

Dynamic Hierarchical Latent Variable Models for Categorical Survey Data

Evan Munro and Serena Ng

February 14, 2019

Abstract

Existing methods for extracting summary indexes from multivariate discrete data are generally ad-hoc, often treating discrete variables as if they were continuous. We provide an overview of existing methods for analyzing cross-sectional categorical survey data and show the benefit of using hierarchical latent variable models instead. Then, we introduce two dynamic models for categorical time series 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 index constructed using the Michigan data is highly correlated to the established index constructed using averaging methods, but has additional economic interpretability. Then, we estimate the model on GPSS data, which does not have an established index, and extract economic insights that are not as clear when studying the data on a question-by-question basis. For example, we document how environmental concern fluctuates with financial insecurity, dropping steeply during the 2008-2009 recession and not recovering until after 2015.

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 polychloric 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 aggregate-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 a population based on a population-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 show how depending on the type of data available, dynamics can be added at an aggregate or an individual level; we choose to focus on methods for time-series survey data which have dynamics at an aggregate level, but point to related methods for panel data with dynamics at an individual level. 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 are not well-represented by a continuous variable. The index of Michigan Consumer sentiment, 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 which questions 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 and $i_c \in \{1, \dots, N_c\}$ indexes the individuals 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 .

2.1 Existing Methods

The first setting that we will examine is where the goal is to recover a K -dimensional real-valued index g_i that is a low dimensional representation of X_i , the vector of survey responses for individual i . Often $K = 1$ 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. So, many researchers use survey data on household consumption such as 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 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. This has been shown to result in poor results compared to alternative methods; for example, 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).

Alternative methods include treating discrete data as continuous and doing PCA on the raw data X directly (Kolenikov & Angeles, 2009). 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. A final alternative method 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 continuous 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 deal with 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 or 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 and have substantial room for improvement.

To address some of these issues, 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), firm management [CITE] and text analysis of central bank communication [CITE].

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 lies in at what location in the hierarchy of individual, group, and population the mixture is located. Each corresponds to a different set of conditional independence assumptions involving the data and the specified latent variables.

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 Filmer-Pritchett or polychoric correlation-

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

Table 1: Hierarchical Latent Variable Models

based PCA. 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:

$$\begin{aligned}
X_{ij} | z_{ij}, \beta_j &\sim \text{Multinomial}(\beta_{j, z_{ij}}) \\
z_{ij} | g_i &\sim \text{Multinomial}(g_i) \\
g_i &\sim \text{Dirichlet}(\alpha) \\
\beta_{jk} &\sim \text{Dirichlet}(\eta_j)
\end{aligned}$$

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 β_{jk} is independent of β_{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 add 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 β and G depend on the hyperparameters α and η 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 assuming that they are 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 very wealthy household has g_{i2} close to 1, and a poor household has g_{i2} close to 0. A middle class 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 “rich” extreme 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 α).

2.3 Latent Dirichlet Allocation Model

GoM, PCA, and polychloric 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 β_k . In some cases, though, applied economists may not require an individual-level index, but are instead interested in an aggregate-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 an aggregate level index simply by averaging an individual-level index. However, if an aggregate level index is the target, the reduction in parameters from estimating an aggregate 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 β_k 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 using a function $c : \{1, \dots, N\} \rightarrow \{1, \dots, C\}$. 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:

$$\begin{aligned} X_i | z_i, \beta &\sim \text{Multinomial}(\beta_{z_i}) \\ z_i | g_{c(i)} &\sim \text{Multinomial}(g_{c(i)}) \\ g_c &\sim \text{Dirichlet}(\alpha) \\ \beta_k &\sim \text{Dirichlet}(\eta) \end{aligned}$$

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.*

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

$$Pr(X_i | g_{c(i)}, z_i, \beta) = Pr(X_i | z_i, \beta)$$

$$Pr(z_i | g_{c(i)}, \beta) = Pr(z_i | g_{c(i)})$$

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

The joint likelihood for LDA can be factorized using the conditional independence assumptions in a similar to how the GoM likelihood was derived.

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

The estimate of g_c provides an aggregate-level index for group c . Previewing the World Management Survey example, an economist could assume that there are two predominant profiles of management practices, one with higher probability on permutations of responses with poor management practices and the other with higher probability on permutations of responses with lots of good management practices. Due to country-specific institutions, there is a country-specific mixture of firms with good vs. bad management characteristics. The weight on the good management profile in g_c can be interpreted as a country-specific index of good management.

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 aggregate-level mixture model. It is also worth pointing out the relationship of these latent variable models to a Bayesian latent class model, which is a population-level mixture model. A Bayesian multinomial mixture model assumes outcomes in a dataset come from a finite set of profiles in a population. These profiles correspond to a multinomial distribution over outcomes.

Below is the specification of the basic multinomial mixture model:

$$X_i | \beta, z_i \sim \text{Multinomial}(\beta_{z_i})$$

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

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

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

In a latent class model, the mixture g is population-level, so is K -dimensional, rather than $N \times K$ as in GoM. Each individual is assumed to belong to a single profile based on the population-level mixture. Each question response is assumed to come from the individual's assigned profile. An individual's responses are generated from one of K profiles, whereas in GoM, an individual is modeled as a mixture of K profiles. In a latent class model, heterogeneity is finite, whereas in GoM there is infinite heterogeneity generated from a finite set of low-dimensional latent variables. The conditional independence assumption for individuals in a latent class model are strict compared to the corresponding assumptions in GoM and LDA.

Latent Class 1 : Conditional Independence of Individuals: *Conditional on the population mixture over profiles g , the probability of a certain response and profile assignment for individual i is independent of individual r 's responses. Furthermore, responses are independent of the population-level mixture given the profile assignments and assignments are independent of profiles given the population-level mixture.*

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

$$Pr(X_i|g, z_i, \beta) = Pr(X_i|z_i, \beta)$$

$$Pr(z_i|g, \beta) = Pr(z_i|g)$$

Latent Class 2 : Independence of Profiles β_k is independent of β_f for $f \neq k$.

The resulting joint likelihood of the model is

$$p(\beta, g, z, X) = Pr(g) \prod_{k=1}^K Pr(\beta_k) \prod_{i=1}^N g_{z_i} \beta_{z_i, X_i}$$

The model assumes that all relationships between individual's responses are captured only by the population-level mixture over profiles, which is not likely to hold in most survey data where there are natural groups with correlated responses within the population. 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 population-level, rather than an individual or aggregate level continuous-valued index. Erosheva *et al.* (2007) proved that individual-level mixture models can be considered a restricted form of a latent class model with an exponential number of profiles.

In summary, in this section we have showed that:

1. GoM is a structural alternative to Filmer-Pritchett and Polychloric correlation-based PCA methods that does not involved unwarranted transformations of the data, and allows ordered, unordered, and missing responses in categorical data
2. GoM is part of a class of hierarchical latent variable models with mixture weights at the individual (GoM), aggregate (LDA), or population-level (Latent Class).
3. Every Bayesian hierarchical latent variable model relies on a set of conditional independence assumptions. Some of the assumptions that they rely on 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 β_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. Then, the hyperparameter η is adjusted so that β_{kj} has high prior probability of being close to zero for the extreme response assigned to the other profile. In the applications, we do this by setting η_{jke} to 0.01, where e is the index of the extreme response for the other profile. With unordered outcome variables, some judgement is involved in determining what the extreme response is.

In addition to the label-swapping issue, there is also the potential issue of the posterior distribution 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 [CITE]. 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, $t : 1, \dots, N \rightarrow \{1, \dots, T\}$. 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

	MS-M	SS-M
Gaussian Version	Markov-Switching	Lin. State Space
Poisson Version	CITE	CITE
Profile Parameter (β)	$K \times P$	$K \times P$
Assignment Parameter (z)	$T \times 1$	$N \times 1$
Mixture Parameter (g)	$K \times K$	$T \times K$
Mixture Hyperparameter (α)	$K \times K$	$T \times K$
Profile Hyperparameter (η)	$K \times P$	$K \times P$

Table 2: Dynamic Hierarchical Latent Variable Models

assignment parameters z in a dynamic version of the latent class model. This is the multinomial version of a gaussian 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 basic linear state space model.

There is a related statistics literature which has derived markov-switching models and state 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 and Poisson 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, rather than coming from a population level mixture, follows a markov-switching process governed by a population-level transition matrix g .

The model specification is as follows:

$$X_i | \beta, z_t \sim \text{Multinomial}(\beta_{z_{t(i)}})$$

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

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

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

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

Assumption Latent Class 1 is replaced by MS-M 1 and MS-M 2. The final conditional independence assumptions is 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 population-level 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 population-level transition matrix given the profile assignments.

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

MS-M 3 : Independence of Profiles β_k is independent of β_f for $f \neq k$.

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

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

$Pr(\beta_k)$ 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 B.

1. Generate z_t conditional on g , X , z_{t+1} , z_{t-1} , and β 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 β 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

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

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

. The model specification is:

$$\begin{aligned} X_i | z_t &\sim \text{Multinomial}(\beta_{z_{t(i)}}) \\ z_{ti} &\sim \text{Multinomial}(\psi(g_t)) \\ g_t &= g_{t-1} + w_t, \quad w_t \sim N(0, \sigma I) \\ \sigma &\sim \text{InverseGamma}(v_0, s_0) \\ \beta_k &\sim \text{Dirichlet}(\eta) \end{aligned}$$

In the previous section, when we added dynamics to the discrete-valued assignments z , the dynamics were captured using a discrete markov process. In this model, where the dynamics are added to the group-specific mixtures, we use a linear markov process. The model is denoted the multinomial state space model, since it involves adding a link function to a multinomial outcome to the dynamics of a Bayesian linear state space model.

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 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.

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

$$Pr(X_i|g_{t(i)}, z_i, \beta) = Pr(X_i|z_i, \beta)$$

$$Pr(z_i|g_{t(i)}, \beta) = Pr(z_i|g_{t(i)})$$

SS-M 2 Markov Property of g_t and Independence of Profiles β_k is independent of β_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:

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

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

1. Generate g_t conditional on σ, g_{t-1}, g_{t+1} , and z_t using Stochastic Gradient Langevin Dynamics
2. Generate z conditional on β, X , and g using a multinomial posterior
3. Generate β conditional on z using a Dirichlet posterior
4. Generate σ 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 β . In the context of LDA, 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 Data

4.1 Michigan Survey Data

The raw 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.

The raw data on the survey respondents 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 C.

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

The data is mapped to the model framework as follows. Each month t has N_t survey respondents. $x_{t,n}$ is the response for the n -th survey respondent for month t and corresponds to an index $v \in \{1, \dots, V\}$. Each index v corresponds to a permutation of potential survey responses. So, for the Michigan data, $x_{t,n} = 1$ corresponds to a survey respondent answering 11111 to the survey, which means the respondent is an optimist who answers “better” or “good” to all of the questions on the U.S. economy.

4.2 Gallup Poll Social Series

The raw 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 repeated each year since 2001. The most recent data available is from 2016. 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 C. The Gallup poll data is mapped into the model framework in the same way as for the Michigan data; each outcome corresponds to a permutation of potential survey responses.

4.3 Existing Work

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 Application

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.

5.1 Michigan Consumer Sentiment Indices

The Michigan ICS is constructed from the five questions listed in Appendix C based on the following ad-hoc procedure, published by the University of Michigan on their website:

1. Compute the relative scores X_i 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.

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

I estimate SS-M on the Michigan Consumer survey data and compare the estimated weight g_{1t} on the positive sentiment to the ICS, rescaled to match the range and variance of the new index, which is between 0 and 1. To deal with the label-swapping identification issue that is inherent in both new models, we identify the “positive” sentiment using the prior distribution on β . In a two-state model, we do this by decreasing the prior probability of the extreme outcome on state 1. For the Michigan data, this corresponds, for example, to 55555, which corresponds to an individual

who selects the most negative responses for each of the five questions making up the ICS. This can be considered as the discrete and Bayesian analogy to the lower triangular restriction that is placed on factor models for identification purposes. The states are then labelled by ordering them based on the posterior probability of the extreme outcome; the one involving a higher probability of responding 55555 is the more “negative” sentiment and the other is the more “positive” sentiment.

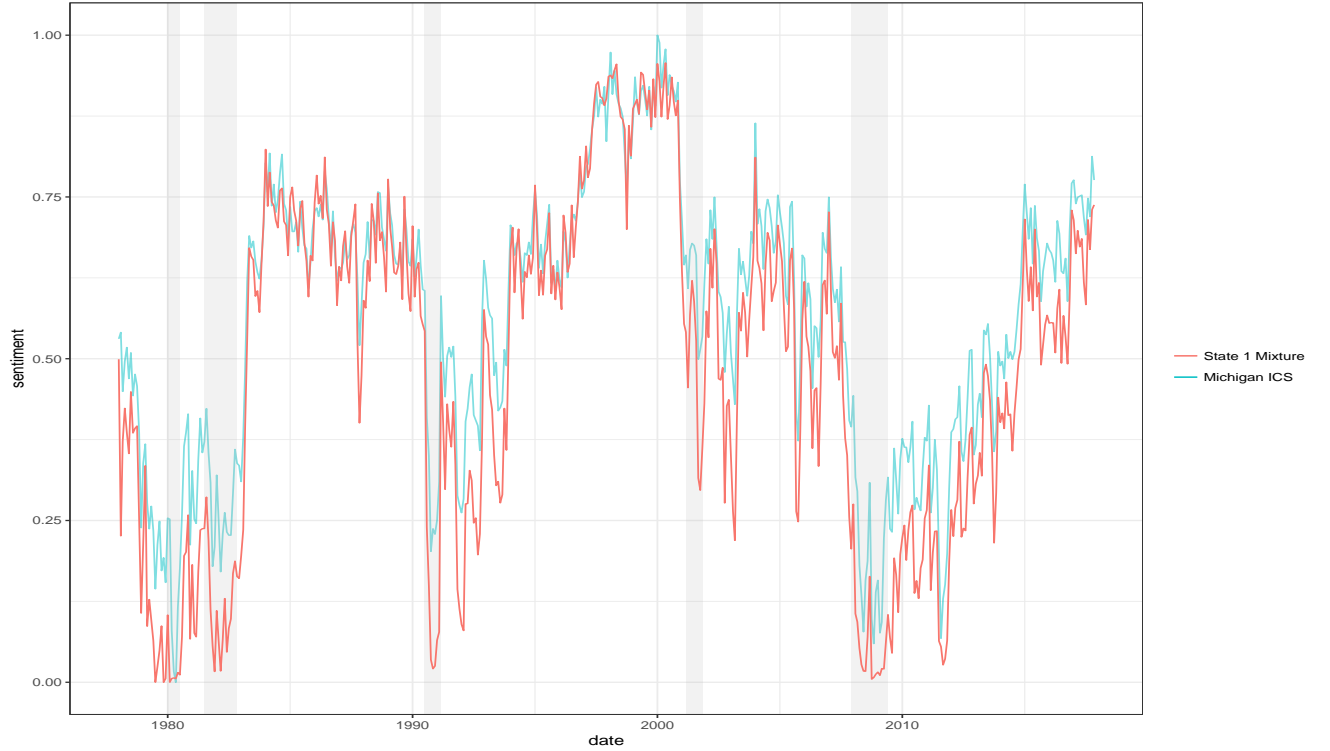


Figure 1: Probability of State 1 vs. ICS

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 3. 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 3 contains the responses with the top 5 weights for each multinomial distribution parameter β_k corresponding to the two estimated states.

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 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.

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 3: Model 1: Top 5 State-Specific Multinomial Probabilities

The latent variable G_t is now discrete and has dynamics corresponding to a transition matrix, compared to g_t which is continuous and has random walk dynamics. The below table 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 4: 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 4 and Table 3 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. The State 2, in the markov-switching model, then, 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 $P(S_{t+1} = 1|S_t = 1) = 0.89$ and $P(S_{t+1} = 2|S_t = 2) = 0.80$. The below figure shows the probability of being in State 2, the more pessimistic state. The Michigan ICS is also plotted and recessions shaded.

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 below.



Figure 2: Probability of State 2

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.

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

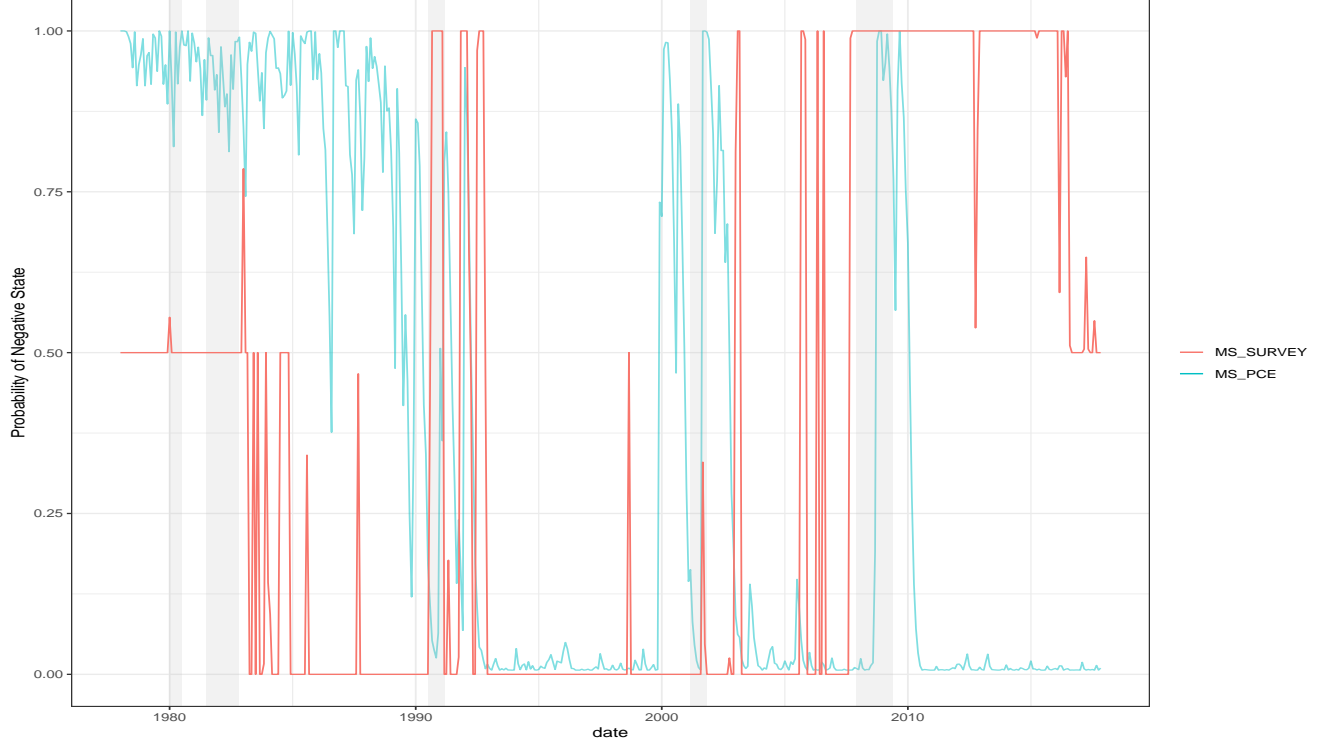


Figure 3: Probabilities of Negative State in DiscreteMS vs. MS on PCE

additional interpretation, since each corresponds to a state 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 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 or values. 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 states, using the six questions in Appendix C on the respondents' views on current and future environmental conditions in the U.S. The weight on the first state 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.

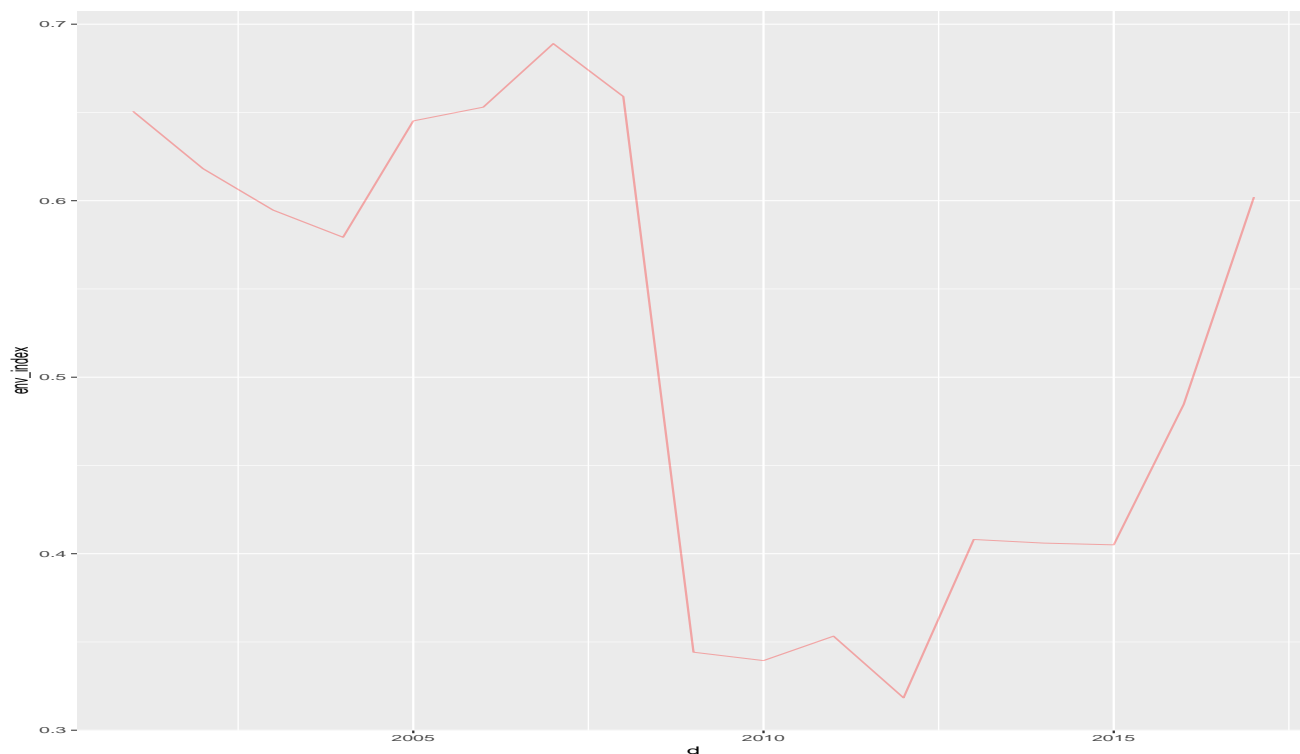


Figure 4: Index of Environmental Concerns

We derive indices from sets of questions from a few other GPSS categories to compare. In Figure 5, we plot the weight on the state 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 state 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, does bottom out in 2009, as shown in Figure 6. 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 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 C to successfully

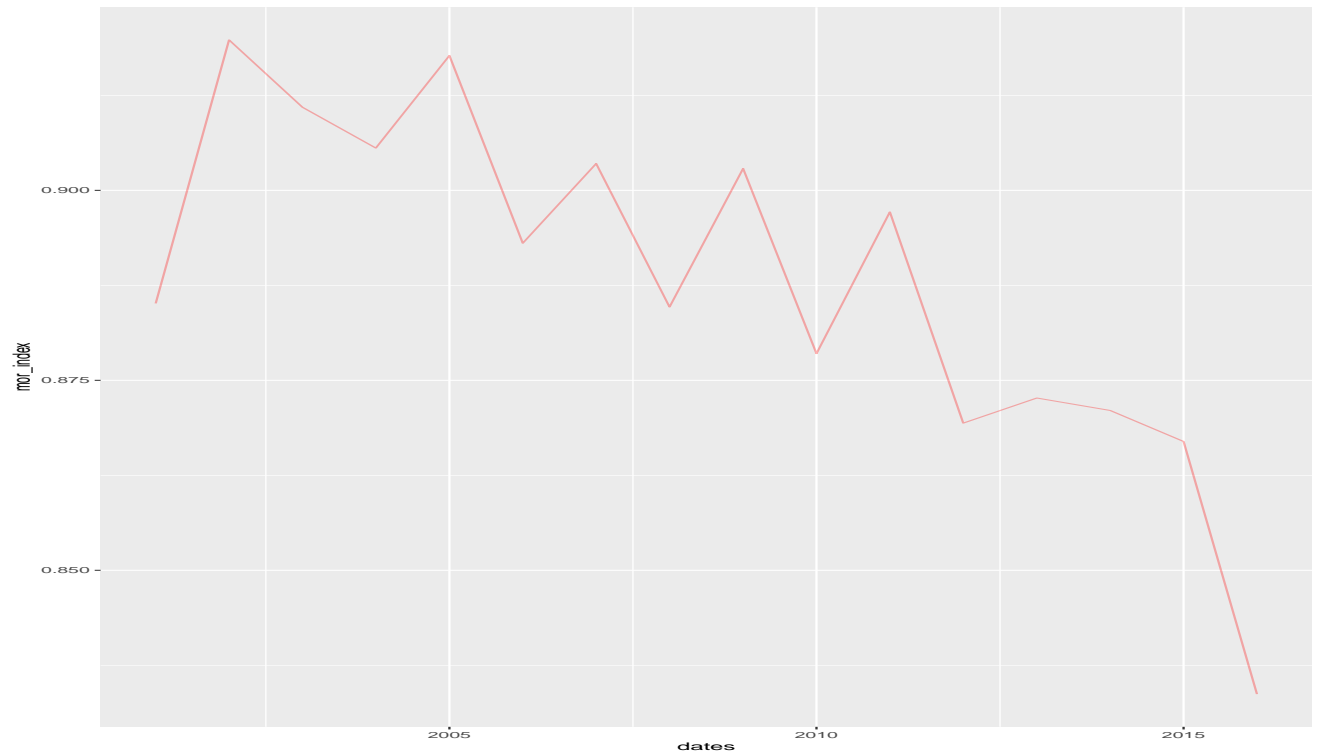


Figure 5: Index of Traditional Values

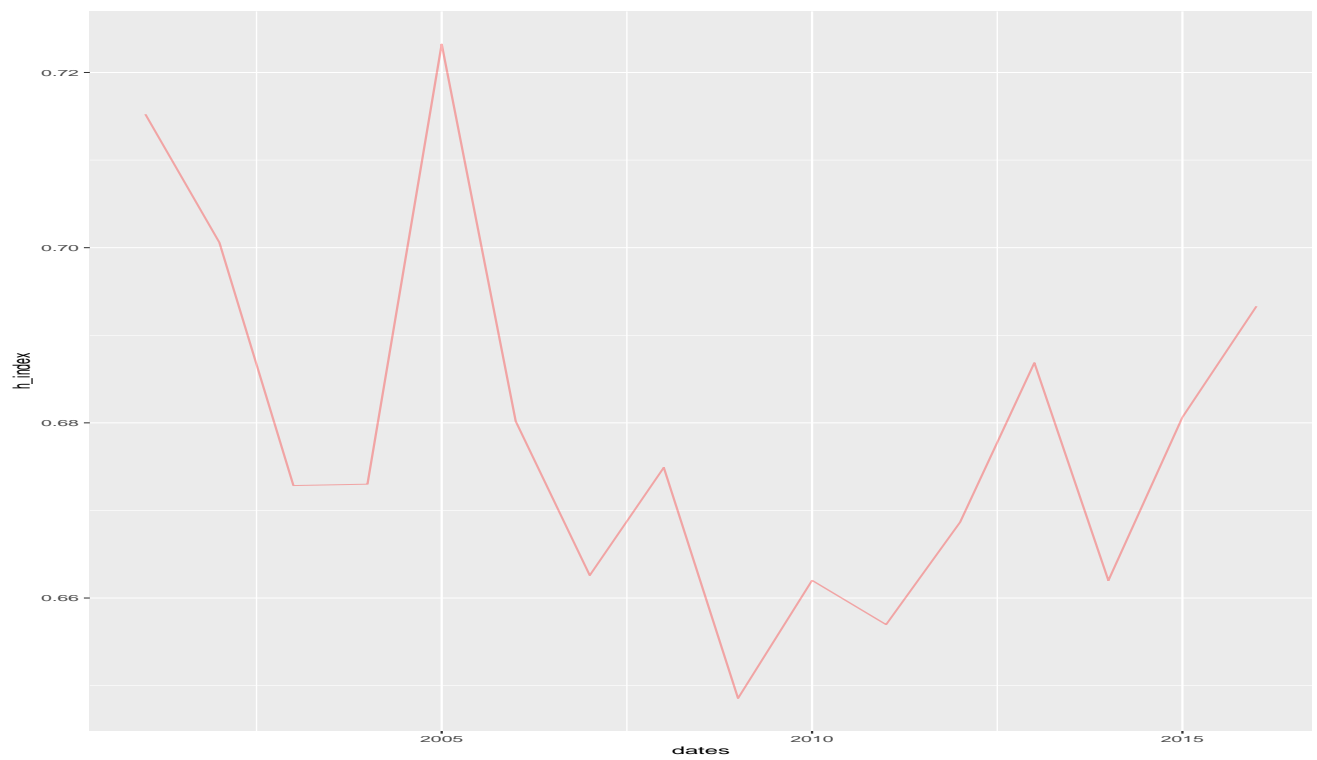


Figure 6: Index of Health Concerns

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 parameteric 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 a variety of ordered and unordered categorical time series data on consumer sentiment towards finance, health, environment, and traditional values.

The relation of these models to the extensive computer science literature on hierarchical latent variable models suggests a variety of extensions for the models, including introducing more complex forms of dynamics and adding the relationship between responses and some outcome variable [CITE sLDA?].

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.
- 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.
- 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.
- 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 Polya-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.
- Nicholson, Sean, & Simon, Kosali. 2010. How did the recession affect health and related activities of americans? *Preliminary and Incomplete Draft*.
- 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.
- 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*.

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 G_t conditional on $R, \hat{X}_t, G_{t+1}, G_{t-1}$, and β

We follow the single-move gibbs sampling procedure of Albert & Chib (1993). We assume that $P(G_1|G_0) = \frac{1}{K}$ and $P(G_{T+1}|G_T) = \frac{1}{K}$.

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

$$p(G_t = i | \hat{X}_t, G_{t-1}, G_{t+1}, \beta, \hat{X}_t) \propto p(G_t | G_{t-1}) p_{mult}(\hat{X}_t; \beta_{G_t}) p(G_{t+1} | G_t)$$

The first and last terms are taken directly from the current estimate for the transition matrix: $p(G_t = i | G_{t-1} = j) = R_{i,j}$.

Generating transition matrix R conditional on G

$$p(R_j | G) \propto p_{dir}(R_j; \alpha) p_{mult}(n_j; G_j)$$

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

The posterior distribution for each column R_j in the transition matrix R is independent Dirichlet:

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

Generating β conditional on \hat{X}, G

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

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

The posterior distribution of the multinomial probabilities is independent Dirichlet for each state 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 σ, 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 ϵ_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 ϵ_s that allows for proper convergence.

Generating z_{ti} conditional on β, \hat{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 $\hat{X}_{ti} = p$.

Generating β 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 state k :

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

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

Generating σ 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 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.2 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