

# Dynamic Bayesian Hierarchical Latent Variable Models for Categorical Survey Data

Evan Munro <sup>\*</sup>      Serena Ng <sup>†</sup>

March 4, 2019

## Abstract

Existing methods for extracting summary indices from multivariate discrete data often treat discrete variables as if they were continuous. We show that Bayesian hierarchical latent variable models, popular with statisticians for text and genomic analysis, provide a good alternative to existing methods for estimating latent variables from categorical survey data. We derive these models from a set of conditional independence assumptions on the joint distribution of survey responses. We illustrate the interpretability and flexibility of these methods in an example on estimating wealth indices from DHS data in Latin America and the Caribbean. Then, we extend the basic static models and introduce two dynamic models for time series survey data, which are the discrete versions of a basic Markov switching model and state space model. We estimate the models on sets of variables from the Michigan Survey of Consumers and the Gallup Poll Social Series. We show that the indexes constructed provide quantitative metrics for difficult-to-measure concepts such as environmental concern, and have additional economic interpretability beyond averaging methods.

---

<sup>\*</sup>Graduate School of Business, Stanford University. Electronic correspondence: munro@stanford.edu,

<sup>†</sup>Department of Economics, Columbia University, and NBER. Electronic correspondence: serena.ng@columbia.edu

# 1 Introduction

As data collection has become cheaper, sources of both traditional and non-traditional discrete data have grown. These data include survey data, unorganized text data, network data, shopping basket data, and unorganized text data. The established methods for extracting low-dimensional factors from multivariate discrete data generally involve treating discrete data as if it was continuous.

For ordered multivariate discrete survey data, typical methods to aggregate response data involve assigning numerical indices to each ordered outcome and taking an average, either of each variable directly (Bloom & Reenen, 2010), or of the relative percentage of respondents selecting the “high” vs the “low” outcome (Ludvigson, 2004). Methods involving averaging, however, rely on the assumption that the distance between each ordered categorical outcome is equal, which may not be an accurate assumption. Furthermore, missing or don’t know responses, which are very common in survey data and could contain important information about uncertainty, must be dropped or imputed. Methods based on polychoric correlation have a stronger theoretical foundation (Bentler *et al.*, 1990), but cannot handle unordered categorical variables.

For extracting individual-level factors from cross-sectional data comprising unordered categorical variables, many applied economists extract an index using the Filmer-Pritchett method, which is a version of PCA for discrete data (Filmer & Pritchett, 2001). The Filmer-Pritchett method involves converting each variable with multiple outcomes to a set of binary variables, and factorizing a matrix of binary responses. However, converting multinomial outcomes to binary variables in this way introduces spurious negative correlations within the multiple columns that are mapped from a single question. The factors extracted from matrix factorization, then, are not an optimal low-dimensional representation of the data (Kolenikov & Angeles, 2009).

This paper makes two major contributions. The first is showing how hierarchical bayesian models popularized in the CS literature for text analysis can be used to extract interpretable factors from multivariate ordered or unordered categorical survey data, without treating discrete outcomes as if they were continuous. We show that Grade of Membership (Erosheva *et al.*, 2007) is an interpretable alternative to the Filmer-Pritchett method for estimating individual-level indices from cross-sectional survey data. We also show that Latent Dirichlet Allocation (Blei *et al.*, 2003), a popular method for text analysis, can be used in place of averaging methods for estimating group-level indices from cross-sectional survey data. Grade of Membership and Latent Dirichlet Allocation are more flexible forms of a latent class model, which estimates groups in the data based on an aggregate-level mixture.

The second major contribution is showing how to introduce dynamics in Bayesian parametric models for time-dependent survey data with categorical outcomes. We also show that depending on the flexibility and computational complexity of the model desired, dynamics can be added using a discrete-valued latent variable or a continuous-valued latent variable.

The first dynamic method introduced is a version of Hamilton (1989)’s markov switching model for categorical response variables, and has been suggested in the statistics literature by MacDonald & Zucchini (1997). The model assumes that there is a single hidden state at each point in time in

the economy that generates outcome variables; this state is one of a finite set of states and follows a markov-switching process. The second is a version of a linear state space model for categorical response variables. The model assumes that at each point in time, there are a mixture of states that generate the outcome variables. The mixture probabilities follow an random walk process and each state corresponds to a different multinomial distribution over categorical outcomes.

It is important to have methods specific to discrete data not only to avoid unwarranted assumptions, but also to avoid losing important economic information that are contained in discrete data that don't have analogous continuous data. The Index of Consumer Sentiment published by the University of Michigan, for example, is constructed from five questions from the survey that have ordered responses based on views about current and future economic conditions. A series of studies have not found that the index provides little to no forecasting power for consumer spending patterns beyond what is found in aggregate statistics on actual consumer purchases (Ludvigson, 2004). As a sanity check for the new methods, in the application section we show that the extracted index from the Michigan data closely correlates with the existing ad-hoc index, but provides additional information on patterns in responses that are driving fluctuations in sentiment. This allows us to distinguish what sort of information the discrete data does provide beyond what is contained in series published on real personal consumption expenditure, for example.

Unlike consumer spending, there are many important economic variables, however, that are not well-measured by spending and other aggregate economic indicators. For example, fluctuations in consumers' environmental sentiments are challenging to measure when the "environment" is not a tangible good, many recreation activities are free and most pollution controls are regulation rather than market-based. In this case qualitative information provided by ordered and unordered categorical variables may provide important insights about aggregate sentiments that are not easily quantified otherwise. For the second application, we extract indices from unordered categorical variables from Gallup Poll Social Series (GPSS) data, which do not have established indices published, and show how the resulting indices provide quantitative insights into consumer's preferences about non-tangible goods around recessions.

## 2 Static Models for Categorical Data

For the rest of the paper, we assume that we have a dataset with  $i = 1, \dots, N$  survey respondents, who each respond to  $j = 1, \dots, J$  survey questions with categorical outcomes. Each survey question has possible outcomes indexed by  $l_j = 1, \dots, L_j$ . A missing or don't know response is included as response  $L_j$  for each question. The number of choices can vary across survey questions. Survey respondents may be grouped into  $c = 1, \dots, C$  groups, in which case  $N_c$  denotes the number of respondents in group  $c$ . Groups can be formed in a variety of ways: respondents who are of the same sex, the same country, or who responded to a survey at the same time, for example. The raw dataset  $X$  is an  $N \times J$  matrix where  $X_{ij} \in \{1, \dots, L_j\}$  corresponds to the index of the outcome that survey respondent  $i$  chose for question  $j$ . For models where respondents are assigned to groups,

we sometimes add a third index  $c$  to denote the group membership of individual  $i$ , and refer to responses  $X_{cij}$ .

## 2.1 Existing Methods

The first setting that we will examine is where the goal is to recover a  $K - 1$ -dimensional real-valued index  $g_i$  that is a low dimensional representation of  $X_i$ , the vector of survey responses for individual  $i$ . Sometimes  $K = 2$  and we assign a single index to each individual based on their survey responses. For example, there is a lack of reliable income data in the developing world. Many researchers use survey data on household consumption from the Demographic and Health Surveys (DHS), or the Living Standard Measurement Surveys (LSMS), to extract a wealth index that proxies for income for each household in the sample.

These wealth indices are often created using the Filmer-Pritchett method (Filmer & Pritchett, 2001), which is a form of PCA for categorical variables. All rows of  $X$  that contain a missing-data response are generally dropped. Each categorical variable  $j$  is transformed into  $j^*$  binary variables, where  $j^* = L_j - 1$ .  $\tilde{J} = \sum_{j=1}^J j^*$ . Each column is normalized by the its mean and standard deviation to derive a transformed  $N \times \tilde{J}$  matrix  $\tilde{X}$  of data.

Let  $S$  be the  $\tilde{J} \times \tilde{J}$  sample covariance matrix of  $\tilde{X}$ . The first  $k$ -principal components are

$$Z = \tilde{X}A$$

where  $A$  is the first  $k$  of the  $\tilde{J}$  total eigenvectors of  $S$ , normalized to a unit vector.

The  $N$ -dimensional first principal component, the first column of  $Z$ , is used as the estimated wealth index  $g_i^{(FP)}$ . The index corresponds to a weighted sum of indicators for each of the possible responses in the survey, corresponding to indicators for asset ownership, for example, which explains the maximum variance in the data. The main issue with Filmer-Pritchett is for binary variables which are derived from categorical variables, there are spurious negative correlations introduced between variables that come from the same categorical variable: the directions of maximum variance may be related to the spurious negative correlations in the augmented data matrix rather than in directions that correspond to differences in wealth between households in the survey. Kolenikov & Angeles (2009) and Ng (2015) have pointed out these issues and examined alternative principal components and factor analysis methods for extracting latent variables from both continuous and categorical data, although neither found clear evidence that a single method is best.

One alternative method involves treating discrete data as continuous and doing PCA on the scaled  $X$  directly. This relies on assuming that the ordinal data has constant distance between different categories, which may not be a good assumption for many survey datasets. Another alternative method which we will examine is based on polychoric correlation (Lee *et al.*, 1990). Polychoric correlation methods are a structural parametric model for ordinal survey data. The assumption is that  $X_j$ , the data for survey question  $j$  is generated from an underlying continuous variable  $w_j$  with thresholds corresponding to each of the outcomes  $L_j$ . A correlation matrix based

on these latent continuous variables is estimated using maximum likelihood with a number of identifying restrictions. Then, the polychoric covariance matrix is used as  $S$  in the PCA procedure on the ordinal data to extract a household wealth index. The method is suitable, however, for ordinal data only and cannot handle unordered categorical outcomes. Furthermore, the methods described so far require dropping all respondents that have missing data or imputing their responses, which can involve a large loss of sample size and loss of important information, since refusing to answer certain questions could be relevant to characterizing a household. Though PCA-based approaches to summarizing categorical data are popular and simple to calculate, they each involve significant limitations. Another simple approach involves constructing a common factor by averaging (Petersen, 2006). For example, an individual-level index can be created by averaging the individual z-score for each of the responses.

We introduce an alternative Bayesian approach, which treats the survey response outcome data as random variables, and estimates directly the parameters of the joint distribution of survey responses  $p(X_1, \dots, X_N)$ . The approaches described provide a structural parametric approach to modeling multinomial data, that does not rely on ad-hoc data transformations and does not treat discrete data as continuous.

In order to estimate the parameters of the joint distribution of survey responses, some simplifying conditional independence assumptions must be made. To see why, imagine a survey dataset with  $N = 100$  respondents,  $J = 5$  questions and  $L_j = 5$  possible responses to each question for each  $j$ . For each individual,  $\prod_{j=1}^J L_j = 5^5$  parameters are required to capture all possible dependencies between questions. For a dataset with  $N = 100$  individuals, in order to capture all possible dependencies between responses for all individuals, we require a total of  $100^{5^5}$  parameters, which is an infeasibly large number. As a result, the approaches involved augment the joint distribution of observed responses with a set of random latent variables, each with their own prior distribution. A series of conditional independence assumptions involving the latent variables and the observed data are made, which allow factorization and estimation of the joint distribution of the data and latent variables. The estimated latent variables provide interpretable low-dimensional summaries of the observed data and replace the PCA-based indexes described earlier in this section. Bayesian hierarchical latent variable models have appeared recently in the economics literature, including in analysis of consumer choice (Ruiz *et al.*, 2018; Athey *et al.*, 2018), CEO management practices (Bandiera *et al.*, 2017), and text analysis of central bank communication (Hansen *et al.*, 2018).

For the remainder of this section, we introduce the three models in Table 1, which each assume multinomial data is generated from a mixture of multinomial distributions. The main difference between the models is where in the hierarchy of aggregate, group, and individual the mixture is located. Each corresponds to a different set of conditional independence assumptions involving the data and the specified latent variables.

	GoM	LDA	Latent Class
Mixture Level	Individual	Group	Aggregate
Dynamic Version		SS-M	MS-M
Profile Parameter ( $\beta$ )	$K \times J \times L_j$	$K \times P$	$K \times P$
Assignment Parameter ( $z$ )	$N \times J$	$N \times 1$	$C \times 1$
Mixture Parameter ( $g$ )	$N \times K$	$C \times K$	$K \times 1$
Mixture Hyperparameter ( $\alpha$ )	$N \times K$	$C \times K$	$K \times 1$
Profile Hyperparameter ( $\eta$ )	$K \times J \times L_j$	$K \times P$	$K \times P$

Table 1: Hierarchical Latent Variable Models

## 2.2 Grade of Membership Model

Grade of Membership (GoM) (Erosheva *et al.*, 2007), is a hierarchical latent variable model for categorical survey data that is a structural alternative to PCA-based methods. The latent variables and hyperparameters of the model with  $K$  hidden profiles and their dimensions are in Table 1. There are two sets of Dirichlet hyperparameters,  $\alpha \in \mathbb{R}^k$  and  $\eta_j \in \mathbb{R}^{L_j}$ , for the multinomial latent variables in the model. There are  $J$  multinomial distributions  $\beta_{jk} \in \Delta^{L_j-1}$  for each profile  $k$  giving profile-specific distributions over responses. There is the mixture parameter for each individual,  $g_i \in \Delta^{K-1}$ , which provides a multinomial distribution over  $K$  profiles for each individual.  $z$  is an indicator variable giving the assigned profile  $z_{ij} \in \{1, \dots, K\}$  for each individual  $i$  and each question  $j$ .

The full specification of the GoM model assumes that the discrete outcome data  $X_{ij}$  is generated as follows:

$$X_{ij}|z_{ij}, \beta_j \sim \text{Multinomial}(\beta_{j,z_{ij}})$$

$$z_{ij}|g_i \sim \text{Multinomial}(g_i)$$

$$g_i \sim \text{Dirichlet}(\alpha)$$

$$\beta_{jk} \sim \text{Dirichlet}(\eta_j)$$

The model described relies on a set of conditional independence assumptions to factorize the joint distribution. These are listed below.

**GoM 1 : Conditional Independence of Questions** *The joint probability of assigning question  $j$  to profile  $k$  and of individual  $i$  selecting response  $l_j$  is independent of that individual's other responses, given the profiles and the individual-level mixture over profiles. In addition, responses are independent of individual-level mixtures given the profile assignments, and assignments are independent of profiles given the individual-level mixtures.*

$$Pr(X_i, z_i|g_i, \beta) = \prod_{j=1}^J Pr(X_{ij}, z_{ij}|g_i, \beta)$$

$$Pr(X_{ij}|g_i, z_{ij}, \beta) = Pr(X_{ij}|z_{ij}, \beta)$$

$$Pr(z_{ij}|g_i, \beta) = Pr(z_{ij}|g_i)$$

**GoM 2 : Conditional Independence of Individuals:** *Conditional on individual  $i$ 's mixture over profiles  $g_i$ , the probability of a certain response for individual  $i$  is independent of other individuals' responses.*

$$Pr(X, z|g, \beta) = \prod_{i=1}^N Pr(X_i, z_i|g_i, \beta)$$

**GoM 3 : Independence of Profiles and Mixtures (GoM)** *Profiles and mixtures are independent of each other. In addition,  $g_i$  is independent of  $g_r$  for  $r \neq i$  and  $\beta_{jk}$  is independent of  $\beta_{mf}$  for  $f \neq k$  and  $m \neq j$ .*

$$Pr(g) = \prod_{i=1}^N Pr(g_i)$$

$$Pr(\beta) = \prod_{k=1}^K \prod_{j=1}^J Pr(\beta_{kj})$$

With these assumptions, the GoM joint likelihood factorizes as follows: First, the factorize the joint distribution into marginal and conditional distributions.

$$Pr(\beta, g, z, X) = Pr(\beta, g)Pr(X, z|\beta, g)$$

Then, apply GoM 3 and GoM 2:

$$Pr(\beta, g, z, X) = \prod_{k=1}^K \prod_{j=1}^J Pr(\beta_{jk}) \prod_{i=1}^N Pr(g_i) \prod_{i=1}^N Pr(X_i, z_i|\beta, g_i)$$

Then, apply GoM 1 and factorize the joint distribution of  $X_{ij}$  and  $z_{ij}$  into marginal and conditional distributions:

$$Pr(\beta, g, Z, X) = \prod_{k=1}^K \prod_{j=1}^J Pr(\beta_{jk}) \prod_{i=1}^N Pr(g_i) \prod_{i=1}^N \prod_{j=1}^J Pr(X_{ij}|\beta, z_{ij})Pr(z_{ij}|\beta, g_i)$$

Lastly, include the explicit probabilities from the model specification, which indicates multinomial conditional distributions for the outcome variables and the assignments.  $Pr(X_{ij}|\beta, z_{ij}) = \beta_{j, z_{ij}, X_{ij}}$  and  $Pr(z_{ij}|\beta, g_i) = g_{i, z_{ij}}$ .

$$Pr(\beta, g, Z, X) = \prod_{k=1}^K \prod_{j=1}^J Pr(\beta_{jk}) \prod_{i=1}^N Pr(g_i) \prod_{i=1}^N \prod_{j=1}^J g_{i, z_{ij}} \beta_{j, z_{ij}, X_{ij}}$$

The probabilities of the latent variables  $\beta$  and  $G$  depend on the hyperparameters  $\alpha$  and  $\eta$  and the Dirichlet density function. For example, for  $g_i$ :

$$p(g_i) = \frac{\Gamma\left(\sum_{k=1}^K \alpha_k\right)}{\prod_{k=1}^K \Gamma(\alpha_k)} g_{i1}^{\alpha_1-1} \dots g_{iK}^{\alpha_K-1}$$

The assumption that the outcome variables are multinomial is robust to both ordered and unordered categorical data. Missing data is simply included as a possible categorical outcome. GoM assumes a more natural underlying structural model for discrete data than assuming categorical variables are normally distributed or derived from thresholded continuous variables. For example, in the household survey data, a researcher could assume that there are two extreme types of households in a developing country. One is a poor household, profile 1, which has high probability of having pit toilets, low probability of owning a bicycle or car, and, low probability of having a savings account. The other is a rich household, profile 2, which has a high probability of having a flush toilet, car, and a savings account. Individual households are modeled as a mixture of these two extreme types; a wealthy household has  $g_{i2}$  close to 1, and a poor household has  $g_{i2}$  close to 0. A middle income household might have a mixture closer to the center of the simplex. This approach allows for infinite heterogeneity in household characteristics (since there are infinite points defined on the simplex over extreme profiles), while maintaining an interpretable and low-dimensional underlying structure (since mixtures of only  $K$  profiles are assumed to generate all household responses). The mixture weight  $g_{i2}$  of a household on the wealthy profile can be considered to be a probability estimate of an individual-level wealth index, and used in further economic analysis. This is a structural alternative to existing PCA-based methods for creating individual-level wealth indexes.

GoM can be estimated using gibbs sampling, or a fast variational inference procedure using the R package `mixedMem`. Due to the conjugacy of the Multinomial and Dirichlet distributions, the sampling procedure is straightforward. See Erosheva *et al.* (2007) for the full details of the sampling procedure for GoM with an added prior distribution on  $\alpha$ .

## 2.3 Latent Dirichlet Allocation Model

GoM, PCA, and polychoric correlation methods produce a decomposition of the  $N \times J$  raw survey data  $X$  into individual-level weights on profiles  $g_i$  and profile-specific weights on responses  $\beta_k$ . In some cases, though, applied economists may not require an individual-level index, but are instead interested in an group-level index. For example, Bloom & Reenen (2010) explain cross-country differences in productivity using a country-wide index of management ability, which is an average of firm-level management score in each country. The firm-level score is derived from averaging ordinal outcome variables from a firm management survey. It is always possible to estimate a group level index simply by averaging an individual-level index. However, if a group level index is the target, the reduction in parameters from estimating a group rather than individual-index directly can be significant, and can be used to offset the increase in parameters that can result from modifying some of GoM’s problematic independence assumptions.

GoM 1, conditional question independence, is not likely to hold in most survey data. There is



likely to be much more complex dependence between an individual's responses to different survey questions than can be captured by the individual-specific mixtures  $g_i$ . One way to loosen GoM 1 is to model the probability  $P(X_{i1}, X_{i2}, \dots, X_{iJ})$  directly, rather than factorizing the joint distribution into a set of  $J$  multinomial distributions conditional on the latent variable  $g$ . Each row  $X_i$  of the  $N \times J$  matrix corresponds to one permutation  $p$  of all possible survey response permutations indexed by  $\{1, \dots, P\}$ , where  $P = \prod_{j=1}^J L_j$ . The multinomial distribution over outcomes  $\beta_k \in \Delta^{P-1}$ , is  $P$ -dimensional. This is in contrast to GoM, where each profile has a separate multinomial distribution for each question  $j$  over  $L_j$  possible responses.

Assuming that the low dimensions latent variable that determines the structure of outcomes is group-level rather than individual-level means that  $g$  is now  $C \times K$  rather than  $N \times K$ . We map each individual  $i$  to group  $c$  and assign a second index to each row  $X_{ci}$  to indicate which group an individual is assigned to. Even for a moderate number of survey questions  $J$ , the number of permutations  $P$  can be quite large. However, the model considered is actually a reformulated Latent Dirichlet Allocation (Blei *et al.*, 2003), designed for the sparse, high-dimensional settings of text analysis, and able to handle large  $P$  and large  $C$ . The full specification of the model is:

$$\begin{aligned} X_{ci}|z_{ci}, \beta &\sim \text{Multinomial}(\beta_{z_{ci}}) \\ z_{ci}|g_c &\sim \text{Multinomial}(g_c) \\ g_c &\sim \text{Dirichlet}(\alpha) \\ \beta_k &\sim \text{Dirichlet}(\eta) \end{aligned}$$

The conditional independence assumptions for the model are below. GoM 1 is now eliminated entirely, since we directly estimate the probability of each possible survey response permutation for an individual. GoM 2 has now been strengthened so that group-level mixtures capture all relationships between individual's responses: the response of an individual is now independent of the response of other individuals given the group-level mixture assignments, rather than individual-level mixture assignments. GoM is maintained in an appropriately modified form.

**LDA 1 : Conditional Independence of Individuals Given Group Structure** *Conditional on group  $c(i)$ 's mixture over profiles  $g_{c(i)}$ , the probability of a certain response and assignment for individual  $i$  is independent of other individual's responses. In addition, responses are independent of group-level mixtures given the profile assignments and assignments are independent of profiles given the group-level mixtures.*

$$\begin{aligned} Pr(X, z|g, \beta) &= \prod_{i=1}^N Pr(X_i, z_i|g_c, \beta) \\ Pr(X_{ci}|g_c, z_{ci}, \beta) &= Pr(X_{ci}|z_{ci}, \beta) \\ Pr(z_{ci}|g_c, \beta) &= Pr(z_i|g_{ci}) \end{aligned}$$

**LDA 2 : Independence of Profiles and Mixtures**  $g_c$  is independent of  $g_d$  for  $c \neq d$  and  $\beta_k$  is independent of  $\beta_f$  for  $f \neq k$ .

The joint likelihood for LDA is factorized using the conditional independence assumptions:

$$p(\beta, g, Z, X) = \prod_{k=1}^K p(\beta_k) \prod_{c=1}^C p(g_c) \prod_{i=1}^N g_{c, z_{ci}} \beta_{z_{ci}, X_{ci}}$$

The estimate of  $g_c$  provides an index for group  $c$ . For the DHS survey example, an economist could assume that there are two profiles of household wealth: one with a higher probability on bundles of purchases of higher quality household assets and the other with higher probability on permutations of responses with low asset ownership. Due to country-specific economic institutions, there is a country-specific probability that a household is assigned to the profile correlated with poverty. The weight on the profile with high asset ownership can be interpreted as a country-specific wealth index.

The Gibbs sampling procedure for estimating LDA is described in detail in Griffiths & Steyvers (2004). Estimation procedures are available in R, for example using the package `topicmodels` or `lda`. As in GoM, the conjugacy of the Dirichlet and multinomial distributions results in conjugate posterior distributions and computationally efficient sampling.

## 2.4 Latent Class Models

In the preceding section, we characterized GoM as an individual-level mixture model and LDA as an group-level mixture model. It is also worth pointing out the relationship of these latent variable models to a Bayesian latent class model. A Bayesian latent class model assumes outcomes in a dataset come from a finite set of profiles, where each individual in a group is assigned to a single profile. A profile corresponds to a multinomial distribution over outcomes.

Below is the specification of a latent class model:

$$X_{ci} | \beta, z_{ci} \sim \text{Multinomial}(\beta_{z_{ci}})$$

$$z_c | g \sim \text{Multinomial}(g)$$

$$g \sim \text{Dirichlet}(\alpha)$$

$$\beta_k \sim \text{Dirichlet}(\eta)$$

In a latent class model, the mixture  $g$  is aggregate-level, so is  $K$ -dimensional, rather than  $N \times K$  as in GoM. All individuals in a group are assigned to the same profile, whereas in LDA individuals in a group can be assigned to different profiles, and in GoM each individual is assigned to a mixture of profiles. The conditional independence assumption of assignments for individuals in a latent class model are strict compared to the corresponding assumptions in GoM and LDA.

**Latent Class 1 : Conditional Independence of Assignments Given Mixture:** *Conditional*

on the aggregate mixture over profiles  $g$ , the profile assignment of group  $c$ ,  $z_c$ , is independent of the profile assignment of group  $d$ ,  $z_d$ .

$$Pr(z|g, \beta) = \prod_{c=1}^C Pr(z_c|g)$$

**Latent Class 2 : Conditional Independence of Individuals Given Assignments:** *Conditional on an individual's assignment based on their group membership, individual responses are independent. Furthermore, responses are independent of the aggregate-level mixture given the profile assignments.*

$$Pr(X|g, z, \beta) = \prod_{i=1}^N Pr(X_{ci}, |z_{ci}, \beta)$$

**Latent Class 3 : Independence of Profiles**  $\beta_k$  is independent of  $\beta_f$  for  $f \neq k$ .

The resulting joint likelihood of the model is

$$p(\beta, g, z, X) = Pr(g) \prod_{k=1}^K Pr(\beta_k) \prod_{c=1}^C p(z_c|g) \prod_{i=1}^N \beta_{z_c, X_{ci}}$$

The model assumes homogeneity within groups in a way that is not likely to hold in most survey data where groups contain more complex forms of heterogeneity. As a result, the interpretation of  $g$  is more limited than  $g_i$  or  $g_c$  in GoM and LDA: the model simply estimates the predominant patterns in the data at a aggregate-level, rather than an individual or group level profiles that can be interpreted as an index.

So far, in this section we have showed that GoM is a structural alternative to Filmer-Pritchett and Polychoric correlation-based PCA methods that does not involve unwarranted transformations of the data, and allows ordered, unordered, and missing responses in categorical data. Furthermore, GoM is part of a class of hierarchical latent variable models with mixture weights at the individual, group (LDA), or aggregate-level (Latent Class). Every Bayesian hierarchical latent variable model relies on a set of conditional independence assumptions, some of which are strong. What remains is a short discussion of identification in these models before we discuss introducing dynamics in the mixture weights for time series survey data where some of the GoM and LDA independence assumptions are not likely to hold.

## 2.5 Identification of Static Hierarchical Latent Variable Models

In the section on GoM, we described a hypothetical “rich” and “poor” profile in the context of household wealth index estimation based on household survey data. Without prior restrictions, the posterior likelihood of a model with Profile 1 as the “poor” profile and Profile 2 as the “rich” profile is the same as a model with the labels swapped. To avoid this label-swapping issue, we use the prior distribution of  $\beta_k$  to label each profile before estimation. For two profiles, each profile is assigned a permutation of survey responses that correspond to an extreme response: for example,

in the DHS example, owning every asset. Then, the hyperparameter  $\eta$  is adjusted so that  $\beta_{kj}$  has high prior probability of being close to zero for the extreme response assigned to the other profile. This can be considered as the discrete and Bayesian analogy to the lower triangular restriction that is placed on factor models for identification purposes.

In addition to the label-swapping issue, there is also the potential issue of the posterior likelihood being flat for multiple values of the model parameters, which corresponds to a frequentist notion of parameter identification. If there are multiple parameter values for which the posterior likelihood is similar, the sampling procedure may converge to different values for parameters depending on the starting point of the algorithm, and interpretation of the results is difficult. In a Bayesian framework the shape of the posterior distribution depends both on the model specified and restrictions implemented, as well as the data available. So, assessing the flatness of the peak in the posterior distribution of the parameters is best done using typical Bayesian methods for checking posterior convergence (Gelman *et al.*, 2013). For example, running chains from multiple starting values and ensuring that posterior estimates converge to similar values.

We have found that for small values of  $K$  and with the label-swapping restrictions on the prior parameters described above, that the parameter posterior distributions generally appear to be single-peaked with consistent estimated mean values for the parameters across model runs. For large values of  $K$ , more prior restrictions on the profiles may be required to avoid flat areas in posterior distributions.

### 3 Dynamic Models for Categorical Time Series

In the previous section, we described models for cross-sectional survey data, for example development wealth surveys and management surveys. There is a large class of survey data, in macroeconomics and political economy, for example, where the same survey is run repeatedly at a constant interval, each time with a different sample of individuals. Measures derived from these surveys such as the confidence indices from the Michigan Survey of Consumers and approval ratings from Gallup polls on political sentiment are frequently cited in the popular press. The models in the previous section are amenable to analysis of time series data; however, the addition of time dependence violates some of the conditional independence assumptions and must be appropriately modified. In this section, we continue to work with  $N \times J$  survey response data  $X$ . However, we also assign each individual  $i$  to the time period in which the response was observed,  $t$ . We denote  $X_t$  as the subset of  $X$  that includes only responses of individuals that were observed at time  $t$ .

We consider two ways of relaxing conditional independence: the first adds dynamics to the assignment parameters  $z$  in a dynamic version of the latent class model. This is the multinomial version of a Bayesian markov-switching model. The second method adds dynamics to the mixture parameters  $g_t$  in a dynamic version of LDA. This can be considered the multinomial version of a linear state space model.

There is a related statistics literature which has derived markov-switching models and state

	MS-M	SS-M
Gaussian Version	Markov-Switching	Lin. State Space
Profile Parameter ( $\beta$ )	$K \times P$	$K \times P$
Assignment Parameter ( $z$ )	$T \times 1$	$N \times 1$
Mixture Parameter ( $g$ )	$K \times K$	$T \times K$
Mixture Hyperparameter ( $\alpha$ )	$K \times K$	$T \times K$
Profile Hyperparameter ( $\eta$ )	$K \times P$	$K \times P$

Table 2: Dynamic Hierarchical Latent Variable Models

space models for discrete time series, mostly time series of counts (Davis *et al.* , 2016). The relationship between the models introduced here and markov switching and state space models for Gaussian data is in Table 2. There is also related work developing dynamic versions of LDA with alternatives to a conditional lognormal distribution for the latent variables for improved sampling speed and convergence (Bradley *et al.* , 2018), (Linderman *et al.* , 2015).

### 3.1 Multinomial Markov-Switching Model (MS-M)

For the first dynamic model introduced, as in the latent class model, each survey response is generated from a profile-specific multinomial distribution. The profile indicator for an individual, determined by the time at which they responded to the survey, follows a markov-switching process governed by an aggregate-level transition matrix  $g$ .

The model specification is as follows:

$$X_{ti}|\beta, z_t \sim \text{Multinomial}(\beta_{z_t})$$

$$z_t|z_{t-1}, g \sim \text{Multinomial}(g_{z_{t-1}})$$

$$g_k \sim \text{Dirichlet}(\alpha_k)$$

$$\beta_k \sim \text{Dirichlet}(\eta)$$

$g$  is now a  $K \times K$  where each column  $g_k \in \Delta^{K-1}$  gives  $Pr(z_t|z_{t-1} = k)$ .

Assumption Latent Class 1 is replaced by MS-M 1. The other conditional independence assumptions are the same as in latent class models.

**MS-M 1 : Markov Process of Assignments:** *Conditional on the assignment of the responses in the previous period  $z_{t-1}$ , the profile assignment of responses in period  $t$ ,  $z_t$ , is independent of those in other periods. In addition, assignments are independent of profiles given the transition matrix.*

$$Pr(z|g, \beta) = \prod_{t=1}^T Pr(z_t|z_{t-1}, g)$$

**MS-M 2 : Conditional Independence of Individuals Given Assignments:** *Conditional on an individual's assignment based on the time period in which the response was observed, individual*

responses are independent. Furthermore, responses are independent of the transition matrix given the profile assignments.

$$Pr(X|g, z, \beta) = \prod_{i=1}^N Pr(X_{ti}, |z_t, \beta)$$

**MS-M 3 : Independence of Profiles**  $\beta_k$  is independent of  $\beta_f$  for  $f \neq k$ .

The likelihood of the markov-switching model for categorical data is as follows:

$$Pr(X, z, \beta, g) = \prod_{k=1}^K Pr(g_k) \prod_{k=1}^K Pr(\beta_k) \prod_{t=1}^T p(z_t | z_{t-1}, g) \prod_{i=1}^N \beta_{z_t, X_{ti}}$$

$Pr(\beta_k)$  and  $Pr(g_k)$  are Dirichlet densities and  $Pr(z_t | z_{t-1})$  is directly from the transition probability matrix. The posterior distribution of the model parameters are estimated via Gibbs Sampling using the below steps, which are described in full detail in Appendix A.

1. Generate  $z_t$  conditional on  $g$ ,  $X$ ,  $z_{t+1}$ ,  $z_{t-1}$ , and  $\beta$  using a posterior multinomial distribution and Albert & Chib (1993)'s single-move sampling procedure.
2. Generate transition matrix  $g$  conditional on  $z$  using a posterior Dirichlet distribution.
3. Generate  $\beta$  conditional on  $X$ ,  $z$  with a posterior Dirichlet distribution.

We call this model the multinomial markov switching model (MS-M), since it is equivalent to a Bayesian markov-switching model for Gaussian models with state-specific means replaced by state-specific multinomial distributions. The time dependence between responses at time  $t$  and responses at time  $t - 1$  is captured through the markov-switching process of  $z_t$ . It is computationally simple to estimate and has a convincing interpretation when analyzing discrete data where it is natural to assume a discrete-valued latent variable can influence responses patterns over time (for example, recessions can generate switching patterns in responses to consumer confidence surveys).

However, as in the latent class model, it is limiting in the form of heterogeneity that it can capture in responses for individuals across time. All responses at time  $t$  are assumed to come from the same profile indexed by  $z_t$ . This is unduly restrictive in many settings where it is more natural to think about evolution in response patterns as coming from proportions of respondents of different types changing, rather than nonlinear switching in the state of the entire survey sample. The next model addresses this by adding time dynamics to a time-specific mixture parameter  $g_t$  instead, which captures time dependence in responses while allowing a more flexible form of heterogeneity in responses in each time period.

### 3.2 Multinomial State Space Model (SS-M)

The second model introduced is a dynamic version of LDA. Like LDA, responses follow a profile specific multinomial distribution and an individual's profile is generated from a group-specific mixture, where groups are indexed by time  $t$  rather than group  $c$ . The group, or time-specific mixture

now follows a lognormal distribution. In static LDA  $g_t \in \Delta^{K-1}$ . In the dynamic version  $g_t \in \mathbb{R}^K$ , follows a random walk process and is converted to a vector of proportions  $\psi(g_t) \in \Delta^{K-1}$  using the softmax function

$$\psi(g_t) = \frac{\exp(g_t)}{\sum_{k=1}^K \exp(g_{tk})}$$

The model specification is:

$$\begin{aligned} X_{ti}|z_{ti} &\sim \text{Multinomial}(\beta_{z_{ti}}) \\ z_{ti} &\sim \text{Multinomial}(\psi(g_t)) \\ g_t &= g_{t-1} + w_t, \quad w_t \sim N(0, \sigma I) \\ \sigma &\sim \text{InverseGamma}(v0, s0) \\ \beta_k &\sim \text{Dirichlet}(\eta) \end{aligned}$$

In the MS-M model, we added dynamics to the discrete-valued assignments  $z$ , and the dynamics were captured using a discrete markov process. In this model, the dynamics are added to the group-specific mixtures, and are represented by a linear markov process. The model is denoted the multinomial state space model, since it involves adding a multinomial outcome to the dynamics of a Bayesian linear state space model via a link function.

Since when we choose to group individuals by time  $t$ , rather than by potentially independent group membership  $c$ , LDA 2, the independence of group mixtures no longer holds, and is adjusted. We maintain the assumption of conditional independence of individuals given the group structure.

**SS-M 1 : Conditional Independence of Individuals Given Group Structure** *Conditional on time  $t$ 's mixture over profiles  $g_t$ , the probability of a certain response and assignment for individual  $i$  is independent of other individual's responses. In addition, responses are independent of time-specific mixtures given the profile assignments and assignments are independent of profiles given the time-specific mixtures.*

$$Pr(X, z|g, \beta) = \prod_{i=1}^N Pr(X_{ti}, z_{ti}|g_t, \beta)$$

$$Pr(X_{ti}|g_t, z_{ti}, \beta) = Pr(X_{ti}|z_{ti}, \beta)$$

$$Pr(z_{ti}|g_t, \beta) = Pr(z_{ti}|g_t)$$

**SS-M 2 Markov Property of  $g_t$  and Independence of Profiles**  $\beta_k$  is independent of  $\beta_f$  for  $f \neq k$ . Given  $g_{t-1}$ ,  $g_t$  is independent of  $g_s$  conditional on  $g_{s-1}$  for  $s \neq t$ .

With these assumptions, the joint probability distribution of the SS-M model factorizes as follows:

$$p(\beta, g, Z, X) = \prod_{k=1}^K p(\beta_k) \prod_{t=1}^T p(g_t | g_{t-1}) \prod_{i=1}^N g_{t, z_{ti}} \beta_{z_{ti}, X_{ti}}$$

The gibbs-sampling steps necessary to estimate the model are as follows, and are explained in detail in Appendix A:

1. Generate  $g_t$  conditional on  $\sigma$ ,  $g_{t-1}$ ,  $g_{t+1}$ , and  $z_t$  using Stochastic Gradient Langevin Dynamics
2. Generate  $z$  conditional on  $\beta$ ,  $X$ , and  $g$  using a multinomial posterior
3. Generate  $\beta$  conditional on  $z$  using a Dirichlet posterior
4. Generate  $\sigma$  conditional on  $g$  using an Inverse-Gamma posterior

Step 1 is computationally intensive and can be difficult to tune. Dynamic LDA (Blei & Lafferty, 2006), is a similar model that allows evolution over time both in the profile mixtures  $g_t$  as well as the actual profiles  $\beta$ . In this paper, however, we are interested in identifying latent profiles in survey respondents that are constant over time; we assume that differences in patterns of survey responses over time is due to changes in prevalence of each latent profile,  $g_t$ , rather than changes in the latent profiles themselves.

In this section, we described a formulation where we assume that at each time period there is a mixture of dominant profiles. We will show that in many applications this assumption is effective in obtaining interpretable and economically meaningful latent states. For example, in the Michigan Survey of Consumers setting, it is natural to assume that there are always both optimistic and pessimistic people in the survey sample each month. Depending on general economic conditions and media reports, there are different proportions of optimistic people in each month. Being optimistic involves a different multinomial distribution over permutations of survey responses compared to being pessimistic.

## 4 Example for Static Models

Using the static models described, we explore the low-dimensional structure of a dataset of 187,616 survey responses from the Demographic and Health Surveys in 5 Caribbean and Latin American countries from 2009-2012. Aggregate economic measures from developing countries from data sources like the Penn World Tables are based on very little hard data (Young, 2012). The DHS survey data, which includes millions of survey responses on demographics, health, and economic outcomes from dozens of developing world households since the 1990s, has become a crucial source of data for deriving consumption and income estimates for households in the developing world. The World Bank publishes wealth indices derived from Filmer-Pritchett type PCA on a variety of asset indicators from the survey; a variety of other researchers have derived their own PCA-based socioeconomic status indicators from the data, for example to train machine learning models to associate local income measures with satellite images (Jean *et al.*, 2016).



The survey data used in our example include variables on water quality and house floor quality (categorized from 1 to 4), toilet quality (scored from 1 to 5) and binary variables on ownership of electricity, radio, tv, fridge, motorbike, car and phone. The variables are mapped from the raw data according to the table in Appendix B. The mean response for each of these variables for each country is in Table 3. Colombia is the wealthiest country by most asset-based measures, and Nicaragua and Haiti are the poorest.

Country	water	toilet	floor	electric	radio	tv	fridge	motorbike	car	phone
Colombia (CO)	3.50	4.66	2.86	0.95	0.72	0.87	0.71	0.22	0.10	0.90
Guyana (GY)	2.78	3.95	2.92	0.71	0.55	0.73	0.55	0.09	0.14	0.80
Haiti (HT)	2.52	2.46	2.15	0.30	0.51	0.23	0.08	0.07	0.04	0.75
Nicaragua (NI)	3.24	1.75	1.52	0.23	0.53	0.19	0.06	0.15	0.04	0.56
Peru (PE)	3.45	3.44	2.24	0.82	0.84	0.73	0.33	0.13	0.09	0.73

Table 3: Mean Response for Each Country

## Household-Level Index

Most economists would agree that the asset ownership indicators are correlated with measures of household wealth and expenditure, but there is little agreement on how best to extract measures of wellbeing from noisy and heterogeneous survey response data. First, we pool responses across countries and standardize by the sample mean and standard deviation for each response. Then, we estimate one individual index which is the simple average of an individual’s Z-scores across all responses, and three indices based on weighted averages, with weights derived from three forms of PCA used most often in this setting (PCA on untransformed data, Filmer-Pritchett, and PCA based on polychoric correlation). Then, we illustrate the benefits and drawbacks of these simple approaches compared to using the probabilistic approach of GoM to extract an individual-level index.

Table TBD provides a table of loadings for each of the PCA methods, and Figure 1 shows the multinomial distributions over each question for each of the two profiles estimated by GoM. The three PCA-based methods estimate a very similar first principal component with positive loadings on all variables. The interpretation is simple but limited: a positive household index means that individual has higher weighted asset ownership than the average household in the sample. The interpretation of the weights profiles for GoM is more detailed than a weighted average of the responses to the survey. A household with the estimated  $g_{i2}$  close indicates that a home likely has high water quality, and has a small but positive probability of having a motorcycle or a car. A household with a weight close to 0 on Profile 2 indicates that while the home might have a radio, it is not likely to have good floor quality or fridge. A medium income household with good water and floor quality but no fridge or phone might have a weight on Profile 2 closer to 0.5.

In Figure 2, the histograms for the household-level indices for Nicaragua and Guyana are plotted for both PCA on the untransformed responses and GoM. The first principal component in Nicaragua

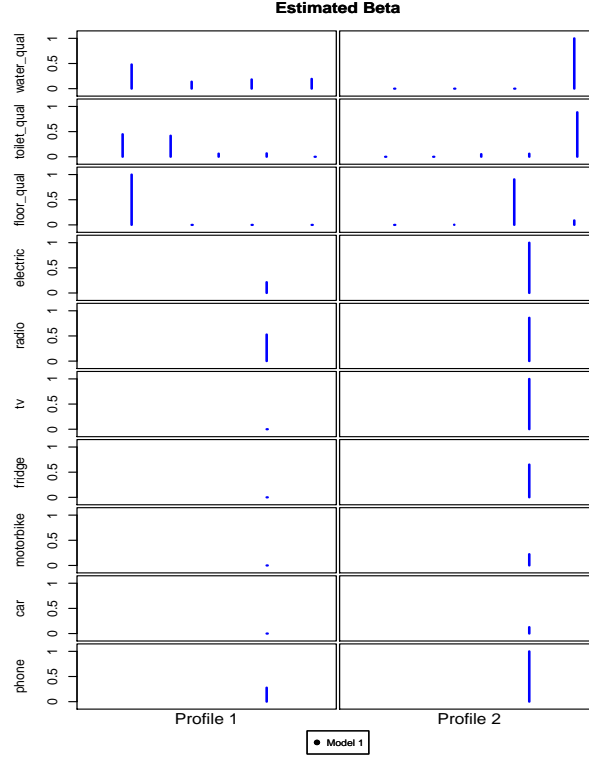


Figure 1: GoM Estimated Profile-Specific Distributions with  $K = 2$

has a right skew, and in Guyana has a left skew; for Guyana, many individuals have average or above average asset ownership, and in Nicaragua many individuals are very poor. For GoM, there is a similar shape to the distribution of household indices. Guyana has a GDP/capita in the middle of the countries in the sample: most individuals are assigned to profile 2, but there is a significant proportion who are poorer and have a higher weight on Profile 1. In Nicaragua, the poorest country in the western hemisphere, most individuals have 0 weight on the wealthier profile, but there is a small proportion of households who are categorized as wealthy.

The benefits of the z-score and PCA based approaches are their simplicity; they are easy to compute, and a higher score means that an individual has higher than average asset ownership in the pooled sample. But, further interpretation and inference is limited. The first principal component for the PCA-based methods explains only 35% of the variance in the data. Increasing the dimension of the index by adding principal components increases the proportion of variance explained, but at the cost of losing clear interpretability of the index.

GoM is more computationally expensive, and doesn't give as smooth of an index as PCA does, since for GoM many household weights are concentrated at 0 and 1. As shown in Figure 1, the posterior estimate of  $g_{ik}$  on the profile corresponding to higher asset ownership, like the PCA methods, does provide a good indicator of the wealth of a household. It is also important to point out that  $g_{ik}$  is the posterior estimate of a probability model, and has standard errors that can be used to conduct inference: for example, it can be used to test whether the latent variable estimate

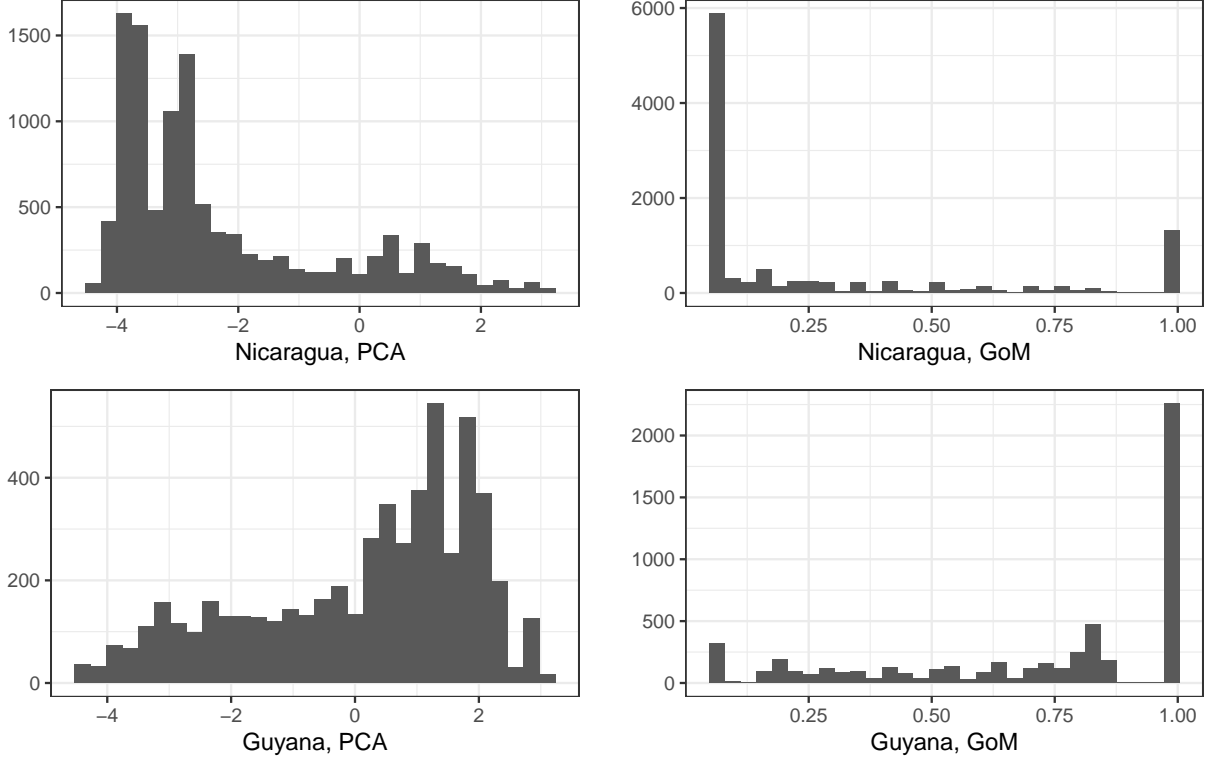


Figure 2: Distribution of Household-Level Indices for GoM vs. PCA

for household  $i$  compared to household  $j$  is equal or not. Furthermore, the GoM index retains interpretability as  $K$  increases, whereas interpretability becomes difficult for PCA as the number of components increase.

To show this explicitly, Figure 3 plots the posterior estimate of  $\beta_{jk}$  for a model with  $K = 4$  profiles and the same 10 questions. The rich heterogeneity of households in the datasets is even more apparent when estimating the model with four profiles. Profile 2 and Profile 4 are the wealthier profiles. A household with a high weight on profile 4 has a high chance of having high water, toilet, and floor quality, and owning all assets, except a car. Profile 2 has high quality water infrastructure and likely owns cheap assets, but does not have a high probability of owning the more expensive assets surveyed. Profile 3 is a rising income household with mixed quality toilet and water and moderate probability of owning cheaper assets. Profile 4 corresponds to a poor household with no assets and the lowest quality infrastructure. The estimated GoM household-specific indices retain a straightforward visualization and interpretation as the latent variable increases in dimension.

### Country-Level Index

With a representative sample, or appropriate sample weights, household data can be aggregated into country-level statistics, which can be useful for countries where there is little GDP and investment data necessary to back out aggregate household consumption (Young, 2012). The five household-level indices discussed in the previous section, for example, can be averaged by country to derive a

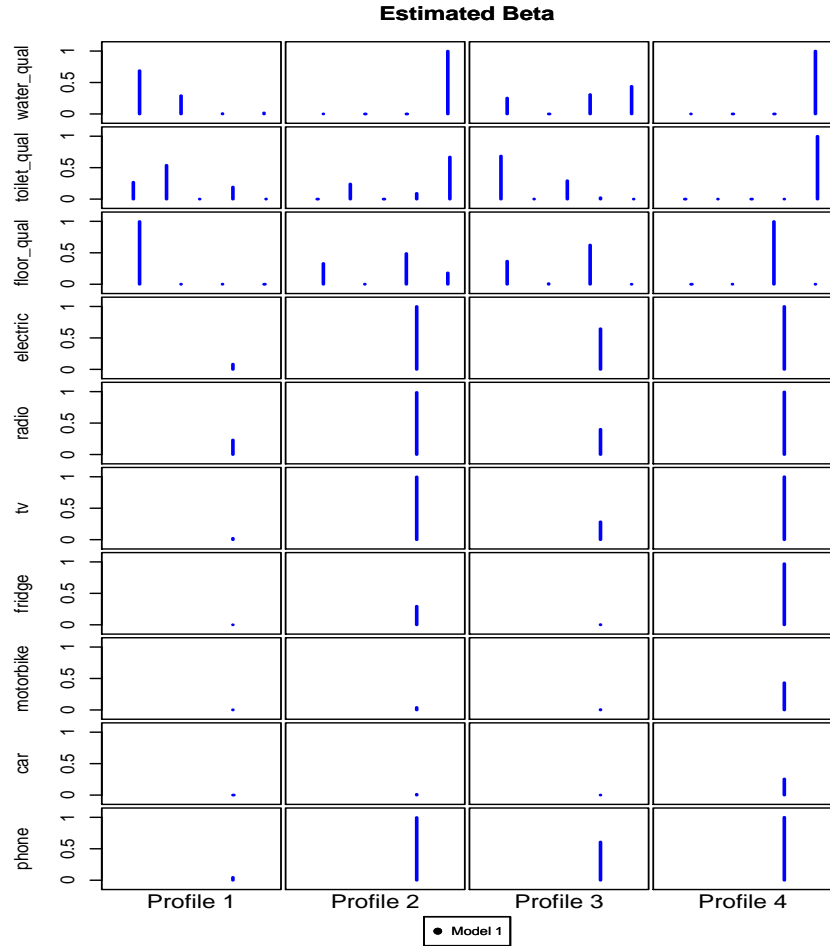


Figure 3: GoM Estimated Profile-Specific Distributions with  $K = 4$

country-level wealth index. Those are in Table 4, along with the estimated country-level indices  $g_c$  directly from LDA and LCA models estimated on the survey response data with  $K = 2$ . The PCA and z-score indices indicate which countries have higher average asset ownership relative to other countries in the sample. The GoM index indicates which countries have a higher average weight on the profile corresponding to higher asset ownership. All methods indicate that Haiti and Nicaragua are the poorest countries, and Colombia is the richest.

Country	pca	fpit	pchlor	zscore	gom	lda	lca
Colombia	1.10	1.20	0.62	1.02	0.86	1.00	1.00
Guyana	0.13	-0.28	0.05	0.30	0.71	0.55	0.00
Haiti	-1.87	-2.42	-1.16	-0.81	0.35	0.01	0.00
Nicaragua	-2.33	-2.48	-1.30	-0.86	0.27	0.01	0.00
Peru	-0.08	-0.03	-0.03	0.35	0.66	0.27	0.00

Table 4: Country-Level Socioeconomic Status Indicators

The LDA index identifies the predominant sets of bundles owned by individuals in the sample, and classifies each country as a mixture over two different multinomial distributions over bundles. Profile 1, identified as the wealthy profile, is the one that has a lower probability on owning no assets. By relaxing GoM 1 in return for strengthening GoM 2, LDA provides a description of which permutations of survey responses are common in the sample, and how the patterns vary within each country. GoM indicates the marginal probability of owning each asset or having an asset of certain quality, given a household’s assignment to a certain profile mixture. LDA assumes households are a member of only a single profile, but by estimating the joint probability of survey responses directly, allows conditional probability queries that are not informative in GoM: for example, the posterior estimates for  $\beta_k$  can be used to address the following question: given a household is in a certain country, what is the probability of having a fridge given a household has electricity? In GoM, the probability of having a fridge is constant given the household specific mixture  $g_i$ .

LCA assigns each country to a single profile that is a probability distribution over bundles of assets. The latent class model isolates Colombia as the wealthy country but lumps together all the rest in the same class. This structural model is limiting, since even among this small sample of countries, there is more complex heterogeneity than LCA can capture with a finite number of classes. Allowing there to be a country-specific mixture, as in LDA, shows that Guyana and Peru to have a significant proportion of individuals who own more expensive bundles, and have a very different household wealth distribution compared to Haiti and Nicaragua. Static hierarchical latent variable models, although computationally expensive compared to PCA and averaging-based methods, estimate similar wealth indices for households and countries, with the benefit of clear interpretation, even as the dimension of the latent variable increases.

## 5 Application of Dynamic Models

The static hierarchical latent variable models are well-explored in the CS and biostatistics literature, but are not familiar to economists. The previous section introduces and shows why they are a compelling alternative or complement to PCA-based methods for summarizing multivariate survey data. This section shows how the basic hierarchical latent variable models can be easily extended to other settings of interest to economics, such as time series data.

First, we estimate SS-M and MS-M on responses to the Michigan Survey of Consumers. We show that the estimated index from the SS-M model corresponds closely to the existing ICS, and that certain recessions are better characterized by switching in consumer sentiment than switching in real expenditure data. Second, we estimate SS-M on GSPSS responses on health, environment, and values and morals. During the 2008-2009 recession, which was characterized by a negative switch in consumer sentiment, concern towards the environment drops steeply, like the Michigan ICS and an index created from the GSPSS health survey. Sentiment toward traditional values, on the other hand, seems unaffected by financial insecurity.

This is related to existing work linking survey data to macroeconomic conditions. Scruggs & Benegal (2012) links poor economic conditions to the decline in public concern about the environment in the U.S. and EU during the Great Recession. Nicholson & Simon (2010) uses Gallup survey data on wellbeing to examine the relationship between the aggregate unemployment rate and health outcomes. Dominitz & Manski (2003) have examined the accuracy of Michigan survey data and concluded that asking questions about subjective probabilities is more accurate than the qualitative questions used as data in this paper. We take the stance that there is good quantitative information available in the patterns of aggregate qualitative responses, but show an improved method of extracting that information.

### 5.1 Michigan Consumer Sentiment Indices

The raw Michigan Survey of Consumers data contains the survey responses for 500 telephoned respondents in continental U.S. each month of the year. For each month's sample, an independent cross-section sample of households is drawn, and some are reinterviewed six months later. The total sample for any one survey is normally made up of about 60% new respondents, and 40% being interviewed for the second time. We don't model the repeated interviewing that occurs in the Michigan data.

The data for every month from January 1978 to November 2017 is publicly available. The Index of Consumer Sentiment (ICS) is made up of five questions of the survey on the respondent's opinion about current and future economic conditions. The questions along with possible responses are in Appendix B. The ICS is constructed from the five questions listed in Appendix B based on the following ad-hoc procedure, published by the University of Michigan on their website:

1. Compute the relative scores  $X_j$  for each of the five questions. The relative scores are the percentage of respondents giving favorable replies minus the percent giving unfavorable replies,

plus 100.

2. Round each relative score to the nearest whole number
3. Using the following formula, sum the relative scores, divide by the 1966 base period value, and add 2.0 to correct for sample design changes in the 1950s.

$$ICS = \frac{X_1 + X_2 + X_3 + X_4 + X_5}{6.7558} + 2.0$$

Under this procedure, responses that are missing or incomplete are dropped. With our procedure, it is not necessary to have adjustments for missing or incomplete data, as refusing to answer a question can be modeled as part of the permutations of potential survey responses. We also do not need to make any adjustment for outliers. Furthermore, for ordered variables, we make no assumption on if the the distance from a responses categorized as 1 and 2 is the same as the distance between 4 and 5.

Each month  $t$  has  $N_t$  survey respondents.  $X_{ti}$  is the response for the  $i$ -th survey respondent in month  $t$  and corresponds to a survey response permutation indexed by  $p \in \{1, \dots, P\}$ .  $X_{ti} = 11111$  corresponds to a survey respondent who answers “better” or “good” to all of the questions on the U.S. economy. We estimate SS-M on the Michigan Consumer survey data and compare the estimated weight  $g_{t1}$  on the more optimistic profile to the ICS, rescaled to match the range and variance of the new index, which is between 0 and 1. We identify the “positive” sentiment using the prior distribution on  $\beta$ , by decreasing the prior probability of responding 55555 for Profile 1.

Though the ad-hoc method used to create the ICS and the sampling procedure to estimate our index are very different, the trends in the two indices are very close, as seen in Figure 4. The index based on SS-M tends to drop more during recessions. The Michigan data has ordered responses so it is relatively straightforward to come up with an ad-hoc method to combine the responses in each month using the percentage of favorable versus unfavorable responses. In this setting, both indexes capture similar low-dimensional representations of the five survey questions. SS-M, however, specifies a probability model that generates the index, which allows additional interpretation of the index. Table 5 contains the responses with the top 5 weights for each multinomial distribution parameter  $\beta_k$  corresponding to the two estimated states.

State 1		State 2	
11111	11.1%	53551	3.9%
13111	8.4%	53555	3.5%
33111	5.1%	55555	3.0%
31111	3.1%	55551	2.8%
51111	2.4%	13551	2.8%

Table 5: Top 5 State-Specific Multinomial Probabilities

When  $g_t$  has a higher weight on State 1, then there is a higher probability of survey respondents in a month selecting all optimistic answers to the survey (11111), or a mix between indifferent and

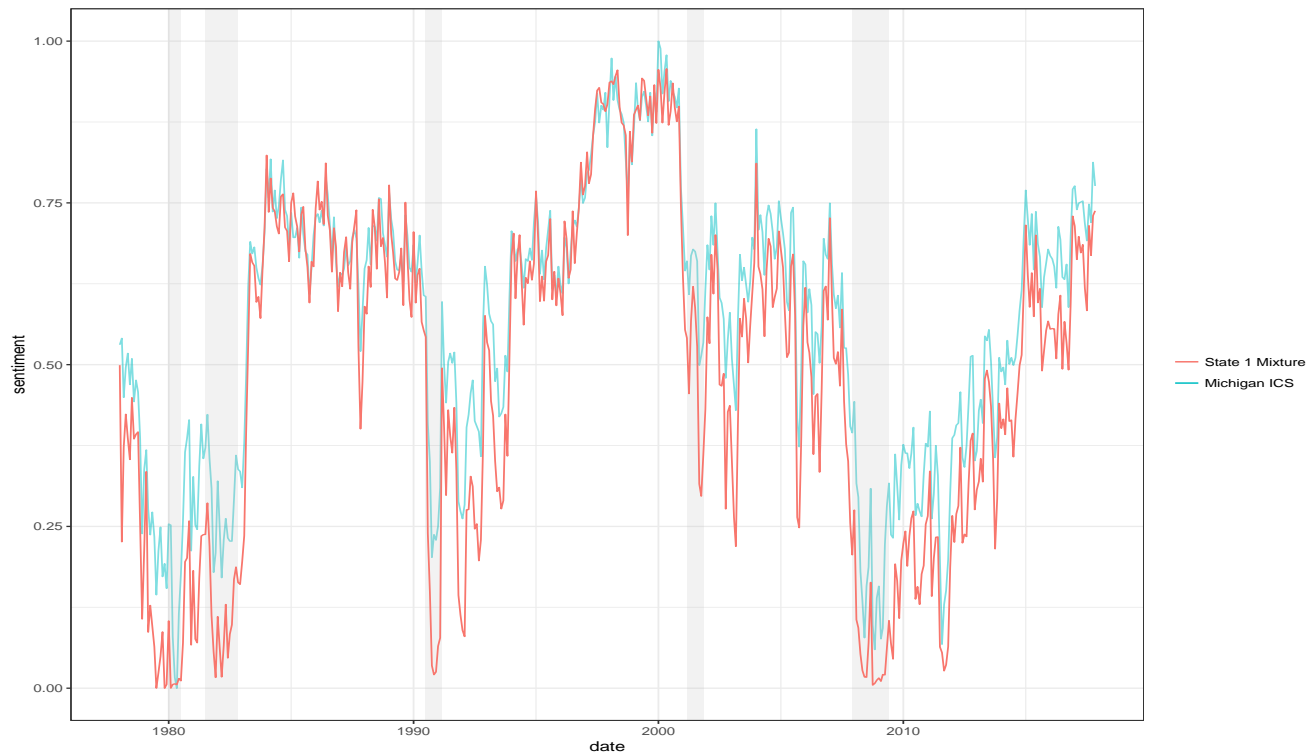


Figure 4: Probability of State 1 vs. ICS

optimistic (for example, 33111). When there is a higher weight on State 2, however, there is a higher probability of many respondents selecting all pessimistic answers to the survey (55555).

We also estimate MS-M with  $K = 2$  on the Michigan Consumer Survey data with two states. The latent variable  $z_t$  is discrete and has dynamics corresponding to a transition matrix. Table 6 gives the top 5 survey permutations and their probabilities for each of the estimated states in the model.

State 1		State 2	
11111	6.8%	11111	3.6%
13111	5.1%	53551	3.3%
33111	3.1%	13111	3.0%
13551	2.0%	53555	2.8%
31111	2.0%	55555	2.5%

Table 6: Model 1: Top 5 State-Specific Multinomial Probabilities

SS-M model has a more complex and computationally intensive estimation procedure than MS-M. Examining the difference between Table 6 and Table 5 indicates one advantage of allowing there to be a mixture of states in each period rather than a single state in each period. In both tables, State 1 corresponds to more positive sentiment towards current and future economy, while State 2 corresponds to a more negative sentiment. However, in all periods, there are many respondents responding “11111” to the survey. State 2, in the markov-switching model, while having more weight



on pessimistic response permutations, also has weight on the common permutations that appear in every month, like 13111. The mixture of states in the SS-M allows a cleaner separation between the optimistic group of respondents and pessimistic group of respondents. In SS-M, there always non-zero weight on State 1, so State 2 corresponds more directly to the group of respondents that is concerned about the economy, which fluctuates with recessions and general economic conditions.

In the discrete markov switching model, the states are highly persistent: the estimate of the transition probabilities between states is  $P(z_{t+1} = 1|z_t = 1) = 0.89$  and  $P(z_{t+1} = 2|z_t = 2) = 0.80$ . Figure 5 shows the probability of being in State 2, the more pessimistic state. The Michigan ICS is also plotted and recessions shaded.

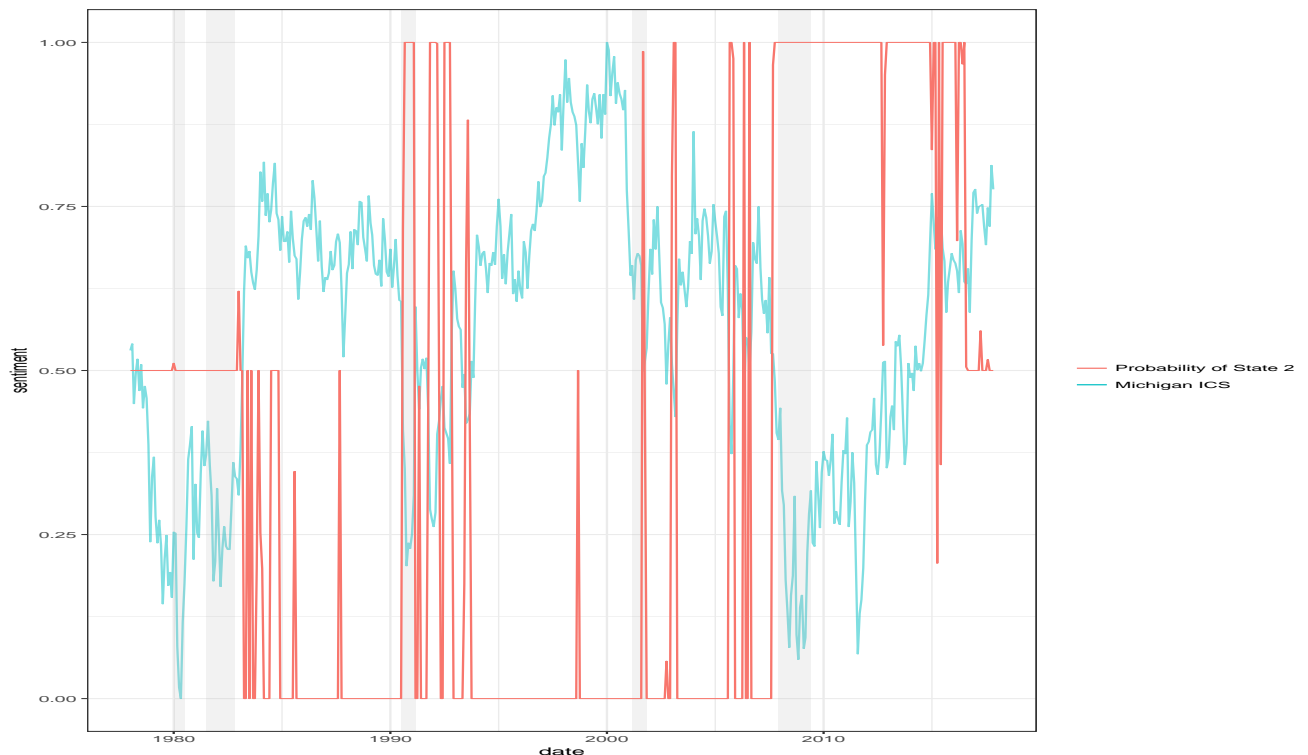


Figure 5: Probability of State 2

The probability of State 2 remains high for long after the 2008-2009 recession, but for other recessions drops quickly back to 0. This shows how consumer confidence remained low for a long time after the recession was declared over, which is unique to the 2008-2009 recession compared to other recessions in the sample period. This characteristic of the 2008-2009 recession is not apparent when examining trends in Personal Consumption Expenditure.

We estimate a standard markov mean-variance switching model on the log-difference of Personal Consumption Expenditure in the United States. The below figure shows the probability of being in the negative state for the model derived from survey data as well as the probability of the low mean state in the Gaussian markov-switching model derived from PCE data. As mentioned previously, the markov switching model on the survey data remains in the pessimistic state post-2008 recession

until recently, but the model on the expenditure data classifies a return to growth state. The model estimated on PCE data tends to lag recessions and misses the 1990s recession. The model estimated on the discrete sentiment data leads the 2008 recession and captures the 1990s recession well, but misses the 2001 recession. Both models have some difficulty classifying the beginning of the period.

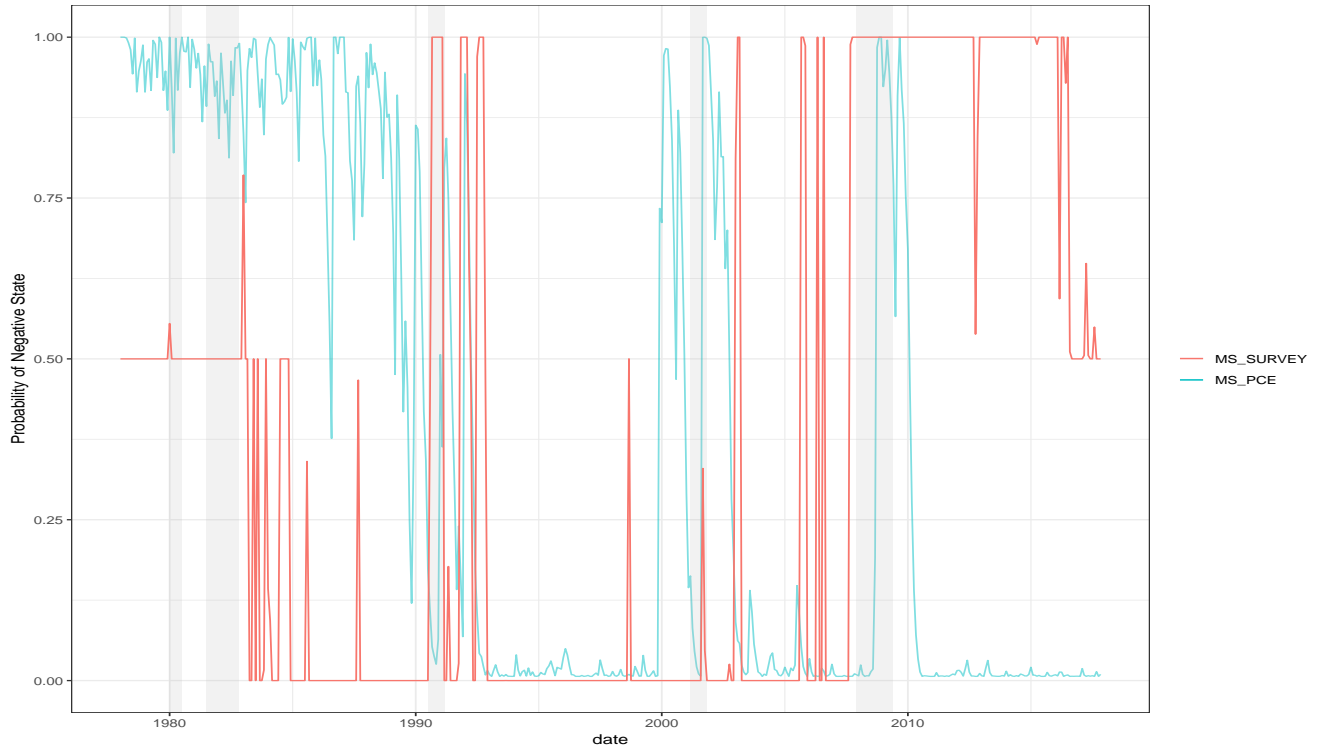


Figure 6: Probabilities of Negative State in MS-M vs. MS on PCE

Most authors have found that the consumer sentiment data does not have predictive value for real economic data like consumption expenditure beyond other macroeconomic variables (Ludvigson, 2004). However, we have shown using a markov-switching model estimated directly on the survey data responses that despite not having general independent predictive power, the survey data still contains distinct economic insights in the nature of certain recessions since 1980 compared to real economic data on consumer purchases.

In the Michigan ordinal survey data setting, where there is an existing ad-hoc sentiment index already constructed, we show that the latent variable  $g_t$  in the discrete space model can be interpreted as a sentiment index with a very similar trend to the existing ICS. Furthermore, both this continuous-valued index and the discrete-valued index from the markov-switching model have additional interpretation, since each corresponds to a profile in the model that involves a multinomial distribution over survey response permutations. This motivates extending these methods to quantify survey data that correspond to consumer characteristics like environmental concern and values where, unlike for consumer purchasing, there is little quantitative data to compare to.

## 5.2 Indexes from Gallup Data on Environment, Health, and Values

The raw Gallup Poll Social Series data contains the survey responses of a minimum of 1000 people in the U.S. each month. The core questions are devoted to a different topic each month, and are available from 2001 to 2016 through a subscription. The analysis in this section uses data from the survey on the environment conducted each March, the survey on health conducted each November, and the survey on values and morals done each May. The survey questions and responses used to create the environmental concerns index, health concerns index, and traditional values index are in Appendix B. The Gallup poll data is mapped into the model framework in the same way as for the Michigan data.

The GPSS data contains a variety of qualitative questions on important aspects of individual's economic condition that may not be well-measured in aggregate economic statistics. For example, consumer's concern towards the environment and their values are aspects of their preferences that affect their purchasing and other life decisions. Unlike consumption expenditure, there are not continuous-valued real economic statistics corresponding to environmental concern. The wealth of discrete data in survey series like Gallup contains important information about trends in these sentiments. However, the questions are not uniform. Some are ordered and others unordered, and there are differing scales and number of responses for each question. We show that the model presented in this paper can address these challenges and extract useful quantitative representations of health, environmental concern, and values during the last recession.

First, we estimate SS-M with two profiles, using the six questions in Appendix B on the respondents' views on current and future environmental conditions in the U.S. The weight on the first profile is plotted, which has higher weight on response permutations that correspond to the respondent thinking that the environment is poor, that global warming is already happening, that the threat of global warming is under exaggerated, and that environment should be prioritized over energy and economic security. The environmental concerns sentiment drops significantly during the recession when respondents were likely be distracted based on financial concerns, and was very slow to recover even to pre-2005 levels. Rather than examining trends in responses on a question by question basis, treating the responses as outcomes of a probability model allows estimation of a single factor corresponding to environmental concern.

We derive indices from sets of questions from a few other GPSS categories to compare. In Figure 7, we plot the weight on the profile that has higher probability on responses to the moral and values survey that correspond to traditional values (i.e. abortion/death penalty/suicide is immoral and religion is important). The weight on the profile corresponding to traditional values slowly decreases from over 90 percent in the early 2000s to under 85 percent in 2016, without showing any break in the trend during the recession.

Estimated sentiment based on satisfaction with personal health and the U.S. health care system, though it does not fluctuate as much as the sentiment towards the environment, does bottom out in 2009. During times when it is more difficult to afford health care, dissatisfaction with the health care system is higher. It is interesting that sentiment toward the environment follows trends in



Figure 7: Indexes Derived from GPSS Data

financial and health insecurity, and not the path of more fixed sentiments on values.

We find that SS-M can be used on the 15 disparate questions listed in Appendix B to successfully extract three low-dimensional continuous-valued representations on environmental concerns, health concerns, and traditional values. As in the literature, we find that recessions correspond to decreases in environmental concern and health satisfaction.

## 6 Conclusion

We derive Bayesian approaches to extracting interpretable discrete and continuous-valued indexes from high-dimensional multinomial time series. The models provide a structural parametric form for extracting low-dimensional summaries of high-dimensional discrete data, in contrast to existing methods for summarizing categorical data which tend to use ad-hoc averaging methods or PCA. The models successfully extract interpretable indexes from ordered and unordered, categorical and binary, and static and time series survey data on economic well-being, values, health, and environmental sentiment.

The relation of these models to the extensive computer science literature on Bayesian hierarchical latent variable models suggests a variety of extensions for the models. For the survey data application, it is possible to introduce more complex forms of dynamics, or to directly model the relationship between responses and some outcome variable (Mcauliffe & Blei, 2008), or to inte-

grate some of the parameter estimates into causal analysis (Wang & Blei, 2018). It would also be straightforward to extend the models presented for survey data that had a mixture of continuous and discrete responses.

Along with developing survey data specific examples, this paper also describes the close relationship between of a variety of seemingly disparate models, developed originally for a variety of different types of data, showing that they are all similar hierarchical latent variable models with different conditional independence assumptions. Understanding these models more fully is important, as they have applications to estimating structural models of a variety of high-dimensional discrete economic data, including text data, network data, and consumer choice.

## References

- Albert, James H, & Chib, Siddhartha. 1993. Bayes Inference via Gibbs Sampling of Autoregressive Time Series Subject to Markov Mean and Variance Shifts. *Journal of Business & Economic Statistics*, **11**(1), 1–15.
- Athey, Susan, Blei, David, Donnelly, Robert, Ruiz, Francisco, & Schmidt, Tobias. 2018. Estimating Heterogeneous Consumer Preferences for Restaurants and Travel Time Using Mobile Location Data. *AEA Papers and Proceedings*, **108**, 64–67.
- Bandiera, Oriana, Hansen, Stephen, Prat, Andrea, & Sadun, Raffaella. 2017. Ceo Behavior and Firm Performance. *NBER Working Paper Series*.
- Bentler, P M, Lee, Sik-Yum, & Poon, Wai-Yin. 1990. Full Maximum Likelihood Analysis of Structural Equation Models with Polytomous Variables. *Statistics and Probability Letters*, **9**, 91–97.
- Bhadury, Arnab, Chen, Jianfei, Zhu, Jun, & Liu, Shixia. 2016. *Scaling up Dynamic Topic Models*.
- Blei, David M, & Lafferty, John D. 2006. Dynamic Topic Models. *Pages 113–120 of: Proceedings of the 23rd international conference on Machine learning*.
- Blei, David M, Ng, Andrew Y, & Jordan, Michael I. 2003. Latent Dirichlet Allocation. *Journal of Machine Learning Research*, **3**, 993–1022.
- Bloom, Nicholas, & Reenen, John Van. 2010. Why Do Management Practices Differ Across Firms and Countries? *Journal of Economic Perspectives*, **24**(1), 203–224.
- Bradley, Jonathan R., Holan, Scott H., & Wikle, Christopher K. 2018. Computationally Efficient Multivariate Spatio-Temporal Models for High-Dimensional Count-Valued Data. *Bayesian Analysis*, **13**(1), 253–310.
- Davis, Richard A, Holan, Scott H, Lund, Robert, & Ravishanker, Nalini. 2016. *Handbook of discrete-valued time series*. CRC Press.
- Dominitz, Jeff, & Manski, Charles F. 2003 (August). *How Should We Measure Consumer Confidence (Sentiment)? Evidence from the Michigan Survey of Consumers*. Working Paper 9926. National Bureau of Economic Research.
- Erosheva, Elena A, Fienberg, Stephen E, & Joutard, Cyrille. 2007. Describing Disability Through Individual-Level Mixture Models for Multivariate Binary Data. *The Annals of Applied Statistics*, **1**(2), 502–537.
- Filmer, Deon, & Pritchett, Lant H. 2001. Estimating Wealth Effects Without Expenditure Data - or Tears. *Demography*, **38**(1), 115–132.

- Gelman, Andrew, Stern, Hal S, Carlin, John B, Dunson, David B, Vehtari, Aki, & Rubin, Donald B. 2013. *Bayesian data analysis*. Chapman and Hall/CRC.
- Griffiths, Thomas L., & Steyvers, Mark. 2004. Finding scientific topics. *Proceedings of the National Academy of Sciences*, **101**(suppl 1), 5228–5235.
- Hamilton, James D. 1989. A New Approach to the Economic Analysis of Nonstationary Time Series and the Business. *Econometrica*, **57**(2), 357–384.
- Hansen, Stephen, McMahon, Michael, & Prat, Andrea. 2018. Transparency and Deliberation within the FOMC: A Computational Linguistics Approach. *Quarterly Journal of Economics*, 801–870.
- Jean, Neal, Burke, Marshall, Xie, Michael, Davis, W. Matthew, Lobell, David B., & Ermon, Stefano. 2016. Combining satellite imagery and machine learning to predict poverty. *Science*, **353**(6301), 790–794.
- Kolenikov, Stanislav, & Angeles, Gustavo. 2009. Socioeconomic Status Measurement With Discrete Proxy Variables: Is Principal Component Analysis A Reliable Answer? *Review of Income and Wealth*, **55**(1), 128–165.
- Lee, Sik-yum, Poon, Wai-yin, & Bentler, P. M. 1990. A Three-Stage Estimation Procedure for Structural Equation Models with Polytomous Variables. *Psychometrika*, **55**(1), 45–51.
- Linderman, By Scott W, Johnson, Matthew J, & Adams, Ryan P. 2015. *Dependent Multinomial Models Made Easy: Stick Breaking with the Poly-Gamma Augmentation*.
- Ludvigson, Sydney C. 2004. Consumer Confidence and Consumer Spending. *Journal of Economic Perspectives*, **18**(2), 29–50.
- MacDonald, Iain L, & Zucchini, Walter. 1997. *Hidden Markov and other models for discrete-valued time series*. Vol. 110. CRC Press.
- Mcauliffe, Jon D, & Blei, David M. 2008. Supervised topic models. *Pages 121–128 of: Advances in neural information processing systems*.
- Ng, Serena. 2015. Constructing Common Factors from Continuous and Categorical Data. *Econometric Reviews*, **34**(6-10), 1141–1171.
- Nicholson, Sean, & Simon, Kosali. 2010. How did the recession affect health and related activities of americans? *Preliminary and Incomplete Draft*.
- Pesaran, M Hashem. 2006. Estimation and inference in large heterogeneous panels with a multi-factor error structure. *Econometrica*, **74**(4), 967–1012.
- Ruiz, Francisco J. R., Athey, Susan, & Blei, David M. 2018 (nov). *SHOPPER: A Probabilistic Model of Consumer Choice with Substitutes and Complements*.

- Scruggs, Lyle, & Benegal, Salil. 2012. Declining public concern about climate change: Can we blame the great recession? *Global Environmental Change*, **22**(2), 505–515.
- Wang, Yixin, & Blei, David M. 2018. The blessings of multiple causes. *arXiv preprint arXiv:1805.06826*.
- Welling, Max, & Teh, Yee Whye. 2011. Bayesian Learning via Stochastic Gradient Langevin Dynamics. *Pages 681–688 of: Proceedings of the 28th International Conference on International Conference on Machine Learning*.
- Young, Alwyn. 2012. The African Growth Miracle. *Journal of Political Economy*, **120**(4), 696–739.



## A Gibbs Sampling Steps

In this section,  $p_{mult}(\cdot)$  is used to signify the multinomial density function,  $p_{dir}(\cdot)$  the Dirichlet density function, and  $p_{norm}(\cdot)$  the Normal density function.

### A.1 Steps for Markov-Switching Model

#### Generating $z_t$ conditional on $g, X_t, z_{t+1}, z_{t-1}$ , and $\beta$

We follow the single-move gibbs sampling procedure of Albert & Chib (1993). We assume that  $P(z_1|z_0) = \frac{1}{K}$  and  $P(z_{T+1}|z_T) = \frac{1}{K}$ .

For  $t = 1, \dots, T$ , sample  $z_t$  from  $\{1, \dots, K\}$  from the posterior distribution, which is multinomial with

$$p(z_t = j | X_t, z_{t-1}, z_{t+1}, \beta, X_t) \propto p(z_t | z_{t-1}) p_{mult}(X_t; \beta_j) p(z_{t+1} | z_t)$$

The first and last terms are taken directly from the current estimate for the transition matrix:  $p(z_t = j | z_{t-1} = k) = g_{jk}$ .

#### Generating transition matrix $g$ conditional on $z$

$$p(g_j | z) \propto p_{dir}(g_j; \alpha) p_{mult}(n_j; g_j)$$

where  $n_{i,j} = \sum_{t=2}^T \mathbb{1}(z_t = i) \mathbb{1}(z_{t-1} = j)$

The posterior distribution for each column  $g_j$  in the transition matrix  $g$  is independent Dirichlet:

$$g_j \sim \text{Dir}(\alpha_1 + n_{1j}, \dots, \alpha_K + n_{Kj})$$

#### Generating $\beta$ conditional on $X, z$

$$p(\beta_k | G, X; \eta) \propto p_{dir}(\beta_k; \eta) p_{mult}(m_k; \beta_k)$$

$$m_{k,v} = \sum_{t=1}^T \sum_{i=1}^{N_t} \mathbb{1}(\hat{X}_{ti} = p) \mathbb{1}(z_t = k)$$

The posterior distribution of the multinomial probabilities is independent Dirichlet for each profile  $k$ :

$$\beta_k \sim \text{Dir}(\eta_1 + m_{k1}, \dots, \eta_P + m_{kP})$$

### A.2 Steps for State Space Model

#### Generating $g_t$ conditional on $\sigma, g_{t-1}, g_{t+1}$ , and $z_t$

We adapt the method from Bhadury *et al.* (2016) for Dynamic LDA, and use Stochastic Gradient Langevin Dynamics (Welling & Teh, 2011) to draw  $g_t$ . SGLD is a form of gradient descent, adding Gaussian noise at each step, which Welling & Teh (2011) shows allows the method to generate samples from the true posterior without a Metropolis-Hastings test, as long as the shrinkage parameter  $\epsilon_n$  fulfils certain conditions.

$$p(g_t|g_{t-1}, g_{t+1}, z_t) \propto p_{norm}(g_t; g_{t-1}, \sigma^2 I) p_{norm}(g_{t+1}; g_t, \sigma^2 I) \prod_{i=1}^{N_t} p_{mult}(z_{ti}; \theta(g_t))$$

In step  $s$  of the gibbs sampler, for each  $k = 1, \dots, K$ ,

$$\Delta g_{t,k}^{(s)} = \frac{\epsilon_s}{2} \nabla_{g_{t,k}} \log p(g_t^{(s-1)} | g_{t-1}^{(s-1)}, g_{t+1}^{(s-1)}, z_t^{(s-1)}) + \psi_i, \quad \psi_i \sim N(0, \epsilon_s)$$

$$\nabla_{g_{tk}} p(g_t^{(s)} | g_{t-1}^{(s-1)}, g_{t+1}^{(s-1)}, z_t^{(s-1)}) = \frac{-1}{\sigma^2} (g_{k,t} - g_{k,t-1}) - \frac{1}{\sigma^2} (g_{k,t+1} - g_{kt}) + n_{tk} - N_t \psi(\theta_t)_k$$

$$n_{tk} = \sum_{i=1}^{N_t} \mathbb{1}(z_{ti} = k)$$

$\epsilon_s = a(b + s)^{-c}$  for gibbs sampling step  $s$ . We choose  $a = 0.1$ ,  $b = 1$  and  $c = 0.5$  for our applications. One downside of this method is for each application having to tune those parameters to get a sequence of  $\epsilon_s$  that allows for proper convergence.

### Generating $z_{ti}$ conditional on $\beta, X$ , and $g_t$

The posterior distribution of  $z_{ti}$  is multinomial with probabilities:

$$p(z_{ti} = k | \beta, g_t) \propto \beta_{k,p} \theta(g_t)_k$$

for  $X_{ti} = p$ .

### Generating $\beta$ conditional on $z$

$$p(\beta | z) \propto p_{dir}(\beta; \eta) p_{mult}(z; \beta)$$

The posterior distribution of the multinomial probabilities is Dirichlet for each profile  $k$ :

$$\beta_k | z \sim \text{Dir}(\eta_1 + m_{k,1}, \dots, \eta_P + m_{k,P})$$

$$m_{kp} = \sum_{t=1}^T \sum_{i=1}^{N_t} \mathbb{1}(X_{ti} = p) \mathbb{1}(z_{ti} = k)$$

### Generating $\sigma$ conditional on $g$

$$\sigma \sim \text{IGamma}(v1, s1)$$

$$v1 = v0 + T$$

$$s1 = s0 + \sum_{t=1}^T (g_t - g_{t-1})^2$$

## B Survey Index Component Variables

### B.1 DHS Data

The DHS survey data years and countries used are in Table ??.

The ten variables used to construct the wealth indices in the paper are mapped from the raw DHS data according to Table ??.

### B.2 Michigan Data

Below are the five questions used to create the Michigan indices.

1. Would you say that you are better off or worse off financially than you were a year ago?
  - (1) Better, (3) Same, (5) Worse, (8) Don't know or missing
2. Now looking ahead—do you think that a year from now you will be better off financially, or worse off, or just about the same as now?
  - (1) Better, (3) Same, (5) Worse, (8) Don't know or missing
3. Now turning to business conditions in the country as a whole—do you think that during the next twelve months we'll have good times financially, or bad times, or what?
  - (1) Good times, (2) Good with qualifications, (3) Pro-con, (4) Bad with qualifications, (5) Bad times, (8) Don't know or missing
4. Looking ahead, which would you say is more likely—that in the country as a whole we'll have continuous good times during the next five years or so, or that we will have period of widespread unemployment or depression, or what?
  - (1) Good times, (2) Good with qualifications, (3) Pro-con, (4) Bad with qualifications, (5) Bad times, (8) Don't know or missing
5. Generally speaking, do you think now is a good or bad time for people to buy major household items?
  - (1) Good, (3) Pro-con, (5) Bad, (8) Don't know or missing

### B.3 Gallup Poll Social Survey

#### Environmental Concerns Index

Below are the six questions used to create the environmental concerns index.

1. How would you rate the overall quality of the environment in this country today ?
  - (1) Excellent, (2) Good, (3) Only fair, (4) Poor, (5) Don't know or missing
2. Right now, do you think the quality of the environment in the country as a whole is :
  - (1) Getting better, (2) Getting worse, (3) Same, (4) Don't know or missing
3. With which one of these statements about the environment and the economy do you most agree?

- (1) Protect environment, even at risk of curbing economic growth, (2) Economic growth priority even if environment suffers to some extent, (3) Equal priority, (4) Don't know or missing
4. With which one of these statements about the environment and energy production do you most agree ?
- (1) Protect environment, even at risk of limiting energy supplies which the U.S. produces, (2) Development of U.S. energy supplies – such as oil, gas and coal – should be given priority, even if the environment suffers to some extent, (3) Equal priority, (4) Other, don't know or missing
5. Which of the following statements reflects your view of when the effects of global warming will begin to happen?
- (1) Already begun to happen, (2) Will start happening within a few years, (3) Will start happening within your lifetime, (4) Will not happen within your lifetime, but they will affect future generations, (5) Will never happen, (6) Don't know or missing
6. Thinking about what is said in the news, in your view is the seriousness of global warming:
- (1) Generally exaggerated, (2) Generally correct, (3) Generally underestimated (4) Missing

### **Traditional Values Index**

Below are the four questions used to create the traditional values index.

1. Regardless of whether or not you think it should be legal, please tell me whether you personally believe that in general abortion is:
  - (1) Morally acceptable, (2), Morally wrong, (3) Depends on the situation, (4) Not a moral issue, (5) Don't know or missing
2. Regardless of whether or not you think it should be legal, please tell me whether you personally believe that in general the death penalty is:
  - (1) Morally acceptable, (2), Morally wrong, (3) Depends on the situation, (4) Not a moral issue, (5) Don't know or missing
3. Regardless of whether or not you think it should be legal, please tell me whether you personally believe that in general suicide is:
  - (1) Morally acceptable, (2), Morally wrong, (3) Depends on the situation, (4) Not a moral issue, (5) Don't know or missing
4. How important would you say religion is in your own life?
  - (1) Very important, (2) Fairly important, (3) Not very important, (4) Don't know or missing

## Health Concerns Index

Below are the five questions used to create the health concerns index.

1. How would you describe your own physical health at this time?
  - (1) Excellent, (2) Good, (3) Only fair, (4) Poor, (5) Don't know or missing
2. How would you describe your own mental health or emotional well-being at this time?
  - (1) Excellent, (2) Good, (3) Only fair, (4) Poor, (5) Don't know or missing
3. Are you generally satisfied or dissatisfied with the total cost of health care in this country?
  - (1) Satisfied, (2) Dissatisfied, (3) Don't know or missing
4. Overall, how would you rate the quality of health care in this country ?
  - (1) Excellent, (2) Good, (3) Only fair, (4) Poor, (5) Don't know or missing
5. Overall, how would you rate health care coverage in this country ?
  - (1) Excellent, (2) Good, (3) Only fair, (4) Poor, (5) Don't know or missing