical mixture priors applied to blood count data,” *Journal of the American Statistical Association*, 92, 1279–1292.
- Regazzini, E., et al. (2002), “Theory and numerical analysis for exact distributions of functionals of a Dirichlet process,” *AnnStat*, 30, 1376–1411.
- Sethuraman, J. (1994), “A constructive definition of the Dirichlet process prior,” *Statistica Sinica*, 2, 639–650.
- Tomlinson, G. and Escobar, M. (1999), “Analysis of densities,” Technical report, University of Toronto.