

PANEL PROBIT WITH FLEXIBLE CORRELATED EFFECTS: QUANTIFYING TECHNOLOGY SPILLOVERS IN THE PRESENCE OF LATENT HETEROGENEITY

MARTIN BURDA^{a,b} AND MATTHEW HARDING^{c*}

^a *Department of Economics, University of Toronto, Ontario, Canada*

^b *IES, Charles University, Prague, Czech Republic*

^c *Department of Economics, Stanford University, CA, USA*

SUMMARY

In this paper, we introduce a Bayesian panel probit model with two flexible latent effects: first, unobserved individual heterogeneity that is allowed to vary in the population according to a nonparametric distribution; and second, a latent serially correlated common error component. In doing so, we extend the approach developed in Albert and Chib (*Journal of the American Statistical Association* 1993; **88**: 669–679; in *Bayesian Biostatistics*, Berry DA, Stangl DK (eds), Marcel Dekker: New York, 1996), and in Chib and Carlin (*Statistics and Computing* 1999; **9**: 17–26) by releasing restrictive parametric assumptions on the latent individual effect and eliminating potential spurious state dependence with latent time effects. The model is found to outperform more traditional approaches in an extensive series of Monte Carlo simulations. We then apply the model to the estimation of a patent equation using firm-level data on research and development (R&D). We find a strong effect of technology spillovers on R&D but little evidence of product market spillovers, consistent with economic theory. The distribution of latent firm effects is found to have a multimodal structure featuring within-industry firm clustering. Copyright © 2012 John Wiley & Sons, Ltd.

Received 11 May 2010; Revised 5 March 2012



Supporting information can be found in the online version of this article.

1. INTRODUCTION

There is broad agreement that individual heterogeneity plays a crucial role in many economic models. In linear models, panel data can be used to identify the effects of interest while at the same time controlling for unobserved individual heterogeneity (Hausman and Taylor, 1981). Nonlinear models with unobserved heterogeneity pose substantial theoretical and computational challenges (Arellano and Hahn, 2006). In particular, in the case of nonlinear panel data models it is in general not possible to remove the unobserved effects by differencing as is commonly done in linear models. Convenient solutions can be obtained in some cases when a specific parametric form is assumed for the distribution of heterogeneity, such as in the negative binomial regression. Nonetheless, relaxing parametric assumptions on the distribution of unobserved heterogeneity in nonlinear models is important, as often such restrictions cannot be justified by economic theory.

One possibility is to treat the unobserved effects as nuisance parameters to be estimated along with the parameters of interest. This approach requires large amounts of data though, as consistency is guaranteed only in the large N and large T limit. In most microeconomic applications, the econometrician only has a small number of repeated cross-sections to work with and the estimation of the individual fixed effects as incidental parameters induces bias. In the logit case Abrevaya (1999) shows that the model with fixed effects and only two time periods leads to severe bias and the estimated coefficients can reach twice their true value. In the parametric setting it is possible in some cases to circumnavigate this problem by

* Correspondence to: Matthew Harding, Department of Economics, Stanford University, 579 Serra Mall, Stanford, CA 94305, USA. E-mail: mch@stanford.edu

redefining the quantity of interest. Fernandez-Val (2007) shows that under certain assumptions the inclusion of fixed effects does not affect the consistency of the marginal effects. More recently, Arellano and Bonhomme (2009) show that well-chosen weights in average or integrated likelihood settings can produce estimators that are first-order unbiased.

Removing the parametric assumptions on the distribution of unobserved heterogeneity is also beneficial since economic models are usually silent on how to formally describe individual heterogeneity. At the same time, recent attempts at estimating nonlinear models nonparametrically are often rather difficult to implement (Berry and Haile, 2009).

Fueled by advances in computation, as well as their flexibility and conceptual simplicity, Bayesian methods provide a powerful alternative to the more traditional approaches to solving these problems. In particular, Bayesian hierarchical models can be readily extended to incorporate inference on latent classes of similar individuals or mixtures of distributions for various objects of interest. This makes Bayesian modeling an extremely flexible tool and a promising avenue to explore relaxing the assumptions discussed.

In some special cases such as the probit model, Bayesian data augmentation completely avoids the need to specify the likelihood in the form of a multivariate integral. This feature was introduced for the probit model in a seminal paper by Albert and Chib (1993). Instead of formulating the likelihood by integrating out the latent utility, the estimation problem is recast in the form of an iterative scheme of linear regressions where the latent utility is explicitly sampled along with other model parameters. Thus the estimation is free from the curse of dimensionality that plagues inference with integral-based likelihoods. The approach was further developed for limited dependent variable (LDV) models to include parametric random latent effects in Albert and Chib (1996), Chib and Carlin (1999), and Gu *et al.* (2009).

In this paper, we further extend this line of research by introducing a model with two latent variables: first, we introduce unobserved individual heterogeneity that is allowed to vary in the population according to a nonparametric distribution; and second, a latent error component that is serially correlated over time. The unobserved individual effects are allowed to be correlated with the observed regressors, in the spirit of Chamberlain (1982, 1984). Our model thus extends beyond the class of traditional random effects models (for a discussion on this issue see, for example, Wooldridge, 2001). We model the distribution of the unobserved heterogeneity component with a nonparametric Dirichlet Process (DP) mixture model. The prior for the latent time component is specified as a parametric autoregressive process but its influence decreases linearly with the amount of data available. Due to its structure we label the proposed model as the 'flexible latent effects probit' (FLEP). We note that individual building blocks of our model have been used in separate settings, such as modeling autoregressive processes in discrete-choice models (Allenby and Lenk, 1994) and implementing the DP prior for studying heterogeneity in choice models (Li and Zheng, 2008; Rossi, 2010). However, the combined model with panel latent effects considered here has not yet been applied in the literature. Our aim in this paper is to show how to account for both flexible forms of unobserved heterogeneity and common latent time effects within the same framework.

We conduct an extensive empirical analysis of the decision to innovate where we suspect that unobserved heterogeneity plays an important role at the firm level. At the same time, patenting activity may also be driven by a common time trend reflecting the macroeconomic environment or the overall stock of scientific knowledge. Without properly accounting for these latent effects it is not possible to correctly identify effects of interest or test hypotheses based on economic theory. We use data from a recent study of firm-level research and development (R&D) by Bloom *et al.* (2010) (henceforth BSV) to estimate a patent equation and test theoretical predictions on R&D spillovers. The dataset captures the majority of the patents granted between 1980 and 2001 in the USA.

We explore the possibility that R&D leads to two major externalities. On the one hand, R&D may increase the productivity of firms using similar technology whereby a firm can benefit from the R&D conducted by another firm in the same technology area. On the other hand, R&D can have a product

market rivalry effect, with a number of firms striving to develop essentially the same product, which is detrimental to social welfare. Economic theory predicts that the marginal effect of technology spillovers on patenting activity is positive, while the marginal effect of product market spillovers on patenting activity is zero. Our econometric approach allows us to additionally account for the two important types of latent effects in the analysis of R&D spillovers mentioned above: firm-level heterogeneity and common time factors. We document the presence of both statistically significant technology spillover effects and firm-level heterogeneity. The estimated distribution of firm-level heterogeneity shows many interesting features and its multimodality suggests the clustering of heterogeneity across different firms. One important advantage of our approach is that it estimates firm-level clustering without having to rely on a priori guesses of the form of heterogeneity. As we shall see, industry classifications, a common proxy for heterogeneity, does a poor job at capturing the measured variation in latent firm-level heterogeneity.

Our paper also introduces a series of computational innovations for the Bayesian estimation of this class of models. A core component in our implementation strategy is the efficient computation of the posteriors using a recent Sequentially Allocated Merge–Split (SAMS) algorithm (Dahl, 2005) that is substantially more efficient than samplers used previously in similar contexts. The SAMS sampler can update in one move large blocks of elements involved in implementation of the DP sampling scheme. It thus avoids a shortcoming of sequential samplers, such as the Polya urn scheme, that can get stuck in particular clustering configurations due to the one-at-a-time nature of their updates. Moreover, the SAMS algorithm is applicable to both conjugate and non-conjugate DP mixture models.

Our approach builds on the stream of literature aiming to relax restrictive assumptions of existing limited dependent variable models. A recent state-of-the-art Bayesian nonparametric analysis was introduced in Chib and Jeliazkov (2006), who study a binary dependent variable model with AR(p) errors and normally distributed unobserved individual heterogeneity. These authors focus on a nonparametric estimation of an unknown function of the model covariates, while we model nonparametrically the distribution of the unobserved individual heterogeneity. Burda *et al.* (2008) analyze a flexible model for multinomial discrete choice with a flexible distribution of several parameters on the observable regressors. Their unobserved error component was fully parametric with an extreme-value type 1 distribution. As a result, their model was based on the logit closed-form solution facilitated by such assumption. Moreover, their model did not incorporate any dynamic element. In contrast, the error component in our model contains both flexible unobserved individual heterogeneity and common latent time effects, which makes our estimation method suitable for panel data with a dynamic latent factor structure. The Normal distribution of the transitory idiosyncratic component stipulates a probit structure here, precluding the closed-form logit likelihood derivation utilized in Burda *et al.* (2008). Instead, here we rely on data augmentation due to Albert and Chib (1993) using an iterative scheme of linear regressions in sampling the latent utility along with other model parameters.

Random error components that induce correlation over alternatives and time can also be accommodated by frequentist procedures. Such an approach would assume a model for the distribution of the latent components and then specify the model likelihood in the form of an integral whose dimensions are formed by the individual unobserved components. Typically, such an integral is analytically intractable and hence is estimated by numerical simulation methods, such as the GHK simulator developed by Geweke (1991), Hajivassiliou (1990), and Keane (1990). The resulting simulated likelihood (SML) is then maximized with respect to the model parameters. The GHK procedure thus numerically approximates the likelihood integral until convergence at every iteration of the model parameters within the optimization procedure. In contrast, Bayesian Gibbs sampling factorizes the high-dimensional multivariate integral into a sequence of low-dimensional conditional density kernels, drawing one dimension at a time until a single convergence state of the resulting Markov chain is attained. In many cases, this implies that Bayesian parameter estimation is substantially faster than SML. For example, in an empirical comparison study for a parametric multinomial probit model Bolduc *et al.* (1997) found the Bayesian approach about twice as fast and much simpler to implement, both conceptually and computationally, than the GHK method.

Moreover, the Bayesian Markov chain of parameter draws can be directly used for inference in analogy to a bootstrap sample. In contrast, frequentist SML procedures including GHK require additional estimation of the shape of the simulated likelihood around the argmax parameter value; this process is fraught with peril as integral likelihoods often suffer from multiple local modes or saddles (Knittel and Metaxoglou, 2008). Dealing with such features is avoided using the Bayesian approach. In a comparison study between a Bayesian approach and the frequentist SML approach for a class of parametric mixed logit models, Train (2001) finds the Bayesian approach to possess theoretical advantages from both a classical and Bayesian perspective. Additional benefits of Bayesian inference in latent variable models are discussed, for example, in Paap (2002).

The advantages of Bayesian methods become even more pronounced with increased dimensionality of the underlying problem. A nonparametric model for the distribution of unobserved heterogeneity, as considered in this paper, if estimated using the GHK approach, would necessitate maximization of a flexible functional form such as a series or kernel estimator involving a large number of parameter iterations. The high-dimensional likelihood integral would need to be numerically approximated to a sufficient degree of precision at each of these iterations, which may become computationally prohibitive for larger sample sizes. In contrast, the Bayesian conditional Gibbs sampling can be performed very accurately along each latent dimension whereby higher dimensionality of the problem does not diminish the precision of inference.

The remainder of the paper is organized as follows. Section 2 introduces our model and discusses the assumptions and sampling procedures. Section 3 presents an application of the method to the estimation of the effect of technological spillovers and product market competition on innovation. Section 4 concludes. A series of Monte Carlo studies comparing the performance of our method with other existing approaches is presented in an online supplement to this article as supporting information.

2. MODEL

Consider a sample of binary responses y_{it} , for N individuals indexed by i , and T time periods indexed by t . We assume that the data are drawn from the following error-components model:

$$\tilde{y}_{it} = \mathbf{x}_{it}\beta + u_{it} \quad (1)$$

$$u_{it} = \tau_i + \lambda_t + \varepsilon_{it}$$

$$y_{it} = 1(\tilde{y}_{it} \geq 0) \quad (2)$$

where \mathbf{x}_{it} is a $(1 \times K)$ vector of explanatory variables, τ_i represents unobserved individual heterogeneity, λ_t captures latent time effects, and $1(\mathcal{C})$ denotes the indicator function, which takes the value one if the condition \mathcal{C} is satisfied and zero otherwise. The term \tilde{y}_{it} can be thought of as a latent utility of individual i at time t . In this error-components model the unobserved error u_{it} is decomposed into three parts: an individual specific error τ_i , a time-specific component λ_t and an idiosyncratic and transitory shock ε_{it} . This structure of u_{it} allows for both the presence of individual heterogeneity and serial correlation in the residual while these components can still be separately identified. In this model we observe the covariates x_{it} , but not τ_i , λ_t or ε_{it} . The model is specified in terms of the latent variable \tilde{y}_{it} , not observed by the econometrician, who only observes the binary outcome variable y_{it} .

Let $\tilde{\mathbf{y}}_i = (\tilde{y}_{i1}, \dots, \tilde{y}_{iT})'$, $\tilde{\mathbf{y}} = (\tilde{\mathbf{y}}'_1, \dots, \tilde{\mathbf{y}}'_N)'$, $\mathbf{X}_i = (x'_{i1}, \dots, x'_{iT})'$, and $\mathbf{X} = (\mathbf{X}'_1, \dots, \mathbf{X}'_N)'$, $\lambda = (\lambda_1, \dots, \lambda_T)'$, $\tau = (\tau_1, \dots, \tau_N)'$, $\varepsilon_i = (\varepsilon_{i1}, \dots, \varepsilon_{iT})'$, $\varepsilon = (\varepsilon'_1, \dots, \varepsilon'_T)'$ and let $\mathbf{1}$ denote a $(T \times 1)$ vector of ones. Model (1) can thus be rewritten more compactly as $\tilde{\mathbf{y}}_i = \mathbf{X}_i\beta + \tau_i\mathbf{1} + \lambda + \varepsilon_i$ for $i = 1, \dots, N$ or simply $\tilde{\mathbf{y}} = \mathbf{X}\beta + \tau + \lambda + \varepsilon$. Following the notation in Geweke (2005), let the set-valued function

$C_{it} = c_{it}(\tilde{y}_{it})$ with $C_{it} = (-\infty, 0]$ if $y_{it} = 0$ and $C_{it} = (0, \infty)$ if $y_{it} = 1$. Denote the collection $\mathbf{C} = \{C_{it} : i = 1, \dots, N; t = 1, \dots, T\}$.

The hierarchical structure of our model allows us to distinguish four different layers of parameters. The first layer corresponds to the structure of the error components τ_i , λ_t and ε_{it} . Its properties are given by the following Assumption.

Assumption 1. The error components τ_i , λ_t , and ε_{it} , for $i = 1, \dots, N$ and $t = 1, \dots, T$, are mutually independent conditionally on the \mathbf{X} .

The second parameter layer characterizes the distributional properties of the first-layer parameters in Assumptions 2–4.

Assumption 2. The variables τ_1, \dots, τ_N are independent, identically distributed:

$$\tau_i \sim F_{\tau}^0$$

where F_{τ}^0 is a continuous unknown distribution, conditionally on X_i and the other model parameters of primary economic interest.

Instead of imposing a parametric family model, F_{τ}^0 will be estimated as an infinite mixture of distributions using a Bayesian Dirichlet Process Mixture (DPM) model which we shall introduce below. Moreover, our sampling mechanism allows for joint posterior correlation of τ_i with other regressors. We do not impose any prior assumptions on this feature explicitly due to the absence of any initial information on this property. Since τ_i is sampled conditional on X_i , such potential relationship is entirely data-driven. The next assumption specifies the prior distribution for the latent time effects.

Assumption 3. λ_t is assumed to follow a stationary Gaussian autoregressive process:

$$\lambda_t = \rho_1 \lambda_{t-1} + \dots + \rho_s \lambda_{t-s} + \eta_t$$

with $\eta_t \sim N(0, \sigma_{\eta}^2)$. Furthermore, η_t is independent of ε_{it} for each $t = 1, \dots, T$.

Failure to account for serial correlation of the error term has potential negative consequences. In the Bayesian framework, the posterior distribution is a weighted average of the prior distribution and the parameter update learned from the data via the likelihood function. The latter is implied by the probit model specification (1). In our sampling scheme detailed below, the prior has weight $1/(T+1)$ while the likelihood information has weight $T/(T+1)$. In samples with very small T this autoregressive specification for the prior impacts inference but the prior influence declines linearly with T . The autoregressive prior specification also facilitates learning about the posterior distribution of the latent time process hyperparameters ρ and σ_{η} that provide information about the nature of the persistence and volatility in the latent time error component.

The parametric nature of Assumption 3 renders the dynamic model specification potentially quite restrictive, especially in cases where the data-generating process follows some other form of dynamics. Nonetheless, panels of data in micro-econometric applications are typically characterized by large N and small T dimensions and hence a parametric model appears as a suitable way to capture the relatively limited amount of information conveyed by the time dimension. Conversely, the relatively rich informational content of the cross-sectional dimension lends itself to non-parametric modeling, which we undertake in this paper.

Assumption 3 is stated conditional on a given lag order s . Model selection criteria can be further employed to determine the optimal lag order for a given dataset. A method of lag selection for the autoregressive model is discussed in Troughton and Godsill (1997).

The following assumption defines the probit structure of the model.

Assumption 4. $\varepsilon_{it} \sim N(0, 1)$ is a stochastic error component uncorrelated with any other regressor.

Our proposed model builds on the traditional error-components framework due to its popularity in applied work. The random error to an observation, $u_{it} = \tau_i + \lambda_t + \varepsilon_{it}$, is given by the sum of an individual effect τ_i , a time effect λ_t and an idiosyncratic shock ε_{it} . Variations on this framework can be readily incorporated into our model.

The third parameter layer in our model is formed by parameters of primary economic interest captured in the vector $\theta = (\beta', \sigma_\eta, \rho')$. The assumptions on the prior distributions for this layer are specified as follows.

Assumption 5.

$$\beta \sim N(\underline{\beta}, \underline{\Sigma}_\beta) \quad (3)$$

$$\sigma_\eta^2 \sim IG(v_0, s_0) \quad (4)$$

$$\rho \sim Uniform(\Omega) \quad (5)$$

where $\Omega \subseteq R^s$ is the stationarity region of the autoregressive process.

The fourth parameter layer is comprised of the remaining hyperparameters introduced in Assumptions 2–5. In order to fully characterize this layer, we will elaborate on the model specified for the distribution of the unobserved heterogeneity component. implies the following model based on Neal (2000):

$$\tau_i | \psi_i \sim F_\tau(\psi_i) \quad (6)$$

$$\psi_i | G \sim G \quad (7)$$

$$G \sim DP(\alpha, G_0) \quad (8)$$

Thus F_τ is specified as an infinite mixture of distributions $F_\tau(\psi)$, with the mixing distribution over ψ being G . Here, ψ_i are hyperparameters of the distribution $F_\tau(\psi_i)$ of τ_i drawn from a random probability measure G , which itself is distributed according to a DP prior. The DP prior for G is indexed by two hyperparameters: a distribution G_0 that defines the ‘location’ of the DP prior and a positive scalar precision parameter α . The distribution G_0 may be viewed as a baseline prior that would be used in a typical parametric analysis. The flexibility of the DP prior model environment stems from allowing G to stochastically deviate from G_0 . The precision parameter α determines the concentration of the prior for G around the DP prior location G_0 and thus measures the strength of belief in G_0 . For large values of α , a sampled G is very likely to be close to G_0 , and vice versa. Early important applications of the DP prior to economics were made in Chib and Hamilton (2002) and Hirano (2002).

By Assumption 2 the distribution F_τ is sampled conditional on the primary parameters of economic interest θ and on the regressors \mathbf{X} . This sampling framework gives us the flexibility to treat τ_i as nuisance parameters while at the same time allowing for the possibility of the individual effects being correlated with other right-hand-side variables. Following Arellano and Bonhomme (2009) we implicitly assume that the support of F_τ contains an open neighborhood of the true parameters θ .

The fourth parameter layer is thus formed by the hyperparameters $\{\psi_i\}_{i=1}^N, G, \alpha, G_0, \underline{\beta}, \underline{\Sigma}_\beta, v_0,$ and s_0 . In our implementation, $G_0, \underline{\beta}, \underline{\Sigma}_\beta, v_0,$ and s_0 are fixed, $\{\psi_i\}_{i=1}^N,$ and α are sampled, while bypassing explicit sampling of G .

Let $\tau = \{\tau_i\}_{i=1}^N, \mu_{it} = \mathbf{x}_{it}\beta + \tau_i + \lambda_t,$ and denote by $\Phi(\mu_{it})$ and $\phi(\mu_{it})$ the cdf and pdf of the Normal random variable with unity variance, respectively. Denote generically by $p(\cdot)$ a probability density or mass function and by $k(\cdot)$ a prior density function. The posterior of our model can then be expressed as

$$p(\tilde{y}, \tau, \beta, \lambda, \sigma_\eta^2, \rho | \mathbf{y}) \propto p(\mathbf{y} | \tilde{y}, \tau, \beta, \lambda, \sigma_\eta^2, \rho, \psi, \alpha) p(\tilde{y} | \tau, \beta, \lambda, \sigma_\eta^2, \psi, \alpha) \times k(\psi | \alpha) k(\alpha) k(\beta) k(\lambda) k(\sigma_\eta^2) k(\rho) \tag{9}$$

with $k(\rho), k(\beta),$ and $k(\sigma_\eta^2)$ given in Assumption 5, $k(\lambda)$ in Section 5.5 in the Appendix, $k(\alpha)$ specified as in Escobar and West (1995), and $k(\psi | \alpha)$ given by (7–8). The remainder of the model is formulated similarly to Albert and Chib (1993) with the single index given by μ_{it} . Specifically, $p(\tilde{y} | \tau, \beta, \lambda, \sigma_\eta^2, \rho, \psi, \alpha) = \prod_i \prod_t \phi(\mu_{it})$ and $p(\mathbf{y} | \tilde{y}, \tau, \beta, \lambda, \sigma_\eta^2, \rho, \psi, \alpha)$ assigns probability mass one to $y_{it} = 1$ if $\tilde{y}_{it} > 0$ and to $y_{it} = 0$ if $\tilde{y}_{it} \leq 0$. Thus y_{it} are independent Bernoulli random variables with $p_{it} = \Phi(\mu_{it})$.

2.1. Average Partial Effects

In nonlinear models the estimated coefficients are only of limited interest by themselves. Instead, the average partial effects (APEs) are particularly useful for computing economic counterfactuals and are widely used in applied work. In this section we describe how they are computed within the setup of our model. We utilize the classical concept of the APEs augmented with the latent variables. Let

$$m_{itk} = \frac{\partial E[y_{it} | \mathbf{x}_{it}\beta, \tau_i, \lambda_t]}{\partial \mathbf{x}_{itk}} = \phi(\mathbf{x}_{it}\beta + \tau_i + \lambda_t) \beta_k$$

denote the marginal effect of a change in \mathbf{x}_k , where $\phi(\cdot)$ denotes the standard normal density function. Define

$$\tilde{\gamma} = \frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T \phi(\mathbf{x}_{it}\beta + \tau_i + \lambda_t) \tag{10}$$

The APE of \mathbf{x}_k on \mathbf{y} is then given by

$$\frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T m_{itk} = \tilde{\gamma} \beta_k \tag{11}$$

We sample explicitly τ_i and λ_t throughout the MC iterations and hence can compute the APEs directly from the definition of $\tilde{\gamma}$, as

$$\gamma = \frac{1}{NTS} \sum_{i=1}^N \sum_{t=1}^T \sum_{s=1}^S \phi(\mathbf{x}_{it}\beta_s + \tau_{is} + \lambda_{ts}) \quad (12)$$

where s is the index over MC steps.

In both the application and the simulation study (included in supporting information), we report the mean bias and means squared error of the estimated ‘APE scale’ coefficient γ defined in (12). To obtain the APEs, γ is simply multiplied by each respective β_k .

3. INNOVATION AND R&D SPILLOVERS

3.1. The Role of Latent Effects in R&D Analysis

An ongoing puzzle in the economic literature on R&D concerns the relationship between innovation as measured by the patenting activity of firms and the spillover effects resulting from the strategic interactions between firms. Firms often interact in geographically delimited markets, which leads to a localization of the spillover effects in terms of geographic distance (Griffith *et al.*, 2011). At the same time firms interact in more abstract spaces, such as the technology space defined by the extent to which two firms are close to each other in terms of the underlying technology, and the product market space defined by the extent to which two firms compete for the same product market (Bloom *et al.*, 2010). Such spillover effects often have contradictory impacts on firm performance. While technological spillover effects may benefit a firm by enhancing the overall stock of knowledge to the firm, product market spillovers can lead to business stealing due to overlapping product offerings to consumers. These issues have been explored in the theoretical literature but have been very difficult to estimate empirically due to the presence of confounding latent effects which are particularly problematic in this setting.

On the one hand, innovation and the patenting activity of firms are likely to be influenced by unobserved firm-level heterogeneity. Firm-specific differences in corporate culture, investment strategies, know-how, or brand name will arguably shape the different degrees of intensity of the innovation activity in firms. Griffith *et al.* (2011) show that ignoring unobserved heterogeneity can have a large quantitative impact on our understanding of innovative activity. On the other hand, the econometric analysis of spillover effects is confounded by the presence of systematic common time effects reflecting the macroeconomic environment, technological trends, or global political events that also affect the resource allocation in firms’ R&D funding. When such common time effects dominate it is easy to falsely attribute the observed correlation in innovative activity to spillover effects.

The econometric analysis of innovation thus presents the econometrician with an important challenge in terms of consistently estimating the effect of spillover effects between firms while accounting for the presence of unobserved individual and time effects. The Bayesian model introduced in this paper presents a useful approach to the consistent estimation of spillover effects while accounting for the presence of latent individual and time effects. As we will show, our proposed model is not only computationally feasible to implement on large datasets but it is also superior to more traditional frequentist approaches in terms of its ability to correctly predict the incidence of innovative activity in firms.

In particular, we apply the method developed in this paper to the estimation of a patent equation on firm-level data and test theoretical predictions on R&D spillovers. Since no economic theory is available that would recommend a particular distributional form for the unobserved heterogeneity we can take advantage of an important feature of our model, namely the ability to specify the distribution of the unobserved individual heterogeneity component nonparametrically. As we will show, this expands the applicability of the analysis substantially since it allows us to investigate the presence of clustering in the unobserved effects and derive new economic insights into the innovative activities of firms in different industries. At the same time, our model is flexible enough to account for the

presence of potentially confounding time factors which, if not properly accounted, will induce spurious correlations in innovative activity not attributable to the spillover effects under consideration.

3.2. Data

We employ data from a recent study of firm-level R&D by Bloom *et al.* (2010) (denoted by BSV for the rest of this section). BSV collected firm-level accounting data, such as sales, from the US Compustat database. These data were then matched to the NBER US Patent and Trademark Office data containing detailed information on granted US patents, yielding an unbalanced panel of 729 firms with observations recorded between 1980 and 2001.

BSV investigate two major spillover effects of R&D: technological and product market spillovers. On the one hand, R&D may increase the productivity of firms using similar technology. A firm can benefit from the R&D conducted by another firm in the same technology area. On the other hand, it can have a product market rivalry effect, which is detrimental to social welfare. Using the firm-level information available, BSV attempt to map the location of each firm in both the technology and product space, by comparing information on patents and information on sales across firms.

Following BSV we measure the technological closeness between firms using information on patents for each firm. All available patents are allocated to $\kappa = 1, \dots, 425$ different technological classes. If we then let $T_i = (T_{i\kappa})$ denote a vector where each element represents the average share of patents of firm i in technological class κ over the period 1980–2001, we can define technological closeness (*Tech*) between two firms i and j by the uncentered correlation between the allocations for the two firms:

$$Tech_{i,j} = \frac{T_i T_j'}{(T_i T_i')^{1/2} (T_j T_j')^{1/2}} \quad (13)$$

The degree of technology spillover *SpillTech* is then measured as the technology distance weighted average of the R&D stock of all other firms at each point in time:

$$SpillTech_{i,t} = \sum_{j,i \neq j} Tech_{i,j} R_{j,t} \quad (14)$$

where $R_{j,t}$ is the stock of R&D of firm j at time t computed from the expenditure on R&D data available in the accounting statements recorded by US Compustat.

Similarly, the distance between firms in the product market can be computed by decomposing each firm's sales by the respective four-digit industry code. Most firms are multi-product firms with reported sales in an average of 5.2 different industry codes. The sample of firms spans a total of 762 different industries. The distance between firms in the product market is then measured as the uncentered correlation between the allocation of sales activity of firms into industries. The degree of product market spillovers (*SpillSIC*) is computed as the product market distance weighted average of the R&D stock of all other firms.

The above definition of technology and product market spillovers are based on the Jaffe (1986) distance measure. BSV note as its drawback the implicit assumption of technology spillovers only occurring within the same product technology class. Patent class categorizations, however, are extremely narrow. As BSV illustrate, the Patent Office distinguishes between 'arithmetic processing and calculating' and 'processing architectures and instruction processing' when they both may refer to very similar computer technology. Moreover, categorization into patent classes may well be subject to measurement error. As such, it is worthwhile to investigate additional distance measures which take into account the fact that technological spillovers may occur across patent classes. One such option is

the Mahalanobis distance, which allows spillovers to occur across multiple patent classes but weighs their importance by the extent to which a firm is active across different patent classes. A Mahalanobis measure of the product market spillovers is constructed similarly. We enrich our analysis by adding a Mahalanobis version of the *SpillTech* and *SpillTech* variables, which will serve as a subsequent robustness check on our baseline specifications and guard against measurement errors resulting from the more narrow Jaffe variable construction.

The dataset contains two additional variables of interest. The first is the R&D stock, which has already been mentioned above. The second is a firm-specific measure of industry sales (*Sales*). This variable uses the same SIC weighting technique as *SpillSIC* but applied to rival firm sales.

The dependent variable of interest *Patenting* is a binary variable denoting whether or not firm i filed at least one patent in year t . The data summary statistics are given in Table 1. All independent variables are expressed in logarithms and have been lagged by one period to remove simultaneity concerns.

3.3. Econometric Implementation

We implement the model developed in Section 2 to the estimation of R&D spillover effects on patenting activity using the data described above. We use a Bayesian Gibbs sampling scheme (the precise implementation details of drawing from individual Gibbs blocks are given in the Appendix). Under the Model (1) and Assumptions 1-5, the joint posterior density can be decomposed into the following Gibbs blocks:

1. $\beta | \tau, \lambda, \psi, \theta / \beta, \tilde{y}, \mathbf{y}, \mathbf{X}$.
2. $\tilde{y} | \tau, \lambda, \psi, \theta, \mathbf{y}, \mathbf{X}$.
3. Update the assignments of τ_i to latent classes by alternating between the SAMS (Dahl, 2005) and Algorithm 7 (Neal, 2000), which includes sampling $\{\psi_i\}_{i=1}^N$.
4. $\tau_i | \psi, \theta, \lambda, \tilde{y}, \mathbf{y}, \mathbf{X}$ for each i .
5. $\lambda | \tau, \psi, \theta, \tilde{y}, \mathbf{y}, \mathbf{X}$.
6. $\sigma_\eta^2 | \tau, \lambda, \psi, \theta / \sigma_\eta^2, \tilde{y}, \mathbf{y}, \mathbf{X}$.
7. $\rho | \tau, \lambda, \psi, \theta / \rho, \tilde{y}, \mathbf{y}, \mathbf{X}$.

Table I. Summary statistics. All variables are in logarithms and lagged by one period

Variable	Mean	SD	Min.	Max.
<i>SpillTech</i> (Jaffe)	9.554	1.142	4.838	11.707
<i>SpillSIC</i> (Jaffe)	7.272	2.323	-4.602	11.154
<i>SpillTech</i> (Mah.)	11.314	0.847	8.235	13.156
<i>SpillSIC</i> (Mah.)	8.513	1.660	-0.356	11.559
<i>R&D Stock</i>	3.030	3.026	-2.513	10.765
<i>Sales</i>	6.230	1.962	0	12.103
<i>Patenting</i>	0.540	0.498	0	1
Firms	729			
Total obs.	12,928			

The Bayesian model described above contains a non-parametric specification of the individual effects and we will denote it as FLEP (flexible latent effects probit). It is possible to estimate a restricted parametric version of the same model by imposing the condition that in the unobserved individual effects are Normally distributed. We shall label this version of the model as PLEP (parametric latent effects probit). By comparing the results of different specifications of the unrestricted nonparametric version with the restricted parametric version of the same model we can gain additional insights into the importance and advantages of using a flexible nonparametric specification over more traditional parametric approaches. Posterior means are reported for these two techniques, obtained from chains of total length of 10,000 MC steps with a 5000 burn-in section.

Additionally we implement two frequentist approaches to the estimation of spillover effects. First, we implement the fixed-effects probit model with time dummies (denoted by FE). As we shall see, this approach suffers from serious computational limitations in large data. Second, we implement the random-effects probit model with time dummies and the Chamberlain (1982) device (which we denote by RE).

Each estimation technique was applied to the two different specifications of the econometric model: one using the Jaffe distance measure and one using the Mahalanobis distance measure. Below we shall discuss the empirical results in detail and perform additional econometric robustness checks.

3.4. Empirical Results

In a simple model of R&D BSV show that it is possible to derive a number of theoretical implications of these two spillover effects. If we assume that the production of knowledge is exogenous then we would not expect to find an effect of market rivalry on patent counts. Empirically, this means that the coefficient on *SpillSIC* should be close to zero. The presence of positive market spillover effects may, however, indicate endogenous patenting activity. Thus we can investigate the extent to which strategic patenting activity is consistent with the evidence in the data.

At the same time we expect the marginal effect of technology spillovers on patent counts to be positive. The production of knowledge benefits from the innovation activity in a firm conditional on its own R&D stock. Empirically this implies that we should expect the coefficient on *SpillTech* to be positive and significant.

We will test these predictions using our model and several alternative benchmark models that are commonly applied in the empirical literature. Recall that we define the dependent variable to be one if the given firm registered a patent during the particular year or not. We can think of this case as an indicator of innovation for a given firm–year dyad. We then regress this indicator on the measures of technological and product market spillovers discussed above: *SpillTech* and *SpillSIC*. In order to control for observed firm-level heterogeneity, we include two additional variables. One corresponds to firm sales $\ln(\text{Sales})$, while the other corresponds to the pre-existing stock of R&D available within the firm $\ln(\text{R\&D stock})$. Furthermore, we lag all right-hand-side variables by one period so as to remove the possibility of contemporaneous effects.

3.4.1. Partial Effects

Estimation results on the partial effects and the latent common time component are reported in Table 2. In the absence of endogenous patenting activity we should see the marginal effect of technology spillovers *SpillTech* on patenting activity to be positive and the marginal effect of product market spillovers *SpillSIC* to be zero.

Across the various model specifications we find the effect of market rivalry to be small and statistically insignificant, with the exception of PLEP for the Mahalanobis distance. Moreover, the effect changes sign depending on which distance measure is used. The evidence presented therefore does not reject the

Table II. Estimation of patent equation

	Jaffe distance				Mahalanobis distance			
	FE probit	RE probit	PLEP	FLEP	FE probit	RE probit	PLEP	FLEP
<i>Ln(SpillTech)</i>	-0.247 (0.233)	0.251* (0.072)	0.389* (0.047)	0.364* (0.046)	-0.125 (0.332)	0.660* (0.104)	0.548* (0.029)	0.641* (0.061)
<i>Ln(SpillSIC)</i>	0.073 (0.074)	0.011 (0.059)	0.021 (0.021)	0.037 (0.019)	-0.179 (0.134)	-0.200 (0.102)	-0.035* (0.013)	-0.024 (0.031)
<i>Ln(R&D Stock)</i>	0.091* (0.044)	0.144* (0.037)	0.304* (0.028)	0.313* (0.027)	0.095* (0.044)	0.143* (0.037)	0.267* (0.016)	0.304* (0.031)
<i>Ln(Sales)</i>	0.377* (0.050)	0.288* (0.044)	0.204* (0.027)	0.204* (0.022)	0.383* (0.050)	0.292* (0.043)	0.100* (0.013)	0.177* (0.023)
APE scale	0.203	0.187	0.169	0.170	0.203	0.184	0.220	0.169
ρ			0.654* (0.142)	0.645* (0.145)			0.670* (0.146)	0.552* (0.164)

Note: Standard errors are reported in parentheses. Coefficients significant at 5% confidence level are marked with an asterisk. All independent variables are lagged by one period. All regressions include a constant and a dummy for observations where lagged R&D stock is zero.

hypothesis of exogenous knowledge production. Moreover, across all specifications we observe positive and significant effects of the lagged R&D stock and the lagged sales, which is consistent with basic economic intuition.

If we use the FE probit model the estimated coefficient on *SpillTech* is not statistically significant. Moreover, the estimate appears to indicate a negative effect of technology spillovers, which contradicts economic theory. RE, PLEP and FLEP predict a positive and statistically significant effect of technology spillovers. The estimated magnitude differs, however, for each method. Note that while the results are qualitatively very similar for both the Jaffe distance measure and the Mahalanobis distance, they are quantitatively different.

It is noteworthy that both RE and FLEP produce results consistent with economic theory. Thus it is worth investigating further which model performs better on other fronts. The quantitative difference in the estimated coefficients indicates that these models may in fact produce very different predictions. From a policy perspective we would like to know which model to use for more accurate predictions. To verify this point we contrasted the outcomes predicted by each method in both models with the actual outcomes observed in the data (Table 3).¹ In the Jaffe distance model, the FLEP predicted correctly 86% of outcomes and incorrectly 14% of outcomes, while the RE predicted correctly 79% of outcomes and incorrectly 21% of outcomes. In the Mahalanobis distance model, the FLEP predicted correctly 86% of outcomes and incorrectly 14% of outcomes, while the RE predicted correctly 81% of outcomes and incorrectly 19% of outcomes. On average, the RE has thus 48% higher prediction error rate than the FLEP.

The latent effects models also indicate the existence of a time factor, measured as having moderate persistence over time with an autocorrelation coefficient of approximately 0.5. Recall that equation (10) implies that the APE scale depends not only on the estimated β coefficients but also on the estimates of the latent individual and time variables. A more precise estimate of the distribution of unobserved heterogeneity should improve the estimates of the APEs. The results in Table 2 show that the flexible latent effects probit model estimates a marginal effect for technology spillovers that is substantially larger than the marginal effect estimated by the random-effects probit model. Such differences in the estimated quantities of interest may lead to very different policy implications.

¹ This tabulation is often termed a 'misclassification' (or 'confusion') matrix. The diagonal elements contain correctly predicted outcomes, while the off-diagonal ones contain incorrectly predicted (confused) outcomes (Kohavi and Provost, 1998).

Table III. Actual vs. predicted outcomes for RE and FLEP in the Jaffe distance model and Mahalanobis distance model

	Predicted	RE		FLEP	
		0	1	0	1
	Actual				
Jaffe distance	0	4473	1470	5026	917
	1	1190	5795	911	6074
Mahalanobis distance	0	4591	1352	5020	923
	1	1148	5837	911	6074

3.4.2. Unobserved Firm Heterogeneity

We have noted above that both FLEP and PLEP produce quantitatively similar results for both distance measures.² A key advantage of FLEP is that it does not impose the normality constraint on the unobserved heterogeneity. Furthermore, FLEP is the only model that allows us to uncover a nonparametric estimate of the distribution of firm heterogeneity. There is no sound economic reason to assume that this distribution is normal and in fact we would expect that different types of production processes have very different forms of unobserved heterogeneity which impact patenting activity. In our application, the distribution of heterogeneity is shown to have a multimodal clustering structure, as plotted in Figure 1. These clusters may reflect the presence of missing variables important for characterizing innovation, such as firm culture or investment strategy. In Figure 1 we can easily discern several major clustering structures in each distance model, labeled by numbers in square boxes, corresponding to the major modes of the distribution. In the Appendix we show that this clustering is robust to the choice of the DPM prior hyperparameter.

In the FLEP output in Figure 1, each clustering structure is composed of draws of the firm-specific unobserved heterogeneity component τ_i , which we can use to further analyze the composition of each cluster. Table 4 lists the SIC code names for 20 firms whose τ_i was most frequently drawn within each given clustering structure. Thus, for the Jaffe distance model, the lowest unobserved heterogeneity component group (Clustering 1) is composed, for example, of ‘meat packing plants’, ‘blowers and fans’, or ‘department stores’; the medium unobserved heterogeneity component group (Clustering 2) includes ‘food and kindred products’, ‘footwear’, and ‘electrical industrial apparatus’; while the high unobserved heterogeneity component group (Clustering 3) features ‘semiconductors and related devices’, ‘electronic components’, or ‘commercial physical research’. The cluster composition is very similar in the Mahalanobis distance model and hence not reported here. Uncovering such cluster structures of firms that behave similarly in terms of their unobserved characteristics can provide important insights for industry analysts and policy makers analyzing firms’ R&D behavior. Below we further investigate the extent to which our theoretical predictions are satisfied for each cluster.

Using the FLEP output we can also explore the evolution of the draws of the unobserved heterogeneity parameters τ_i for specific individual companies in order to investigate whether the draws are concentrated or show large variances as the Markov chain progresses. Overall, we have found the draws to be remarkably stable, indicating a clear tendency of the model to associate each firm with a narrow range of draws of τ_i . If a draw of τ_i jumps to a different cluster, it does not stay there long and returns shortly back to its long-term average. This is an important indicator that the clustering of τ_i values observed in the estimated distribution of unobserved heterogeneity may contain relevant information since it establishes a fairly tight

² Nonetheless, PLEP estimated statistically significant product market spillovers, which was not confirmed by any other model specification.

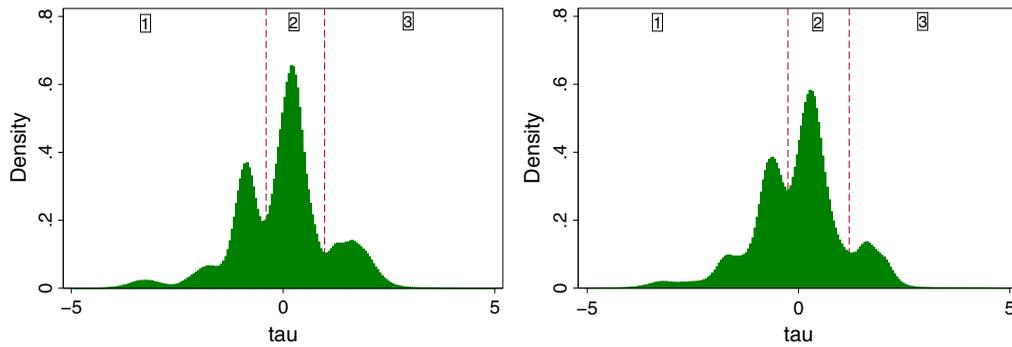


Figure 1. Flexible latent effects probit (FLEP) distribution of unobserved heterogeneity for the Jaffe distance model (left) and Mahalanobis distance model (right)

link between firms and different modes of the distribution of heterogeneity. To exemplify, we plot the draws of τ_i for the first five firms for the model of positive patents model in Figure 2.

The estimated distribution of the unobserved heterogeneity provides valuable economic information which can be used to analyze the data further and test additional economic hypotheses of interest. One

Table IV. Most frequent members of clustering structures for the Jaffe distance model

Clustering 1		Clustering 2		Clustering 3	
4923	Gas Transmission and Distribution	3841	Surgical and Medical Instruments	2835	Diagnostic Substances
2851	Paints and Allied Products	3661	Telephone and Telegraph Apparatus	2840	Cosmetics
3060	Fabricated Rubber Products	3825	Instruments To Measure Electricity	3714	Motor Vehicle Parts and Accessories
3440	Fabricated Structural Metal	3577	Computer Peripheral Equipment	2761	Manifold Business Forms
2011	Meat Packing Plants	4011	Railroads, Line-haul Operating	2390	Misc. Fabricated Textile
2731	Book Publishing	3590	Misc. industrial machinery	3823	Process Control Instruments
3561	Pumps and Pumping Equipment	3537	Industrial Trucks and Tractors	3572	Computer Storage Devices
3640	Electric Lighting and Wiring Equipment	6324	Hospital and Medical Service Plans	3674	Semiconductors and Related Devices
2030	Canned, Frozen, and Preserved Fruit	2000	Food and Kindred Products	3679	Electronic Components
3533	Oil and Gas Field Machinery	3669	Communications Equipment	3420	Handtools
3663	Radio and TV Communications Eqpt	3310	Steel Works, Blast Furnaces	2842	Sanitation Goods
3944	Games, Toys, and Children's Vehicles	3140	Footwear, Except Rubber	3990	Misc. Manufacturing Industries
3621	Motors and Generators	3620	Electrical Industrial Apparatus	8731	Commercial Physical Research
2253	Knit Outerwear Mills	2834	Pharmaceutical Preparations	3613	Switchgear and Switchboard Apparatus
3579	Office Machines	3711	Motor Vehicles and Car Bodies	3670	Electronic Components and Accessories
3743	Railroad Equipment	4931	Electric and Other Services Combined	3530	Material Handling Equipment
3490	Miscellaneous Fabricated Metal Products	3021	Rubber and Plastics Footwear	3829	Measuring and Controlling Devices
3564	Blowers and Fans	2911	Petroleum Refining	8731	Commercial Physical Research
2522	Office Furniture, Except Wood	3569	General Industrial Machinery	2821	Plastics Materials and Resins
5311	Department Stores	3690	Misc. Electrical Machinery	3861	Photographic Equipment and Supplies

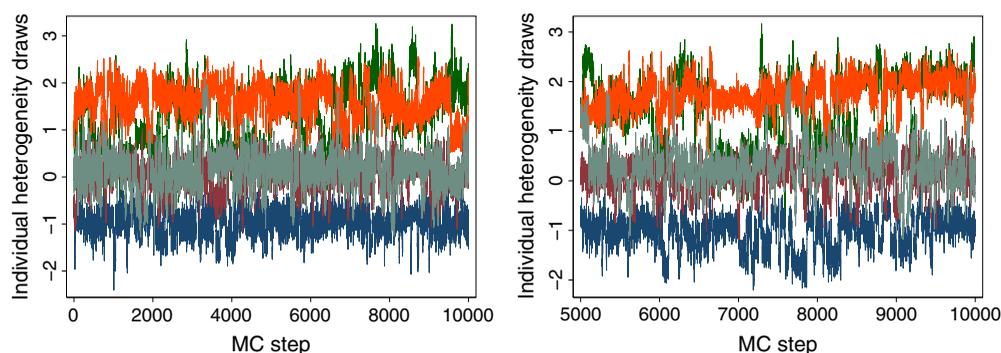


Figure 2. Draws of τ_i for the first five companies. Jaffe distance model (left) and Mahalanobis distance model (right)

such hypothesis is that there are unobserved industry-level factors driving innovation. In order to investigate this hypothesis, we plot the average sampled value of the unobserved heterogeneity component τ_i by each SIC code for the two estimated models of innovation in Figure 3. The absence of any discernible pattern in the graphs suggests that the unobserved heterogeneity component is not driven by industry factors but is rather firm-specific at the individual level. Indeed, further examination of individual τ_i revealed large differences in company types even within SIC categories. It appears that associating the unobserved heterogeneity with industry categories and attempting to capture it, for example by industry indicator variables, may obscure important differences among firms regarding their innovation activity. When we re-estimated the models using industry dummies in addition to the variables introduced above, the resulting changes were negligible. The nonparametric density estimates of the unobserved heterogeneity were almost identical to the ones previously discussed. This further highlights the benefit of the FLEP model in tracking unobserved heterogeneity at the individual level.

3.4.3. Cluster-Based Partial Effects

Given that we have established the presence of three major clusters in the distribution of firm heterogeneity, we can revisit our original model and re-estimate it separately for each cluster. This allows us to investigate the extent to which the strength of the spillover effects varies across groups of firms. As BSV emphasize, an important robustness check for the economic model is to verify the extent to which the results hold across groups of firms. If they do not, this may indicate that the estimated spillover effects are spuriously generated by pooling across different types of firms. The summary statistics for the firms in each cluster are given in Table 5. It is interesting to note that the extent of patenting activity varies substantially across

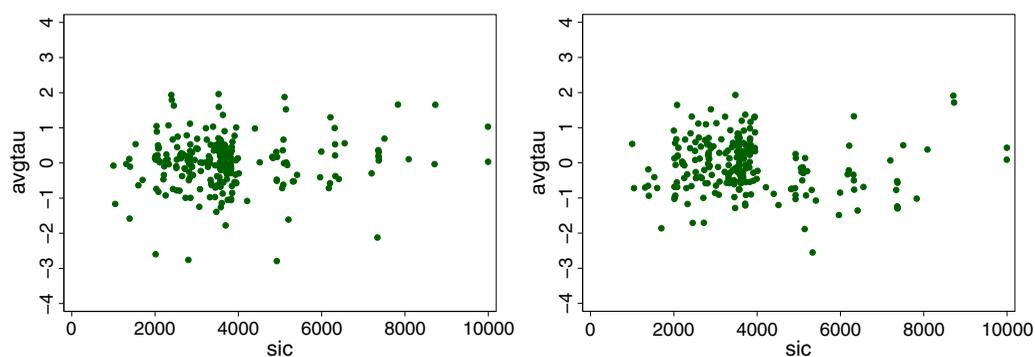


Figure 3. Average τ_i of companies for each SIC. Jaffe distance model (left) and Mahalanobis distance model (right)

Table V. Summary statistics for individual clusters. All variables are in logarithms and lagged by one period

Variable	Jaffe distance				Mahalanobis distance			
	Mean	SD	Min.	Max.	Mean	SD	Min.	Max.
Cluster 1								
<i>SpillTech</i>	9.516	1.029	4.838	11.707	11.252	0.801	8.519	13.156
<i>SpillSIC</i>	7.157	2.420	-4.602	11.154	8.417	1.685	1.220	11.559
<i>R&D Stock</i>	4.112	1.996	-0.888	10.231	2.378	2.909	-1.482	10.231
<i>Sales</i>	6.143	1.922	1.098	11.760	6.144	1.946	1.098	11.760
<i>Patenting</i>	0.262	0.440	0	1	0.291	0.454	0	1
Firms	219				260			
Total obs.	3916				4593			
Cluster 2								
<i>SpillTech</i>	9.621	1.181	4.935	11.531	11.377	0.868	8.252	13.110
<i>SpillSIC</i>	7.388	2.283	-4.321	11.053	8.585	1.656	-0.356	11.355
<i>R&D Stock</i>	4.694	2.219	-2.513	10.765	3.511	3.082	-2.513	10.765
<i>Sales</i>	6.286	2.023	0.693	12.103	6.320	2.000	0	12.103
<i>Patenting</i>	0.622	0.484	0	1	0.649	0.477	0	1
Firms	415				402			
Total obs.	7361				7155			
Cluster 3								
<i>SpillTech</i>	9.346	1.195	5.017	11.411	11.172	0.860	8.235	12.840
<i>SpillSIC</i>	7.032	2.233	-2.700	10.981	8.4523	1.565	2.631	11.268
<i>R&D Stock</i>	4.318	1.698	0.085	8.306	4.0012	1.530	0.085	8.161
<i>Sales</i>	6.187	1.764	0	10.430	6.0130	1.752	1.386	10.430
<i>Patenting</i>	0.832	0.373	0	1	0.8457	0.361	0	1
Firms	95				67			
Total obs.	1651				1180			

clusters. The first cluster corresponding to negative values for the firm-level heterogeneity has a low degree of patenting activity, while the third cluster corresponding to positive values of the firm-level heterogeneity has a high degree of patenting activity. The summary statistics for the observable variables are, however, fairly similar across clusters, which indicates that the unobserved heterogeneity has an important role to play. The results of the patent equation estimation with FLEP are given in Table 6 for each cluster, respectively. The presence of technology spillover effects is confirmed for each cluster individually. The effect of product market rivalry continues to be statistically negligible for each cluster. These results show that the economic model is thus robust to unobserved heterogeneity.

Table VI. Estimation of patent equation by cluster with FLEP

	Jaffe distance			Mahalanobis distance		
	Cluster 1	Cluster 2	Cluster 3	Cluster 1	Cluster 2	Cluster 3
<i>Ln(SpillTech)</i>	0.681* (0.077)	0.529* (0.043)	0.720* (0.097)	0.852* (0.069)	0.610* (0.040)	0.897* (0.124)
<i>Ln(SpillSIC)</i>	0.026 (0.026)	0.035 (0.019)	0.016 (0.054)	-0.031 (0.027)	-0.008 (0.017)	-0.066 (0.069)
<i>Ln(R&D Stock)</i>	0.219* (0.035)	0.345* (0.022)	0.290* (0.088)	0.198* (0.033)	0.332* (0.024)	0.385* (0.117)
<i>Ln(Sales)</i>	0.157* (0.028)	0.182* (0.019)	0.347* (0.068)	0.169* (0.024)	0.177* (0.017)	0.205* (0.086)
APE scale	0.192	0.170	0.067	0.191	0.167	0.054
ρ	0.671* (0.147)	0.620* (0.150)	0.274* (0.214)	0.468* (0.197)	0.689* (0.130)	0.514* (0.196)

Note: Standard errors are reported in parentheses. Coefficients significant at 5% confidence level are marked with an asterisk. All independent variables are lagged by one period. All regressions include a constant and a dummy for observations where lagged R&D stock is zero.

Table VII. Actual vs. predicted outcomes for RE and FLEP in the Jaffe distance model and Mahalanobis distance model for each cluster

	Predicted	RE		FLEP	
		0	1	0	1
	Actual				
Cluster 1	0	2668	219	2674	213
	1	412	617	426	603
Cluster 2	0	2212	567	2195	584
	1	480	4102	497	4085
Cluster 3	0	213	64	214	63
	1	26	1348	26	1348
Cluster 1	0	2987	269	2987	270
	1	481	856	471	866
Cluster 2	0	1959	546	1922	583
	1	463	4187	443	4206
Cluster 3	0	145	37	147	35
	1	20	978	19	979

We can also perform one additional robustness check. If the FLEP model has correctly identified each cluster we should be able to estimate the model reasonably well using RE by sub-setting the data for each cluster. If heterogeneity is driving the results of the model once we condition on a cluster RE should perform similar to FLEP.³ The robustness check results are reported in Table 7, using subsets of data corresponding to each cluster. RE and FLEP perform similarly in terms of predictions. It is important to remember, however, that this exercise can only be performed in post estimation, conditional on the given cluster. A priori, we can never be sure about the structure of the distribution of the unobserved heterogeneity, which emphasizes the importance of using a flexible model when addressing unobserved heterogeneity.

4. CONCLUSION

This paper introduced a new Bayesian semi-parametric approach to the estimation of the probit model in panel data with unobserved heterogeneity. The proposed model substantially improved on current benchmark methods by relaxing three assumptions that are often either ignored or treated in an ad hoc fashion in empirical work. First, we modeled unobserved individual effects using a flexible nonparametric form with desirable local adaptability properties. Second, we allowed for the unobserved heterogeneity to be correlated with the observables. Finally, our model incorporated common latent time effects.

We employed a combination of recent powerful sampling algorithms in order to draw from a DP Mixture model specified for the unobserved heterogeneity component. We evaluated the proposed model in a number of Monte Carlo simulations along with existing fixed and random effects model alternatives. The underlying parameters are shown to be estimated with high precision in the proposed model, unlike for the benchmark cases. The simulations presented in the online supporting information highlight the benefit of using the flexible proposed model when the underlying heterogeneity is not well approximated by a parametric distributional form.

We applied the proposed method to the estimation of a patent equation in the presence of both technological and product market spillover effects. We showed that technological innovation is subject

³ Note that FLEP also controls for the presence of time factors which are ignored by the RE model. If, however, firm-specific heterogeneity dominates in the data, then we would expect both models to have similar performance.

to substantial firm-level heterogeneity which persists within individual industries. We have shown that innovation depends in an important way on technology spillovers but that there is little evidence in favor of product market spillover effects. On the basis of the estimated firm-level heterogeneity we also showed that unobserved heterogeneity is heavily clustered and that the clustering matters when making in-sample predictions of patenting activity.

ACKNOWLEDGEMENTS

We are grateful to Nick Bloom for sharing the data for the empirical application with us. We also thank Siddhartha Chib, David Dahl, Christian Gourieroux, Jerry Hausman, Ivan Jeliazkov, Andriy Norets, seminar participants at UC Berkeley and UC Irvine, and conference audiences at the North American Summer Meetings of the Econometric Society in Boston 2009, and the Canadian Econometrics Study Group 2009 Ottawa meetings for useful comments. This work was made possible by the facilities of the Shared Hierarchical Academic Research Computing Network (SHARCNET: www.sharcnet.ca).

REFERENCES

- Abrevaya J. 1999. Leapfrog estimation of a fixed-effects model with unknown transformation of the dependent variable. *Journal of Econometrics* **93**(2): 203–228.
- Albert J, Chib S. 1993. Bayesian analysis of binary and polychotomous response data. *Journal of the American Statistical Association* **88**(422): 669–679.
- Albert J, Chib S. 1996. Bayesian modeling of binary repeated measures data with application to crossover trials. In *Bayesian Biostatistics* Berry DA, Stangl DK (eds). Marcel Dekker: New York; 577–600.
- Allenby G, Lenk P. 1994. Modeling household purchase behavior with logistic normal regression. *Journal of the American Statistical Association* **89**: 1218–1231.
- Arellano M, Bonhomme S. 2009. Robust priors in nonlinear panel data models. *Econometrica* **77**(2): 489–536.
- Arellano M, Hahn J. 2006. A likelihood-based approximate solution to the incidental parameter problem in dynamic nonlinear models with multiple effects. Working paper, CEMFI.
- Berry ST, Haile PA. 2009. Nonparametric identification of multinomial choice demand models with heterogeneous consumers. Cowles Foundation Discussion Papers 1718, Cowles Foundation, Yale University.
- Bloom N, Schankerman M, Van Reenen J. 2010. Identifying technology spillovers and product market rivalry. NBER Working Paper 13060.
- Bolduc D, Fortin B, Gordon S. 1997. Multinomial probit estimation of spatially interdependent choices: an empirical comparison of two new techniques. *International Regional Science Review* **20**(1–2): 77101.
- Burda M, Harding MC, Hausman JA. 2008. A Bayesian mixed logit–probit model for multinomial choice. *Journal of Econometrics* **147**(2): 232–246.
- Burda M, Liesenfeld R, Richard J-F. 2011. Bayesian analysis of a probit panel data model with unobserved individual heterogeneity and autocorrelated errors. *International Journal of Statistics and Management Systems* **6**(1–2): 1–21.
- Chamberlain G. 1982. Multivariate regression models for panel data. *Journal of Econometrics* **18**: 5–46.
- Chamberlain G. 1984. Panel data. In *Handbook of Econometrics*, Vol. 2, Griliches Z, Intriligator M (eds). North-Holland: Amsterdam; 1247–1318.
- Chib S. 1993. Bayes estimation of regressions with autoregressive errors: a Gibbs sampling approach. *Journal of Econometrics* **58**: 275–294.
- Chib S, Carlin B. 1999. On MCMC sampling in hierarchical longitudinal models. *Statistics and Computing* **9**: 17–26.
- Chib S, Hamilton B. 2002. Semiparametric Bayes analysis of longitudinal data treatment models. *Journal of Econometrics* **110**: 67–89.
- Chib S, Jeliazkov I. 2006. Inference in semiparametric dynamic models for binary longitudinal data. *Journal of the American Statistical Association* **101**(474): 685–700.
- Dahl DB. 2005. Sequentially-allocated merge–split sampler for conjugate and nonconjugate Dirichlet process mixture models. Technical report, Texas A&M University.
- Escobar MD, West M. 1995. Bayesian density estimation and inference using mixtures. *Journal of the American Statistical Association* **90**: 577–588.
- Fernandez-Val I. 2007. Fixed effects estimation of structural parameters and marginal effects in panel probit models. Working paper, Boston University.

- Geweke J. 1991. Efficient simulation from the multivariate normal and Student-t distributions subject to linear constraints. In *Computing Science and Statistics: Proceedings of the Twenty-Third Symposium on the Interface*, Keramidas EM (ed.). Interface Foundation of North America: Fairfax, VA.
- Geweke J. 2005. *Contemporary Bayesian Econometrics and Statistics*. Wiley: Chichester.
- Griffith R, Lee S, Van Reenen J. 2011. Is distance dying at last? Falling home bias in fixed effects models of patent citations. *Quantitative Economics* **2**(2): 211–250.
- Gu Y, Fiebig DG, Cripps E, Kohn R. 2009. Bayesian estimation of a random effects heteroscedastic probit model. *Econometrics Journal* **12**(2): 324–339.
- Hajivassiliou V. 1990. Smooth estimation simulation of panel data LDV models. Working paper, Yale University.
- Hausman JA, Taylor WE. 1981. Panel data and unobservable individual effects. *Econometrica* **49**(6): 1377–1398.
- Hirano K. 2002. Semiparametric Bayesian inference in autoregressive panel data models. *Econometrica* **70**: 781–799.
- Jaffe AB. 1986. Technological opportunity and spillovers of R&D: evidence from firms' patents, profits, and market value. *American Economic Review* **76**(5): 984–1001.
- Keane M. 1990. A computationally efficient practical simulation estimator for panel data, with applications to estimating temporal dependence in employment and wages. Working paper, University of Minnesota.
- Knittel CR, Metaxoglou K. 2008. Estimation of random coefficient demand models: challenges, difficulties and warnings. Working paper, NBER.
- Kohavi R, Provost F. 1998. Glossary of terms. Editorial for the special issue on application of machine learning and the knowledge of discovery process. *Machine Learning* **30**: 271–274.
- Lancaster T. 2004. *An Introduction to Modern Bayesian Econometrics*. Blackwell: Malden, MA.
- Li T, Zheng X. 2008. Semiparametric Bayesian inference for dynamic Tobit panel data models with unobserved heterogeneity. *Journal of Applied Econometrics* **23**: 699–728.
- Neal R. 2000. Markov chain sampling methods for Dirichlet process mixture models. *Journal of Computational and Graphical Statistics* **9**(2): 249–265.
- Paap R. 2002. What are the advantages of MCMC based inference in latent variable models? *Statistica Neerlandica* **56**(1): 2–22.
- Rossi P. 2010. Bayesm: Bayesian inference for marketing/micro-econometrics. R package version 2.2-3.
- Train K. 2001. A comparison of hierarchical Bayes and maximum simulated likelihood for mixed logit. Working paper, Department of Economics, University of California, Berkeley.
- Train K. 2003. *Discrete Choice Methods with Simulation*. Cambridge University Press: Cambridge, UK.
- Troughton P, Godsill S. 1997. Bayesian model selection for time series using Markov chain Monte Carlo. In *IEEE International Conference on Acoustics, Speech, and Signal Processing, 1997 (ICASSP-97)*, Vol. **5**: 3733–3736.
- Wooldridge JM. 2001. *Econometric Analysis of Cross Section and Panel Data*. MIT Press: Cambridge, MA.

APPENDIX A

Sampling β

In this block we apply the method of Albert and Chib (1993) to the recentered latent variable $\tilde{y}_{it}^* = \tilde{y}_{it} - \tau_i - \lambda_t$. The joint conditional density of $(\beta, \tilde{\mathbf{y}}^*)$ is given by

$$p(\beta, \tilde{\mathbf{y}}^* | \tau, \lambda, \psi, \theta_{/\beta}, \tilde{\mathbf{y}}, \mathbf{y}, \mathbf{X}) \propto \exp\left[-\frac{1}{2}(\beta - \underline{\beta})' \underline{\Sigma}_{\beta}^{-1} (\beta - \underline{\beta})\right] \exp\left[-\frac{1}{2}(\tilde{\mathbf{y}}^* - \mathbf{X}\beta)' (\tilde{\mathbf{y}}^* - \mathbf{X}\beta)\right]$$

yielding a closed form of the conditional posterior for β which facilitates direct sampling from $\beta | \cdot \sim N(\bar{\beta}, \bar{\Sigma})$ where

$$\begin{aligned} \bar{\beta} &= \bar{\Sigma} \left(\underline{\Sigma}^{-1} \underline{\beta} + \mathbf{X}' \tilde{\mathbf{y}}^* \right) \\ \bar{\Sigma} &= \left(\underline{\Sigma}^{-1} + (\mathbf{X}' \mathbf{X}) \right)^{-1} \end{aligned}$$

In the application, we specify the hyperparameter values $\underline{\beta} = 0$ and $\underline{\Sigma}_{\beta} = 10I$, where I is the identity matrix. This specification is aimed at rendering the prior for β sufficiently diffuse.

Sampling \tilde{y}_{it}

Here we benefit from the second step of the Albert and Chib (1993) procedure, augmented by τ_i and λ_t . Thus we sample directly

$$\begin{aligned} \tilde{y}_{it} | \cdot &\sim N(v_{it}, 1) \\ v_{it} &= \mathbf{x}_{it} \beta + \tau_i + \lambda_t \end{aligned}$$

truncated by 0 from the left if $y_{it} = 1$ and from the right if $y_{it} = 0$.

Updating Latent Class Assignments

For this block we utilize a hybrid sampler that alternates between the non-conjugate version of the SAMS sampler of Dahl (2005) and Algorithm 7 of Neal (2000). This approach is suggested by Dahl (2005) as optimally combining the virtues of each method: the ability to move large blocks of elements among latent classes in one step for the former, and one-at-a-time allocations of individual elements among latent classes for the latter.

The SAMS sampler is based on an alternative expression of the model (6)–(8) in terms of a set partition $\pi = \{S_1, \dots, S_q\}$ for $S_0 = \{1, \dots, n\}$ in addition to the latent class parameters $\phi = \{\phi_{S_1}, \dots, \phi_{S_q}\}$, where ϕ_S is associated with component S . The set partition π for S_0 is a set of subsets S_1, \dots, S_q such that (1) $\cup_{S \in \pi} S = S_0$, (2) $S^i \cap S^j = \emptyset$ for all $S^i \neq S^j$, and (3) $S \neq \emptyset$ for all $S \in \pi$. Using this notation, the model (6)–(8) can be recast as (Dahl, 2005)

$$\tau_i | \pi, \phi \sim F_\tau(\phi_S^i) \tag{15}$$

$$\psi | \pi \sim \prod_{S \in \pi} G_0(\phi_S) \tag{16}$$

$$\pi \sim b \prod_{S \in \pi} \eta_0 \Gamma(|S|) \tag{17}$$

where $|S|$ is the number of elements of the component S . The sampling scheme works as follows. In each MC iteration, uniformly select a pair of distinct indices i and j . If i and j belong to the same component in π , say S , propose π^* by splitting S . Otherwise, i and j belong to different components in π , say S^i and S^j . Propose π^* by merging S^i and S^j . In each case, compute the Metropolis–Hastings (MH) ratio $a(\pi^*, \phi^* | \pi, \phi)$ and accept the new latent class configuration π^* with probability given by this ratio. We derive the MH ratio for our model in the following section.

Algorithm 7 of Neal (2000), which we utilize in every alternate MC step, is based on limiting probabilities of a latent class finite mixture model, with the number of classes tending to infinity. The sampling procedure itself is built around drawing with a stochastic number of mixture components or classes whose number and size varies at each MC iteration. Denote by c a label of a generic latent class with membership count N_c . Given the current state of the system, τ_i are first reassigned to latent classes with labels c_i , whereby new classes can be created and old ones may vanish. The probabilities of class assignment for the τ_i are proportional to the likelihood of τ_i conditional on the current draw of the class parameters ψ_c . Second, the class parameters ψ_c are updated in a standard way for each class separately. If we specify F_τ as an infinite mixture of Normals, then $\psi_c = (\mu_{\tau c}, \sigma_{\tau c}^2)$ are the moments of the Normal density.

For updating ψ in the Algorithm 7 scan, we specify F_τ as a mixture of Normals with $\psi = (\mu_\tau, \sigma_\tau^2)$. Since for all τ_i that fall into one latent class it holds that $\tau_i \sim N(\mu_{\tau c}, \sigma_{\tau c}^2)$ we can apply result B (p. 300) of Train (2003) to each latent class separately: for an $IG(s_0, v_0)$ prior, the posterior of $\sigma_{\tau c}^2$ is given by $IG(s_1, v_1)$ with $v_1 = v_0 + N_c$ and $s_1 = (v_0 s_0 + N_c \bar{s}_c) / (v_0 + N_c)$ where $\bar{s}_c = N_c^{-1} \sum_{i=1}^{N_c} \tau_i^2$. We utilize a diffuse IG prior. Analogously, to sample $\mu_{\tau c}$ we use result A of Train (2003) applied to each latent class. The hyperparameter of the DP prior α is sampled according to the scheme of Escobar and West (1995).

The iteration between the samplers of Dahl (2005) and Neal (2000) alleviates the influence of particular starting values. The SAMS sampler is capable of reallocating large blocks of data to one of the latent classes, while the Neal algorithm addresses the individual by individual allocation to latent classes. This allows us to initialize the procedure with a unique parametric component. This is then rapidly split into classes by the SAMS sampler before the Neal procedure continues to fine-tune the posterior draws.

Sampling τ_i

Let $\tilde{\mathbf{y}}_i^{**} = \tilde{\mathbf{y}}_i - \mathbf{X}_i \beta - \lambda$. Then

$$\tilde{\mathbf{y}}_i^{**} = \tau_i \mathbf{1} + \varepsilon_i$$

Consider for the moment the case $\tau_i \sim N(\underline{\tau}, \sigma_\tau^2)$; it will be used as a building block in the DP prior sampling. In this case, for every i we have one latent regression with one parameter τ_i and a $(T \times 1)$ vector of ones as explanatory variables in place of a hypothetical \mathbf{X}_i . Using standard latent regression results (see, for example, Lancaster, 2004)

$$p(\tau_i|\cdot) = \phi(\bar{\tau}_i, \bar{\sigma}_{\tau_i}^2) \tag{18}$$

$$\bar{\tau}_i = \bar{\sigma}_{\tau_i}^2 \left(\sigma_{\tau}^{-2} \underline{\tau} + \sum_{i=1}^T \tilde{\mathbf{y}}_i^{**} \right)$$

$$\bar{\sigma}_{\tau_i}^2 = (\sigma_{\tau}^{-2} + T)^{-1}$$

Since $\tau_i \sim N(\mu_{\tau c}, \sigma_{\tau c}^2)$ given a previous assignment to the latent class c , let $\underline{\tau} = \mu_{\tau c}$, $\sigma_{\tau}^2 = \sigma_{\tau c}^2$ and sample τ_i directly from (18).

Sampling λ

Let

$$\tilde{\mathbf{y}}_{i\lambda} = \tilde{\mathbf{y}}_i - \mathbf{X}_i \beta - \tau_i t$$

Then the joint density implied for $\tilde{\mathbf{y}}_{\lambda} = (\tilde{\mathbf{y}}_{1\lambda}, \dots, \tilde{\mathbf{y}}_{N\lambda})$ by the recentered probit model conditional on λ is

$$f(\tilde{\mathbf{y}}_{\lambda}|\lambda, \cdot) = (2\pi)^{-NT/2} \det(I_T)^{-N/2} \exp \left\{ -\frac{1}{2} \sum_{i=1}^N \left[\tilde{\mathbf{y}}_{i\lambda}' I_T^{-1} \tilde{\mathbf{y}}_{i\lambda} - 2\lambda' I_T^{-1} \tilde{\mathbf{y}}_{i\lambda} + \lambda' I_T^{-1} \lambda \right] \right\}$$

while the prior density specified by Assumption 3 takes the form

$$f(\lambda) = (2\pi)^{-T/2} \det(I_T)^{-1/2} \exp \left\{ -\frac{1}{2} \lambda' \Omega_{\lambda}^{-1} \lambda - 2\lambda' \Omega_{\lambda}^{-1} \Lambda \rho + \rho' \Lambda' \Omega_{\lambda}^{-1} \Lambda \rho \right\}$$

where $\Lambda_t = (\lambda_{t-1}, \dots, \lambda_{t-s})$, $\Lambda = (\Lambda'_1, \dots, \Lambda'_T)'$, $\rho = (\rho_1, \dots, \rho_s)$ and Ω_{λ} is the covariance matrix associated with the autoregressive process. Hence we can sample λ directly from $N(\bar{\lambda}, \bar{\Sigma}_{\lambda})$ where

$$\bar{\lambda} = \bar{\Sigma}_{\lambda} \left(\Omega_{\lambda}^{-1} \underline{\lambda} + \sum_{i=1}^N \tilde{\mathbf{y}}_{i\lambda} \right)$$

$$\bar{\Sigma}_{\lambda} = (\Omega_{\lambda}^{-1} + N \times I_T)^{-1}$$

For ease of implementation we restrict ourselves to the AR(1) specification with a single autoregressive parameter ρ in the application. In this case

$$\Omega_\lambda = \gamma_0 \begin{bmatrix} \rho^0 & \rho^1 & \rho^2 & \cdots & \rho^{T-1} \\ \rho^1 & \rho^0 & \rho^1 & \cdots & \rho^{T-2} \\ \rho^2 & \rho^1 & \rho^0 & \cdots & \rho^{T-3} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \rho^{T-1} & \rho^{T-2} & \rho^{T-3} & \cdots & \rho^0 \end{bmatrix}$$

$$\gamma_0 = \frac{\sigma_\eta^2}{1 - \rho^2}$$

Sampling ρ

Note that for the AR(1) process

$$p(\lambda_t | \lambda_{t-1}, \cdot) \propto \begin{cases} \exp\left(-\frac{(1 - \rho^2)}{2\sigma_\eta^2} \lambda_1^2\right), & t = 1 \\ \exp\left(-\frac{1}{2\sigma_\eta^2} (\lambda_t - \rho\lambda_{t-1})^2\right), & t = 2, \dots, T \end{cases}$$

and hence

$$p(\rho | \lambda) = \exp\left(-\frac{1}{2\sigma_\eta^2} \left[(1 - \rho^2) \lambda_1^2 + \sum_{t=2}^T (\lambda_t - \rho\lambda_{t-1})^2 \right]\right)$$

$$= \exp\left(\frac{1}{2\sigma_\eta^2} \left[\rho^2 \left(\sum_{t=2}^T \lambda_{t-1}^2 - \lambda_1^2 \right) - 2\rho \sum_{t=2}^T \lambda_t \lambda_{t-1} + \sum_{t=1}^T \lambda_t^2 \right]\right)$$

Matching this expression with a Gaussian kernel $\exp(-\frac{1}{2\bar{\sigma}_\rho^2} [\rho^2 - 2\rho\bar{\mu}_\rho + \bar{\mu}_\rho^2])$ yields

$$\bar{\sigma}_\rho^2 = \sigma_\eta^2 \left(\sum_{t=2}^{T-1} \lambda_t^2 \right)^{-1}$$

$$\bar{\mu}_\rho = \frac{\bar{\sigma}_\rho^2}{\sigma_\eta^2} \sum_{t=2}^T \lambda_t \lambda_{t-1}$$

$$= \left(\sum_{t=2}^{T-1} \lambda_t^2 \right)^{-1} \sum_{t=2}^T \lambda_t \lambda_{t-1}$$

We can therefore sample ρ directly from $N(\bar{\mu}_\rho, \bar{\sigma}_\rho^2)$ truncated at -1 and 1 to preserve stationarity. Extension to AR(p) will amend the likelihood function $p(\lambda_t | \lambda_{t-1}, \cdot)$ but the derivation would be similar. The approach for sampling the AR(p) parameters conditional on the initial observations is presented in Chib (1993).

Sampling σ_η^2

For this block we use the result derived in Burda *et al.* (2011), which adapts the standard result on sampling univariate variances (given, for example, by result B, p. 300 of Train, 2003) to the likelihood of the variance of the AR process. Conditional on λ and ρ , the likelihood function of σ_η^2 takes the form

$$L\left(\sigma_\eta^2 \mid \lambda, \theta / \sigma_\eta^2\right) \propto \frac{\sqrt{1-\rho^2}}{\sigma_\eta \sqrt{2\pi}} \exp\left[-\frac{1-\rho^2}{2\sigma_\eta^2} \lambda_1^2\right] \prod_{t=2}^T \frac{1}{\sigma_\eta \sqrt{2\pi}} \exp\left[-\frac{1}{2\sigma_\eta^2} (\lambda_t - \rho\lambda_{t-1})^2\right]$$

An $IG(v_0, s_0)$ prior has density

$$k\left(\sigma_\eta^2\right) = \frac{1}{m_0 \sigma_\eta^{(v_0+1)/2}} \exp\left[-\frac{v_0 s_0}{2\sigma_\eta}\right]$$

where m_0 is a normalizing constant. We can then sample directly from the posterior

$$\begin{aligned} L\left(\sigma_\eta^2 \mid \lambda, \theta / \sigma_\eta^2\right) &\propto L\left(\sigma_\eta^2 \mid \lambda, \theta / \sigma_\eta^2\right) k\left(\sigma_\eta^2\right) \\ &\propto \frac{1}{\sigma_\eta^{(T+v_0+1)/2}} \exp\left[-\frac{(1-\rho^2)\lambda_1^2 + \sum_{t=2}^T (\lambda_t - \rho\lambda_{t-1})^2 + v_0 s_0}{2\sigma_\eta^2}\right] \\ &= IG(v_1, s_1) \end{aligned}$$

where

$$\begin{aligned} v_1 &= v_0 + T \\ s_1 &= \frac{v_0 s_0 + (1-\rho^2)\lambda_1^2 + \sum_{t=2}^T (\lambda_t - \rho\lambda_{t-1})^2}{v_0 + T} \end{aligned}$$

In the application, the prior for σ_η^2 will be specified as diffuse with $s_0 \rightarrow 0$ and $v_0 = 0$.

The SAMS Sampler

In this section, we explicitly derive the form of the MH ratio for our case. For a general description, see Dahl (2005). Let k be the successive values in random permutations of the indices in S . In our model, the MH ratio is given by

$$a(\pi^*, \phi^* \mid \pi, \phi) = \min\left[1, \frac{p(\pi^*, \phi^* \mid y) q(\pi, \phi \mid \pi^*, \phi^*)}{p(\pi, \phi \mid y) q(\pi^*, \phi^* \mid \pi, \phi)}\right]$$

If the proposal involves a split, $q(\pi^*, \phi^* \mid \pi, \phi)$ is the split probability and $q(\pi, \phi \mid \pi^*, \phi^*) = 1$ is the merge probability. If the proposal involves a merge, the roles of $q(\pi^*, \phi^* \mid \pi, \phi)$ and $q(\pi, \phi \mid \pi^*, \phi^*)$ are reversed. Consider, for example, proposal for a split:

$$q(\pi^*, \phi^* | \pi, \phi) = \prod_{k=1}^N P(k \in S^i | S^i, S^j, \phi, y) P(\phi_{S^i})$$

The first term is given in equation (13) in Dahl (2005). The second term $P(\phi_{S^i})$ is the proposal density of the new ϕ_{S^i} . The merge probability is

$$q(\pi, \phi | \pi^*, \phi^*) = 1$$

By Bayes' theorem

$$p(\pi, \phi | y) \propto p(y | \pi, \phi) p(\pi, \phi) \tag{19}$$

where $p(y | \pi, \phi)$ is the likelihood

$$p(y | \pi, \phi) = \prod_{i=1}^n p(y_i | \phi_{S^i})$$

and $p(\pi, \phi)$ is the prior

$$p(\pi, \phi) = p(\phi | \pi) p(\pi) \tag{20}$$

where

$$\begin{aligned} p(\phi | \pi) &= \prod_{S \in \pi} F_0(\phi_S) \\ p(\pi) &= b \prod_{S \in \pi} \eta_0 \Gamma(|S|) \\ b^{-1} &= \prod_{i=1}^n \Gamma(\eta_0 + i - 1) \end{aligned}$$

Note that for a split of a class S^s into S^i and S^j

$$\frac{p(y | \pi^*, \phi^*)}{p(y | \pi, \phi)} = \frac{\prod_{t=1}^{|S^i|} p(y_t | \phi_{S^i}) \prod_{t=1}^{|S^j|} p(y_t | \phi_{S^j})}{\prod_{t=1}^{|S^s|} p(y_t | \phi_{S^s})} \tag{21}$$

where the index t in $p(y_t | \phi_{S^i})$ refers to elements of the class S^i . Similarly, for a merge of classes S^i and S^j into S^s

$$\frac{p(y | \pi^*, \phi^*)}{p(y | \pi, \phi)} = \frac{\prod_{t=1}^{|S^s|} p(y_t | \phi_{S^s})}{\prod_{t=1}^{|S^i|} p(y_t | \phi_{S^i}) \prod_{t=1}^{|S^j|} p(y_t | \phi_{S^j})}$$

i.e. the inverse of the ratio of split probabilities. Note that for a split we can use the stored values of the likelihood evaluations from the allocation of k into S^i and S^j . Hence only two additional likelihood evaluations $p(y_i | \phi_{S^i})$ and $p(y_j | \phi_{S^j})$ that initiated the split need to be performed for obtaining the ratio $\frac{p(y | \pi, \phi)}{p(y | \pi^*, \phi^*)}$. For a merge, only $2|S^i| + |S^j|$ likelihood evaluations need to be performed, which for small classes can be substantially less than the sample size n . In the same spirit, for computing prior components for a split

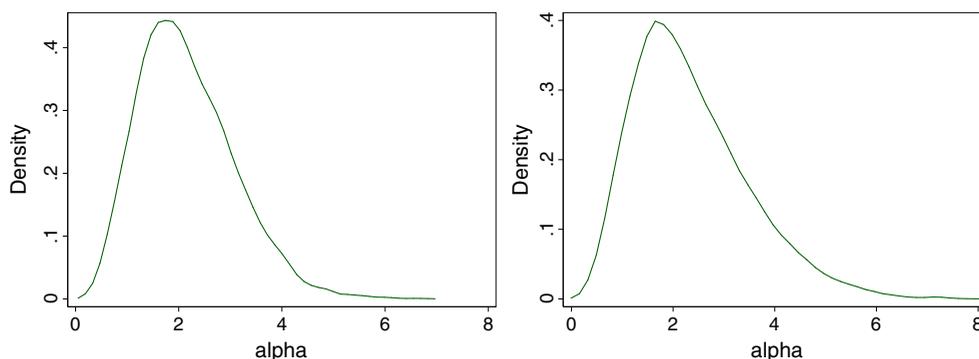


Figure 4. Density of draws of α , Jaffe distance model (left) and Mahalanobis distance model (right)

$$\frac{pt(\phi^*|\pi^*)}{pt(\phi|\pi)} = \frac{F_0(\phi_{S^i})F_0(\phi_{S^j})}{F_0(\phi_{S^s})}$$

and

$$\frac{p(\pi^*)}{p(\pi)} = \frac{\Gamma(|S^i|)\Gamma(|S^j|)}{\Gamma(|S^s|)}$$

while for a merge

$$\frac{pt(\phi^*|\pi^*)}{pt(\phi|\pi)} = \frac{F_0(\phi_{S^s})}{F_0(\phi_{S^i})F_0(\phi_{S^j})}$$

and

$$\frac{p(\pi^*)}{p(\pi)} = \frac{\Gamma(|S^s|)}{\Gamma(|S^i|)\Gamma(|S^j|)}$$

Thus, using (19), the ratio of the p -terms for a split becomes

$$\frac{pt(\pi^*, \phi^*|y)}{pt(\pi, \phi|y)} = \frac{\prod_{t=1}^{|S^i|} pt(y_t|\phi_{S^i}) \prod_{t=1}^{|S^j|} pt(y_t|\phi_{S^j}) F_0(\phi_{S^i}) F_0(\phi_{S^j}) \Gamma(|S^i|)\Gamma(|S^j|)}{\prod_{t=1}^{|S^s|} pt(y_t|\phi_{S^s}) F_0(\phi_{S^s}) \Gamma(|S^s|)} \quad (22)$$

while for a merge this ratio is given by the inverse of the expression in (22).

DPM Prior Hyperparameter

In order to explore the distribution of unobserved heterogeneity, τ_i , in our patent models we need to make sure that its behavior is not implicitly restricted by the estimation procedure or some other deep model parameters. One parameter that is of concern to us is the smoothing parameter α that controls the extent to which the DP draws mixture distributions that are more or less ‘similar’ to the Normal baseline parametric distribution G_0 . In the limiting case of $\alpha \rightarrow \infty$ the mixture distribution becomes equivalent to G_0 , while in the other extreme $\alpha \rightarrow 0$ the mixture distribution limits to a convolution of density kernels centered at each data point without any influence of the DP prior. The posterior distribution estimate for both models is plotted in Figure 4. The distributions are concentrated around a mode of 2, indicating a strong influence of data relative to the baseline prior distribution.