

Battaglia, Laura; Christensen, Timothy; Hansen, Stephen; Sacher, Szymon

Working Paper

Inference for Regression with Variables Generated from Unstructured Data

CESifo Working Paper, No. 11119

Provided in Cooperation with:

Ifo Institute – Leibniz Institute for Economic Research at the University of Munich

Suggested Citation: Battaglia, Laura; Christensen, Timothy; Hansen, Stephen; Sacher, Szymon (2024) : Inference for Regression with Variables Generated from Unstructured Data, CESifo Working Paper, No. 11119, CESifo GmbH, Munich

This Version is available at:

<https://hdl.handle.net/10419/300047>

Standard-Nutzungsbedingungen:

Die Dokumente auf EconStor dürfen zu eigenen wissenschaftlichen Zwecken und zum Privatgebrauch gespeichert und kopiert werden.

Sie dürfen die Dokumente nicht für öffentliche oder kommerzielle Zwecke vervielfältigen, öffentlich ausstellen, öffentlich zugänglich machen, vertreiben oder anderweitig nutzen.

Sofern die Verfasser die Dokumente unter Open-Content-Lizenzen (insbesondere CC-Lizenzen) zur Verfügung gestellt haben sollten, gelten abweichend von diesen Nutzungsbedingungen die in der dort genannten Lizenz gewährten Nutzungsrechte.

Terms of use:

Documents in EconStor may be saved and copied for your personal and scholarly purposes.

You are not to copy documents for public or commercial purposes, to exhibit the documents publicly, to make them publicly available on the internet, or to distribute or otherwise use the documents in public.

If the documents have been made available under an Open Content Licence (especially Creative Commons Licences), you may exercise further usage rights as specified in the indicated licence.

Inference for Regression with Variables Generated from Unstructured Data

Laura Battaglia, Timothy Christensen, Stephen Hansen, Szymon Sacher

Impressum:

CESifo Working Papers

ISSN 2364-1428 (electronic version)

Publisher and distributor: Munich Society for the Promotion of Economic Research - CESifo GmbH

The international platform of Ludwigs-Maximilians University's Center for Economic Studies and the ifo Institute

Poschingerstr. 5, 81679 Munich, Germany

Telephone +49 (0)89 2180-2740, Telefax +49 (0)89 2180-17845, email office@cesifo.de

Editor: Clemens Fuest

<https://www.cesifo.org/en/wp>

An electronic version of the paper may be downloaded

- from the SSRN website: www.SSRN.com
- from the RePEc website: www.RePEc.org
- from the CESifo website: <https://www.cesifo.org/en/wp>

Inference for Regression with Variables Generated from Unstructured Data

Abstract

The leading strategy for analyzing unstructured data uses two steps. First, latent variables of economic interest are estimated with an upstream information retrieval model. Second, the estimates are treated as “data” in a downstream econometric model. We establish theoretical arguments for why this two-step strategy leads to biased inference in empirically plausible settings. More constructively, we propose a one-step strategy for valid inference that uses the upstream and downstream models jointly. The one-step strategy (i) substantially reduces bias in simulations; (ii) has quantitatively important effects in a leading application using CEO time-use data; and (iii) can be readily adapted by applied researchers.

JEL-Codes: C110, C510, C550.

Keywords: unstructured data, information retrieval, topic modeling, Hamiltonian Monte Carlo, measurement error.

Laura Battaglia
Oxford University / UK
laura.battaglia@keble.ox.ac.uk

Timothy Christensen
UCL University College London / UK
t.christensen@ucl.ac.uk

Stephen Hansen
UCL University College London / UK
stephen.hansen@ucl.ac.uk

Szymon Sacher
Stanford University / CA / USA
sacher@stanford.edu

May 7, 2024

Authors are listed in alphabetical order. This paper first circulated without TC as co-author under the title “Hamiltonian Monte Carlo for Regression with High-Dimensional Categorical Data” (<https://doi.org/10.48550/arXiv.2107.08112>) which was the second chapter of SS’s PhD thesis. SH acknowledges funding from ERC Consolidator Grant 864863, which supported his and LB’s time. We thank Nick Bloom, Germain Gauthier, Evan Munro, David Rossell, and Leif Thorsrud for feedback, as well as seminar and workshop participants at Aarhus University, Barcelona School of Economics, Bates College, Columbia University, ETH Zurich, UC San Diego, University of Southern California, University of Warwick, and the 3rd Monash-Warwick-Zurich Text-as-Data Workshop. The authors also thank the NumPyro development team for their outstanding work.

1 Introduction

The amount of unstructured data is growing rapidly and empirical work in economics is increasingly using it. The leading example of such data is text (Gentzkow et al. 2019a, Ash and Hansen 2023), but others include surveys, images, and audio recordings. In economics, unstructured data is primarily used to measure variables of interest that aren't observed in conventional quantitative data sources. Examples abound: Baker et al. (2016) measure economic policy uncertainty with newspaper text; Hoberg and Phillips (2016) infer firms' latent industries with corporate filings; Hansen et al. (2018) construct measures of policy deliberation from Federal Open Market Committee (FOMC) transcripts; Magnolfi et al. (2022) use survey data to measure product differentiation; Compiani et al. (2023) measure substitutability between products using Amazon text and image data; Gorodnichenko et al. (2023) measure tone-of-voice from audio recordings of FOMC press conferences; Gabaix et al. (2023) impute firm characteristics from investor holdings data; Einav et al. (2022) infer patients' health status from surveys; Vafa et al. (2023) construct a measure of labor market experience based on CVs. These derived measures are rarely an end in themselves. Rather, the goal is to study how the latent variables they measure interact with the economic environment. As such, they are typically plugged-in to downstream econometric models whose parameters are the main object of study. In practice, the *upstream* information retrieval (IR) model used to extract measurements from unstructured data and the *downstream* econometric model are almost always taken as wholly separate: the output of the upstream model is treated as observed "data" in the downstream model. We call this the *two-step strategy*.

While clearly a pragmatic initial approach, the two-step strategy has largely unknown statistical properties. On one hand, ignoring the upstream estimation step in downstream inference suggests a generated regressor problem (Pagan 1984). On the other, results in the time-series literature suggest plugging-in estimated latent variables need not lead to inference problems (Stock and Watson 2002, Bernanke et al. 2005, Bai and Ng 2006). More generally, characterizing the statistical guarantees—or lack thereof—of the two-step strategy is an important step in establishing a more mature understanding of reliable inference methods for unstructured data, an area that is still in its infancy.

Our first contribution is to provide theoretical arguments for why the dominant two-step strategy leads to biased inference on regression parameters in empirically plausible settings. We consider a set of n observations of quantitative and unstructured data. As many unstructured datasets can be represented as high-dimensional categorical data, we treat each unstructured observation as a high-dimensional vector of feature counts.¹ We

¹For example, one of the simplest representations of text is the *bag-of-words* model. Each document is represented as a vector of integer counts over the unique vocabulary terms in the corpus. Even relatively small corpora contain thousands of unique dimensions.

allow the amount of unstructured data to vary across observations, letting C_i denote the amount of unstructured data for observation i . The relative magnitudes of n versus moments of C_i play a key role in our analysis. We next specify a statistical model with three parts: a distribution over the feature-count vectors; a low-dimensional, latent variable representation for each such distribution; and a regression of an observed outcome variable onto the latent variables.² The two-step strategy (i) estimates the latent variables from the observed feature counts, then (ii) regresses the outcome variable on these estimates. Our primary theoretical question is: under what conditions do the estimated coefficients and standard errors from (ii) allow for valid inference?

The basic problem is measurement error: the regressors in step (ii) contain estimated rather than true latent variables. As is well known, measurement error leads to inconsistent estimates and distorted standard errors, both of which are present as the number of observations n grows with a fixed amount of unstructured data per observation. Because both the sample size and the amount of unstructured data per observation are typically large in applications, we allow n and the distribution of C_i to evolve together so that sampling error and measurement error are both relevant for inference.³ Our main finding is that, whenever $\sqrt{n} \times \mathbb{E} \left[\frac{1}{C_i} \right]$ tends to a constant $\kappa > 0$, the estimator is consistent but there is a bias present in the asymptotic distribution of the regression coefficients. Larger values of κ give relatively greater importance to measurement error, and hence a larger bias. However, the asymptotic variance is the same as that from regression onto the true latent variables and the usual OLS standard errors are consistent. Hence, treating the estimated latent variables as observed data in the regression does not distort the width of confidence intervals, but centers them away from the truth. This contrasts with the generated regressor literature, which emphasizes the variance distortion arising from plugging-in estimates as data.⁴ Only when $\kappa = 0$, so that sampling error dominates measurement error asymptotically, does the two-step strategy allow for valid inference.

Of course, κ describes limiting behavior and cannot be used to directly compute the magnitude of the bias in a given finite dataset. But our theory provides insights into when this bias is potentially problematic. Take, for example, job postings data as recorded by Lightcast (formerly Burning Glass), which has been used in dozens of papers. In 2022, there were 45 million individual job postings in the United States, with an

²This model is the basis for a large empirical literature. Examples include Nimczik (2017), Hansen et al. (2018), Mueller and Rauh (2018), Larsen and Thorsrud (2019), Bandiera et al. (2020), Thorsrud (2020), Bybee et al. (2020), Adams et al. (2021), Draca and Schwarz (2021), Olivella et al. (2021), Ash et al. (2022), and Munro and Ng (2022).

³Our use of sequences of DGPs to better approximate the finite-sample behavior of estimators is similar in spirit to the weak instrument literature (Staiger and Stock 1997), earlier work on measurement error (Chesher 1991), unit root testing (Phillips 1987), and large n, T panels (Hahn and Kuersteiner 2002).

⁴In the classic generated regressor problem (Pagan 1984) there is a common finite-dimensional parameter estimated in the first stage whereas here all n latent covariates are estimated.

average inverse posting length of 0.003. The empirical analogue of κ is $\sqrt{45,000,000} \times 0.003 \approx 20$, suggesting that measurement error may be distorting inference. Another example is the popular Nielsen Homescan database, where the empirical analogue of κ is approximately 4. A second insight is that the magnitude of the bias arises from the average *inverse* amount of unstructured data per observation. So, if a dataset has a long tail of observations with little data, a bias can arise even if there is a substantial amount of unstructured data per observation on average. Again taking the Lightcast data, the average document size is 575. If one used $1/575$ in place of $\mathbb{E}[C_i^{-1}]$ to compute the analogue of κ , 20 would fall to below 12, highlighting the role of tail behavior in driving measurement error bias. There are other cases, such as US patent data, where the empirical analogue of κ takes a smaller value around 1 because each observation has a relatively large amount of unstructured data. Because the exact magnitude of the problem is hard to assess in any given setting, it is important to develop *robust* inference methods that correctly account for measurement error.

Our second contribution is to propose such an inference method: directly use the model’s joint distribution over unstructured data, latent variables, and numeric outcomes to perform maximum likelihood estimation. We refer to this as the *one-step strategy*.

While implementing the one-step strategy is straightforward theoretically, it presents a major computational challenge due to the large number of observation-specific latent variables that must be integrated out. To address this, we use Hamiltonian Monte Carlo (HMC; MacKay 2003, Neal 2012), a Markov Chain Monte Carlo algorithm that uses information on the gradient of a distribution to sample from it. Implementation is greatly simplified with the use of modern probabilistic programming languages: one simply specifies the likelihood in code, which is then “automatically” compiled to perform sampling. This paradigm is useful for applied researchers because it allows one to focus on model development without the need to re-write the estimation and inference algorithms each time the model is changed. To this end, we use the NumPyro package (Bingham et al. 2018, Phan et al. 2019), which requires only a few lines of code to implement and is well suited to handling models with a large number of latent variables.⁵

Third, we compare the performance of the two-step and one-step strategies in an applied setting. To this end, we introduce the Supervised Topic Model with Covariates (STMC) which combines elements of existing models (Blei et al. 2003, Roberts et al. 2014, Ahrens et al. 2021) but is, to the best of our knowledge, a new statistical model of unstructured data. The model reduces the dimensionality of feature-count vectors by projecting them onto a set of latent factors (or topics), as in Probabilistic Latent

⁵Previous papers that have performed inference using the joint likelihood approach with unstructured data include Gentzkow et al. (2019b), Ruiz et al. (2020), and Munro and Ng (2022). These typically require extensive code to estimate, which makes adapting the model difficult for non-specialists.

Semantic Analysis (Hofmann 1999) and Latent Dirichlet Allocation (Blei et al. 2003). The dependence of outcome variables on latent factor loadings and observed covariates is captured by a downstream regression model. Additionally, the factor loadings can depend on a potentially different set of covariates via the upstream model. All components are woven together by a joint likelihood. Specifying the model in code takes fewer than 25 lines, illustrating how one can perform automatic inference in a new setting that would previously have required a bespoke and complex codebase.

Many important research questions can be addressed with STMC. Suppose each unstructured observation is a monetary policy speech. One latent topic might have an interpretation as price rises, so its loadings represent how much each speech discusses price rises. A first research question, which can be addressed with the downstream model, might ask how speakers’ attention to price rises is related to their policy actions. A second research question might ask how policymakers’ backgrounds relate to the attention they devote to price rises. That question can be addressed with the upstream model.

In simulated data, we show that the two-step strategy produces estimates that exhibit a bias which is increasing in κ . Moreover, two-step confidence interval widths are similar to those obtained using the true latent variables as covariates. Both of these findings reinforce the main predictions of our theory. By contrast, the one-step strategy produces estimates that appear unbiased and confidence intervals that have both the correct width and the correct centering.

Next, we revisit the empirical application from Bandiera et al. (2020) which uses the two-step strategy to first estimate latent CEO behaviors from a CEO time-use survey, then explains firm performance using the estimated behaviors. The one-step strategy substantively changes estimates compared to the two-step strategy. For instance, the estimated effect of having an MBA degree on behavior more than doubles when using the one-step instead of two-step strategy. To further test our theory, we next reduce the amount of unstructured data per observation and again deploy both inference strategies. This increases measurement error in latent behavior, and hence should increase the bias of the two-step strategy. Since the one-step strategy is always unbiased (asymptotically), one should observe larger differences in estimates, which is what we find. The estimated impact of behavior on firm performance, equivalent for both methods in the original data, is now half as large under the two-step strategy. Moreover, under the two-step strategy, the estimated effects of CEO characteristics on behavior reduce by a factor of three.

The settings we consider have an upstream unsupervised learning problem where unstructured data is used to recover latent variables of interest. A parallel set of recent papers considers the upstream generation of regressors via supervised learning (Fong and Tyler 2021, Allon et al. 2023, Egami et al. 2023, Zhang et al. 2023). These papers rely on the strong assumption that the latent variable is observed with no error for a subset of

observations. Their proposed solution uses the correctly labeled subset of data to build IV/GMM estimators. In our baseline setting, in common with much of the economics literature, all latent variables remain unobserved, but our framework could be extended easily to incorporate noisy labels.

Gentzkow et al. (2019b) highlights measurement error in the estimation of group differences in observed speech. This error emerges when the number of observed draws per group (analogous to C_i in our setting) is small relative to the number of possible vocabulary terms. To examine this situation, we extend our baseline model to consider regressions on text-derived similarity measures. Analogously to Gentzkow et al. (2019b) and our baseline setting, there is no bias when average inverse document size grows sufficiently fast with the sample size n . But when it doesn't, so that measurement error is also relevant asymptotically, a bias again emerges.

Our overall message is that the dominant approach for using unstructured data in empirical work in economics may suffer from measurement error which biases inference. We illustrate this formally with theoretical arguments in an empirically relevant albeit restrictive setting, but the take-away applies much more broadly. On a more positive note, a solution exists that is easy to implement and computationally feasible. We therefore see the one-step strategy as a robust and widely applicable starting point for empirical analysis. For instance, an emerging line of research uses text-derived sentiment indices as inputs into forecasting models with a vector autoregressive or dynamic factor structure. Straightforward extensions of our theoretical arguments can be used to show how error in the indices will bias coefficient estimates and limit the effectiveness of these forecasting methods. More constructively, the one-step strategy can be used to enhance the performance of these forecasting methods. Likewise, the industrial organization literature is increasingly using embedded representations of firms and products to characterize market behavior and demand with structural models. Our one-step strategy can be used to mitigate bias introduced by measurement error in the embeddings in these. Going forward, it is important to establish for which specific IR methods and econometric models does measurement error most severely affect inference. More generally, our belief is that inference problems arising from the analysis of unstructured data should be better recognized and taken more seriously in order to fully harness its potential value.

The rest of the paper proceeds as follows. Section 2 provides a simple setting in which the inference problems associated with the two-step strategy emerge. Section 3 further develops these arguments and presents our main theoretical results. Section 4 discusses instead the one-step strategy, associated computational tools, and introduces the Supervised Topic Model with Covariates. Section 5 presents simulation and empirical results comparing the two strategies. Section 6 concludes.

2 Motivating Example

This section presents a stylized model to illustrate clearly how the standard two-step strategy leads to biased inference in both the downstream and upstream models. The main take-aways from the stylized model are borne out in our empirical application.

2.1 Stylized Model

The stylized model is loosely based on the seminal work of Baker et al. (2016), which develops text-based measures of economic policy uncertainty (EPU) and investigates the relationship between EPU indices and economic outcomes. Suppose we are interested in the effect of θ_i (policy uncertainty in month i) on Y_i (employment or investment, say, in month $i + 1$). We are primarily concerned with inference on γ_1 in the regression model

$$Y_i = \gamma_0 + \gamma_1 \theta_i + \varepsilon_i. \quad (1)$$

Policy uncertainty itself is a nebulous concept that is difficult to precisely define let alone observe. The key innovation of Baker et al. (2016) is to construct EPU indices based on monthly counts of articles in ten newspapers containing certain terms, then convert to index form. Their EPU index is then introduced as a covariate in regressions and VARs. But it's arguably the case that their measure, while a strong signal of policy uncertainty, is not numerically the same as policy uncertainty. For instance, one could change the set of newspapers surveyed and obtain a quantitatively different (but related) measure. We therefore adopt the specification

$$X_i \sim \text{Binomial}(C_i, \theta_i), \quad (2)$$

where X_i is the number of counts observed out of a sample of size C_i and θ_i is the rate at which counts are expected. In the terminology of Baker et al. (2016), X_i is the number of articles containing certain key terms in month i , C_i is the total number of articles that month, and θ_i is policy uncertainty that month. The variables X_i , Y_i , and C_i are observed but θ_i is not. One can estimate θ_i using $\hat{\theta}_i = X_i/C_i$, which is what Baker et al. (2016) do to construct their policy uncertainty measure.⁶

To facilitate the theoretical derivations below, let $\mathbb{E}[\varepsilon_i | \theta_i, X_i, C_i] = 0$ and $\text{Var}(\theta_i) > 0$, so the OLS estimator of γ_1 would be consistent if θ_i were observed, and $\mathbb{E}[\varepsilon_i^2] < \infty$. To simplify derivations, we also assume (i) Y_i and (X_i, C_i) are independent conditional on θ_i , and (ii) C_i and θ_i are independent. These assumptions, which are credible in the context of Baker et al. (2016), are made primarily for convenience and can be relaxed. We assume

⁶See p. 1599 of Baker et al. (2016).

the data are a random sample $(X_i, Y_i, C_i)_{i=1}^n$. Our analysis and findings extend easily to time-series data, though we stick to the IID case to simplify presentation.

2.2 Two-Step Strategy

In the context of this example, the usual two-step strategy would regress Y_i on $\hat{\theta}_i$ and perform standard OLS inference for γ_1 . This approach overlooks the fact that $\hat{\theta}_i$ is a noisy estimate of θ_i . Failing to account for this measurement error problem may lead to biased estimates and inference.

Let $\hat{\gamma}_1$ denote the OLS estimator of γ_1 from regressing of Y_i on $\hat{\theta}_i$. By standard OLS algebra, as the sample size $n \rightarrow \infty$ we have

$$\begin{aligned} \hat{\gamma}_1 &\xrightarrow{p} \gamma_1 \frac{\text{Cov}(\theta_i, \hat{\theta}_i)}{\text{Var}(\hat{\theta}_i)} \\ &= \gamma_1 \frac{\text{Var}(\theta_i)}{\text{Var}(\theta_i) + \mathbb{E}[C_i^{-1}] \mathbb{E}[\theta_i(1 - \theta_i)]}, \end{aligned}$$

because $\mathbb{E}[\hat{\theta}_i | \theta_i, C_i] = \theta_i$ and $\text{Var}(\hat{\theta}_i) = \text{Var}(\theta_i) + \mathbb{E}[C_i^{-1}] \mathbb{E}[\theta_i(1 - \theta_i)]$ by the law of total variance and independence of C_i and θ_i . Evidently, there is an attenuation bias caused by measurement error in $\hat{\theta}_i$ which makes $\hat{\gamma}_1$ inconsistent.

The key determinant of bias is the average reciprocal amount of unstructured data per observation $\mathbb{E}[C_i^{-1}]$. If the amount of unstructured data per observation is large so that $\mathbb{E}[C_i^{-1}]$ is small, we have

$$\text{plim}(\hat{\gamma}_1) \approx \gamma_1 - \mathbb{E}\left[\frac{1}{C_i}\right] \frac{\mathbb{E}[\theta_i(1 - \theta_i)]}{\text{Var}(\theta_i)} \gamma_1$$

because $(1 + x)^{-1} \approx 1 - x$ for small x . Hence, the bias is of the order of $\mathbb{E}[C_i^{-1}]$.

In many empirical settings, both measurement error and sampling error may play important roles. To shed light on the behavior of $\hat{\gamma}_1$ in this scenario, we consider a sequence of populations indexed by the sample size n . The distribution of (Y_i, X_i, θ_i) conditional on C_i is fixed but the distribution of C_i is changing with n so that

$$\sqrt{n} \times \mathbb{E}\left[\frac{1}{C_i}\right] \rightarrow \kappa \in [0, \infty). \quad (3)$$

This should not be interpreted literally as the data-generating process. Rather, it is a thought experiment to provide insights about how $\hat{\gamma}_1$ behaves when both measurement and sampling error are present. The parameter κ controls the relative importance of measurement error and sampling error: $\kappa = 0$ means sampling error swamps measurement error, while larger κ gives relatively greater importance to measurement error.

Proposition 1. *Consider the sequence of populations just described. Then*

$$\sqrt{n}(\hat{\gamma}_1 - \gamma_1) \rightarrow_d N \left(-\kappa \gamma_1 \frac{\mathbb{E}[\theta_i(1 - \theta_i)]}{\text{Var}(\theta_i)}, \frac{\mathbb{E}[\varepsilon_i^2(\theta_i - \mathbb{E}[\theta_i])^2]}{\text{Var}(\theta_i)^2} \right).$$

Proposition 1 shows that two-step inference is *valid* when $\kappa = 0$. In this case, measurement error vanishes faster than sampling error and the estimated $\hat{\theta}_i$ can be treated as if they are the true θ_i .

However, Proposition 1 also shows that two-step inference is *invalid* when $\kappa > 0$. In this case, $\hat{\gamma}_1$ is consistent and its asymptotic variance is the same as if Y_i were regressed on the true θ_i , but the center of the asymptotic distribution is shifted due to the effect of measurement error. Confidence intervals based on standard OLS inference will therefore have approximately correct width but incorrect centering, meaning that their coverage rates will be below nominal coverage.⁷

2.3 Upstream Inference

So far we have focused on the “downstream” regression model. Other research questions might involve inference in an “upstream” model linking variation in θ_i (policy uncertainty) to variation in an observed covariate Z_i (legislative gridlock, say). In that context, θ_i or some transformation of θ_i is the dependent variable in a regression on Z_i . Because θ_i is not observed, the two-step strategy would replace θ_i with $\hat{\theta}_i$ in the regression. As before, the two-step strategy causes a measurement error problem, but now one that affects the *dependent* variable rather than the independent variable. As the measurement error $\hat{\theta}_i - \theta_i$ is uncorrelated with Z_i , there would be no bias if $\hat{\theta}_i$ were regressed on Z_i . But there can be a bias if a nonlinear transformation of $\hat{\theta}_i$ is used as the dependent variable.

To illustrate this, consider the following setup. Because θ_i is supported on $[0, 1]$ it is natural to transform it to have support \mathbb{R} using the log-odds ratio (or similar). Suppose we are concerned with inference on ϕ_1 in the regression model

$$\log \left(\frac{\theta_i}{1 - \theta_i} \right) = \phi_0 + \phi_1 Z_i + u_i.$$

We again assume $\mathbb{E}[u_i | Z_i] = 0$ so that OLS would be unbiased if the true θ_i were observed. Because θ_i is latent, one could instead regress the empirical log odds ratio

$$\log \left(\frac{\hat{\theta}_i}{1 - \hat{\theta}_i} \right)$$

⁷It follows from the general treatment in Section 3 that Eicker–Huber–White standard errors computed from $\hat{\theta}_i$ instead of θ_i are consistent.

on Z_i . Let $\hat{\phi}_1$ denote the corresponding OLS estimator. To understand the forces at play, we study the behavior of $\hat{\phi}_1$ in a sequence of populations where the distribution of $(X_i, Y_i, Z_i, \theta_i)$ conditional on C_i is fixed but the distribution of C_i varies with n so that (3) holds. Like before, to facilitate derivations we assume (X_i, C_i) and Z_i are independent conditional on θ_i , and C_i and (θ_i, Z_i) are independent.

Proposition 2. *Suppose that Assumption 3 in Appendix A holds. Then*

$$\sqrt{n}(\hat{\phi}_1 - \phi_1) \rightarrow_d N \left(\kappa \frac{\text{Cov} \left(\frac{2\theta_i - 1}{2\theta_i(1-\theta_i)}, Z_i \right)}{\text{Var}(Z_i)}, \frac{\mathbb{E} [u_i^2(Z_i - \mathbb{E}[Z_i])^2]}{\text{Var}(Z_i)^2} \right).$$

Proposition 2 shows that two-step inference in the upstream model is valid when $\kappa = 0$ but invalid when $\kappa > 0$. In the latter case, confidence intervals based on standard OLS inference will again have approximately correct width but incorrect centering, and will therefore have coverage below nominal coverage. The degree to which standard OLS confidence intervals under-cover depends partly on the size of $\text{Cov}(\frac{2\theta_i - 1}{2\theta_i(1-\theta_i)}, Z_i)$. Because the function $x \mapsto \frac{2x-1}{2x(x-1)}$ diverges to $\pm\infty$ as x approaches 0 and 1, this covariance can be very large when the distribution of θ_i puts mass near zero and/or one. Thus, first-order bias can be large even when κ is small provided θ_i has sufficient mass in its tails.

3 Full Analysis of the Two-Step Strategy

One limiting feature of the stylized model in the previous section is that the observed data are not high-dimensional. In this section, we allow each unstructured observation to lie in a high-dimensional space, which requires some dimensionality reduction prior to regression analysis. We first describe the statistical framework linking unstructured data and the downstream regression model. We then analyze the two-step strategy and show why it leads to biased inference in empirically plausible settings.

3.1 Statistical Framework

We begin by specifying a statistical model that, broadly speaking, has two parts. The first computes low-dimensional numerical representations of the unstructured data. The second introduces these numerical representations as covariates, potentially along with other quantitative data, into a linear regression model. There is a wide array of methods for dimensionality reduction used in the literature. We focus on factor modeling, which, in the context of high-dimensional discrete data, is also called topic modeling. We make this choice for two reasons. First, topic models have a well-defined statistical structure which facilitates theoretical analysis. Second, there is a large empirical literature which

uses topic models as part of the two-step strategy, most commonly using textual data. Examples include Hansen et al. (2018), Mueller and Rauh (2018), Larsen and Thorsrud (2019), Thorsrud (2020), Bybee et al. (2020), Adams et al. (2021), and Ash et al. (2022). However, the use of topic models is not limited to textual data. For example, Bandiera et al. (2020), Draca and Schwarz (2021), and Munro and Ng (2022) use topic models to analyze survey data. Meanwhile, Nimczik (2017) and Olivella et al. (2021) use topic models for network-structured data.

3.1.1 Model

Many types of unstructured data are high-dimensional and discrete. We therefore let each unstructured observation i be described by $\mathbf{x}_i = (x_{i,v})_{v=1}^V$, a V -dimensional vector of count variables. The value $x_{i,v}$ is the number of times a feature v appears in observation i . For example, in the bag-of-words model V is the number of unique terms in a textual corpus, typically in the thousands, and $x_{i,v}$ is the count of term v in document i .

The first part of the model generates a K -dimensional representation $\boldsymbol{\theta}_i$ of \mathbf{x}_i , where $K \ll V$. The second part introduces these low-dimensional representations as covariates, potentially along with other quantitative data \mathbf{q}_i , into a linear regression model:

$$Y_i = \boldsymbol{\gamma}^T(\mathbf{S}\boldsymbol{\theta}_i) + \boldsymbol{\alpha}^T\mathbf{q}_i + \varepsilon_i, \quad \mathbb{E}[\varepsilon_i | \boldsymbol{\theta}_i, \mathbf{q}_i] = 0, \quad (4)$$

where \mathbf{S} is a $H \times K$ selection matrix which is chosen by the researcher. Each row of \mathbf{S} selects a component of $\boldsymbol{\theta}_i$ to be included in the regression. In most applications in economics and finance, $\boldsymbol{\gamma}$ is the key parameter of interest. In other applications, however, $\boldsymbol{\alpha}$ is the focus and $\mathbf{S}\boldsymbol{\theta}_i$ plays the role of a text-derived control variable.

The model we consider for the unstructured data is widely used in practice and tractable enough that we can develop clean characterizations for the two-step strategy. As \mathbf{x}_i is a vector of counts, it is without loss of generality to model it as Multinomial. We impose some structure on the count probabilities for interpretability. The model is based on Probabilistic Latent Semantic Analysis (Hofmann 1999), a widely used factor model for discrete data, and its close cousin Latent Dirichlet Allocation (Blei et al. 2003, LDA).

There are K separate distributions over the V features denoted $\boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_K$ where each $\boldsymbol{\beta}_k$ lies in the $(V - 1)$ -dimensional simplex. The distributions are called *topics* in text applications. More generally, they represent common factors from which individual observations are built. We collect the factors into a $K \times V$ row-stochastic matrix \mathbf{B} where $\mathbf{B}^T = [\boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_K]$. Each observation i is characterized by the latent vector $\boldsymbol{\theta}_i$ which lies in the $(K - 1)$ -dimensional simplex. Its elements $\theta_{i,k}$ represent the weight attached to $\boldsymbol{\beta}_k$ in generating \mathbf{x}_i . Hence, the count probabilities for observation i are $\mathbf{p}_i = \sum_{k=1}^K \boldsymbol{\beta}_k \theta_{i,k}$.

Note that \mathbf{B} is a matrix of common parameters whereas $\boldsymbol{\theta}_i$ is an observation-specific latent random vector. Putting these elements together, the distribution of \mathbf{x}_i can be written as

$$\mathbf{x}_i | (C_i, \boldsymbol{\theta}_i) \sim \text{Multinomial}(C_i, \mathbf{B}^T \boldsymbol{\theta}_i), \quad (5)$$

where $C_i = \sum_{v=1}^V x_{i,v}$ is the count of all features in observation i —a measure of the amount of unstructured data for observation i —and the count probabilities \mathbf{p}_i have the factor structure $\mathbf{p}_i = \mathbf{B}^T \boldsymbol{\theta}_i$. This model nests as a special case a *pure multinomial* model where $K = V$, $\mathbf{B} = \mathbf{I}$, and $\boldsymbol{\theta}_i = \mathbf{p}_i$. The quantity C_i determines the degree of precision with which we can infer $\boldsymbol{\theta}_i$ from \mathbf{x}_i . The interplay between C_i and the number of observations n plays an important role in our theory.

Example: Monetary Policy Speeches. Suppose each unstructured observation is a monetary policy speech. One distribution β_k might put high weight on words like ‘inflation’, ‘prices’, and ‘cpi’, so β_k would have an interpretation as price rises. The corresponding $\theta_{i,k}$ represents how much speech i discusses price rises. One research question might ask how attention paid to price rises, along with other economic conditions captured by other topics and numeric data, affects policy actions. This could be captured by the γ coefficients in (4) where Y_i is the policy action of speaker i , \mathbf{S} selects the price rises topic weight from $\boldsymbol{\theta}_i$ and discards irrelevant topics (e.g., words used in generic conversation), and \mathbf{q}_i measures quantitative information like market forecasts for growth and inflation at the time the speech was made.

The main point beyond this specific example is that many research questions that seek to map variation across high-dimensional count observations as captured by a topic model into variation in some numeric variable will involve inference on $\boldsymbol{\gamma}$.

3.1.2 Data and Maintained Assumptions

The data are a random sample $(Y_i, \mathbf{q}_i, \mathbf{x}_i, C_i)_{i=1}^n$ satisfying (4) and (5). To simplify derivations, we further assume that C_i is independent of $(\boldsymbol{\theta}_i, \mathbf{q}_i, Y_i)$, and that (\mathbf{q}_i, Y_i) and \mathbf{x}_i are independent conditional on $(C_i, \boldsymbol{\theta}_i)$. We also assume that \mathbf{B} is identified. That is, there is a unique decomposition $\mathbf{P} = \mathbf{B}^T \boldsymbol{\Theta}$ with $\mathbf{P} = [\mathbf{p}_1, \dots, \mathbf{p}_n]$ collecting the vectors of count probabilities across observations and $\boldsymbol{\Theta} = [\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_n]$ collecting the topic weights across observations. Identification is commonly achieved in text applications by assuming the existence of anchor words that are known to appear in some topics but not others. We assume identifiability because our objective is to analyze the consequences of the two-step inference approach in a transparent way. Adding partial identification will significantly complicate the analysis but may be an interesting extension in future research.

3.2 Theory for the Two-Step Strategy

The standard two-step strategy can be summarized as follows:

- (i) Estimates $\hat{\boldsymbol{\theta}}_i$ of $\boldsymbol{\theta}_i$ are computed from the unstructured observations, e.g. by LDA.
- (ii) Y_i is regressed on $\mathbf{S}\hat{\boldsymbol{\theta}}_i$ and \mathbf{q}_i . Inference is performed treating the $\hat{\boldsymbol{\theta}}_i$ as if they are regular numeric data.

Evidently there is a measurement error problem: the estimates $\hat{\boldsymbol{\theta}}_i$ are noisy proxies for the true $\boldsymbol{\theta}_i$ appearing in the regression model (4). But Step (ii) overlooks this problem and treats the first-stage estimates $\hat{\boldsymbol{\theta}}_i$ as regular numeric data. This raises the possibility that OLS estimates may be biased due to measurement error introduced in Step (i). Moreover, conventional standard errors are typically reported. But these do not account for any additional variation introduced by using $\hat{\boldsymbol{\theta}}_i$ instead of $\boldsymbol{\theta}_i$.

We now analyze the two-step strategy and show how it can lead to biased estimates and inference. Let

$$\boldsymbol{\xi}_i = \begin{bmatrix} \mathbf{S}\boldsymbol{\theta}_i \\ \mathbf{q}_i \end{bmatrix}, \quad \hat{\boldsymbol{\xi}}_i = \begin{bmatrix} \mathbf{S}\hat{\boldsymbol{\theta}}_i \\ \mathbf{q}_i \end{bmatrix}.$$

The OLS estimator of $\boldsymbol{\psi} = [\boldsymbol{\gamma}, \boldsymbol{\alpha}]^T$ in the two-step strategy is given by

$$\hat{\boldsymbol{\psi}} = \left(\frac{1}{n} \sum_{i=1}^n \hat{\boldsymbol{\xi}}_i \hat{\boldsymbol{\xi}}_i^T \right)^{-1} \left(\frac{1}{n} \sum_{i=1}^n \hat{\boldsymbol{\xi}}_i Y_i \right). \quad (6)$$

3.2.1 Fixed Population

We first consider the large-sample properties of $\hat{\boldsymbol{\psi}}$ where the number of observations becomes large ($n \rightarrow \infty$) but the distribution of $(Y_i, \mathbf{q}_i, \mathbf{x}_i, C_i)_{i=1}^n$ is held fixed. This asymptotic framework mimics many empirical designs where there is a relatively small amount of unstructured data per observation but a large number of observations.

There are many different ways of estimating \mathbf{B} and $\boldsymbol{\theta}_i$ in (5). For instance, one could use LDA (Blei et al. 2003) or more recent methods developed by Bing et al. (2020), Wu et al. (2023), Ke and Wang (2022), and many others. As our objective is to focus on the consequences of the two-step strategy, we abstract from algorithmic-specific details and instead impose mild conditions on the estimators $\hat{\mathbf{B}}$ of \mathbf{B} and $\hat{\boldsymbol{\theta}}_1, \dots, \hat{\boldsymbol{\theta}}_n$ of $\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_n$. Let $\hat{\mathbf{p}}_i = \mathbf{x}_i/C_i$, $i = 1, \dots, n$. Let \rightarrow_p denote convergence in probability as $n \rightarrow \infty$.

Assumption 1. (i) \mathbf{B} has full rank.

(ii) $\hat{\mathbf{B}} \rightarrow_p \mathbf{B}$.

(iii) $\max_{1 \leq i \leq n} \|\hat{\boldsymbol{\theta}}_i - (\hat{\mathbf{B}}\hat{\mathbf{B}}^T)^{-1}\hat{\mathbf{B}}\hat{\mathbf{p}}_i\| \rightarrow_p 0$.

(iv) $\mathbb{E} [\|\mathbf{q}_i\|^2] < \infty$ and $\mathbb{E} [\boldsymbol{\xi}_i \boldsymbol{\xi}_i^T]$ has full rank.

Assumption 1(i) says that there are no fewer than K topics. This is a weak restriction, as K is typically much smaller than V in applications. Assumption 1(ii) says that $\hat{\mathbf{B}}$ is consistent for the topic weights \mathbf{B} . This is a mild condition satisfied by many estimators for topic models. Assumption 1(iii) imposes some structure on the estimators $\hat{\boldsymbol{\theta}}_i$. This condition is not vacuous: we have $\boldsymbol{\theta}_i = (\mathbf{B}\mathbf{B}^T)^{-1}\mathbf{B}\mathbf{p}_i$ (by Assumption 1(i)) so, given any consistent estimator $\hat{\mathbf{B}}$ of \mathbf{B} , one could estimate $\boldsymbol{\theta}_i$ simply by setting $\hat{\boldsymbol{\theta}}_i = (\hat{\mathbf{B}}\hat{\mathbf{B}}^T)^{-1}\hat{\mathbf{B}}\hat{\mathbf{p}}_i$. In that case, $\max_{1 \leq i \leq n} \|\hat{\boldsymbol{\theta}}_i - (\hat{\mathbf{B}}\hat{\mathbf{B}}^T)^{-1}\hat{\mathbf{B}}\hat{\mathbf{p}}_i\| = 0$. Assumption 1(i)-(iii) holds trivially for the pure multinomial model because $\hat{\mathbf{B}} = \mathbf{B} = \mathbf{I}$ and $\hat{\mathbf{p}}_i = \boldsymbol{\theta}_i$. Finally, Assumption 1(iv) imposes a standard moment condition on the numeric regressors \mathbf{q}_i and ensures that these are not perfectly collinear with the included latent variables $\mathbf{S}\boldsymbol{\theta}_i$.

Our first main result shows that $\hat{\boldsymbol{\psi}}$ is inconsistent when the amount of unstructured data per observation is small relative to the sample size. Let $\text{diag}(\mathbf{v})$ denote a diagonal matrix whose diagonal elements are the elements of the vector \mathbf{v} . Let $\mathbf{Q}_\mathbf{B} = (\mathbf{B}\mathbf{B}^T)^{-1}\mathbf{B}$ and let $\mathbf{0}$ denote a conformable matrix of zeros.

Theorem 1. *Suppose that Assumption 1 holds. Then*

$$\hat{\boldsymbol{\psi}} \rightarrow_p \left(\mathbb{E} [\boldsymbol{\xi}_i \boldsymbol{\xi}_i^T] + \begin{bmatrix} \mathbb{E} \left[\frac{1}{C_i} \right] \mathbf{S} (\mathbf{Q}_\mathbf{B} \text{diag}(\mathbf{B}^T \mathbb{E}[\boldsymbol{\theta}_i]) \mathbf{Q}_\mathbf{B}^T - \mathbb{E} [\boldsymbol{\theta}_i \boldsymbol{\theta}_i^T]) \mathbf{S}^T & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \right)^{-1} \mathbb{E} [\boldsymbol{\xi}_i \boldsymbol{\xi}_i^T] \boldsymbol{\psi}.$$

In particular,

$$\text{plim}(\hat{\boldsymbol{\psi}}) \approx \boldsymbol{\psi} - \mathbb{E} \left[\frac{1}{C_i} \right] \mathbf{b}$$

for $\mathbb{E} [C_i^{-1}]$ small, where

$$\mathbf{b} = \mathbb{E} [\boldsymbol{\xi}_i \boldsymbol{\xi}_i^T]^{-1} \begin{bmatrix} \mathbf{S} (\mathbf{Q}_\mathbf{B} \text{diag}(\mathbf{B}^T \mathbb{E}[\boldsymbol{\theta}_i]) \mathbf{Q}_\mathbf{B}^T - \mathbb{E} [\boldsymbol{\theta}_i \boldsymbol{\theta}_i^T]) \mathbf{S}^T \boldsymbol{\gamma} \\ \mathbf{0} \end{bmatrix}.$$

Theorem 1 shows that $\hat{\boldsymbol{\psi}}$ is inconsistent due to the measurement error in $\hat{\boldsymbol{\theta}}_i$. The reason is that each $\hat{\boldsymbol{\theta}}_i$ has a measurement error that doesn't disappear when the number of observations n becomes large. As a consequence, $\hat{\boldsymbol{\psi}}$ is asymptotically biased.

More constructively, Theorem 1 shows that bias is proportional to the average inverse amount of unstructured data per observation, $\mathbb{E} [C_i^{-1}]$. Because this is an inverse relationship, increasing the amount of unstructured data for observations with small C_i will reduce bias by more than if the additional data was collected for the observations with large C_i . Consequently, the across-observation distribution of C_i matters beyond its mean. In particular, bias may be large if most observations have large C_i but a small mass have small C_i .

3.2.2 Sequence of Populations

We now consider a sequence of populations where the amount of unstructured data per observation becomes larger as the sample size n increases. This asymptotic framework is designed to shed light on how $\hat{\psi}$ behaves in empirically realistic settings where there are a relatively large number of observations and there is a large amount of unstructured data per observation. In this scenario, the measurement errors across observations are small but their cumulative effect may not be ignorable relative to sampling error.

Formally, we consider a sequence of populations indexed by sample size n . In each population, we keep the distribution of $(\mathbf{x}_i, Y_i, \boldsymbol{\theta}_i, \mathbf{q}_i)$ conditional on C_i fixed and as described in Section 3.1. However, we let the marginal distribution of C_i change with the sample size n to allow the amount of unstructured data per observation to become large as the sample size n increases. Specifically, we consider a framework in which

$$\sqrt{n} \times \mathbb{E} \left[\frac{1}{C_i} \right] \rightarrow \kappa \in [0, \infty) \quad (7)$$

as $n \rightarrow \infty$. The quantity κ plays a key role in the following analysis. Loosely speaking, κ represents the relative magnitudes of sampling error and measurement error.

The case $\kappa = 0$ corresponds to a setting in which the amount of unstructured data per observation is of much larger order than sample size. Consequently, measurement error is of smaller order than sampling error. In this case, our theory implies that the two-step strategy leads to *valid* inference. That is, the measurement error in $\hat{\boldsymbol{\theta}}_i$ can effectively be ignored and standard inference can proceed treating the $\hat{\boldsymbol{\theta}}_i$ as if they are the true $\boldsymbol{\theta}_i$.

The case $\kappa \in (0, \infty)$ is the critical case in which there is a large, but not overwhelming, amount of unstructured data per observation. This case mimics many empirically realistic designs where measurement error and sampling error are both small but non-negligible. We show in this case that $\hat{\psi}$ is consistent but standard two-step inference is *invalid*. In particular, the asymptotic distribution of $\hat{\psi}$ has the correct variance but its center is shifted due to measurement error bias. Consequently, confidence intervals based on the usual two-step strategy have the correct width but incorrect centering, and therefore have a coverage rate that is smaller than nominal coverage.⁸

In what follows, notions of convergence in probability and distribution should be understood as holding along this sequence of populations satisfying (7).

Assumption 2. (i) \mathbf{B} has full rank.

(ii) $\sqrt{n}(\hat{\mathbf{B}} - \mathbf{B}) \rightarrow_p 0$.

⁸The case $\kappa = +\infty$ corresponds to a setting where measurement error is of larger order than sampling error. Here $\hat{\psi}$ is consistent provided $\mathbb{E}[C_i^{-1}] \rightarrow 0$ but two-step inference is invalid because bias is of *larger* order than sampling uncertainty. In that case, the coverage rates of standard OLS confidence intervals asymptote to zero as the sample size n becomes large.

(iii) $\sqrt{n} \max_{1 \leq i \leq n} \|\hat{\boldsymbol{\theta}}_i - (\hat{\mathbf{B}}\hat{\mathbf{B}}^T)^{-1}\hat{\mathbf{B}}\hat{\mathbf{p}}_i\| \rightarrow_p 0$.

(iv) $\mathbb{E}[\|\mathbf{q}_i\|^4] < \infty$ and $\mathbb{E}[\boldsymbol{\xi}_i \boldsymbol{\xi}_i^T]$ has full rank.

(v) $\mathbb{E}[\varepsilon_i^4] < \infty$.

(vi) $C_i \gtrsim (\log n)^{1+\epsilon}$ almost surely for some $\epsilon > 0$.

Assumption 2(i) is the same as Assumption 1(i). Assumption 2(ii)-(iii) strengthens Assumption 1(ii)-(iii) to require convergence at a faster-than-root- n rate. We believe Assumption 2(ii) is broadly satisfied in view of known convergence rates for estimators of \mathbf{B} .⁹ Assumption 2(iii) is made to simplify derivations but is not vacuous: given any estimator $\hat{\mathbf{B}}$, one could set $\hat{\boldsymbol{\theta}}_i = (\hat{\mathbf{B}}\hat{\mathbf{B}}^T)^{-1}\hat{\mathbf{B}}\hat{\mathbf{p}}_i$, in which case $\sqrt{n} \max_{1 \leq i \leq n} \|\hat{\boldsymbol{\theta}}_i - (\hat{\mathbf{B}}\hat{\mathbf{B}}^T)^{-1}\hat{\mathbf{B}}\hat{\mathbf{p}}_i\| = 0$. As before, Assumption 2(i)-(iii) holds trivially for the pure multinomial model because $\hat{\mathbf{B}} = \mathbf{B} = \mathbf{I}$ and $\hat{\mathbf{p}}_i = \hat{\boldsymbol{\theta}}_i$. Assumptions 2(iv) and 2(v) are standard. Assumption 2(vi) is made to simplify technical derivations and can be relaxed. It implies that C_i is supported on $[c(\log n)^{1+\epsilon}, \infty)$ for some $c, \epsilon > 0$. This is weaker than the conventional assumption that all C_i grow at the same rate C (Bing et al. 2020, Wu et al. 2023, Ke and Wang 2022) which, in view of (7), would imply that C_i is supported on $[cn^{1/2}, \infty)$ for some $c > 0$. This condition is only used to establish consistency of standard errors.

Our second main result shows that $\hat{\boldsymbol{\psi}}$ is consistent, derives its asymptotic distribution, and establishes consistency of standard errors. Recall the definition of \mathbf{b} from Theorem 1. Let $\hat{\varepsilon}_i = Y_i - \hat{\boldsymbol{\xi}}_i^T \hat{\boldsymbol{\psi}}$.

Theorem 2. *Suppose that Assumption 2 holds. Then*

$$\sqrt{n}(\hat{\boldsymbol{\psi}} - \boldsymbol{\psi}) \rightarrow_d N\left(\kappa \mathbf{b}, \mathbb{E}[\boldsymbol{\xi}_i \boldsymbol{\xi}_i^T]^{-1} \mathbb{E}[\varepsilon_i^2 \boldsymbol{\xi}_i \boldsymbol{\xi}_i^T] \mathbb{E}[\boldsymbol{\xi}_i \boldsymbol{\xi}_i^T]^{-1}\right), \quad (8)$$

and

$$\left(\frac{1}{n} \sum_{i=1}^n \hat{\boldsymbol{\xi}}_i \hat{\boldsymbol{\xi}}_i^T\right)^{-1} \left(\frac{1}{n} \sum_{i=1}^n \hat{\varepsilon}_i^2 \hat{\boldsymbol{\xi}}_i \hat{\boldsymbol{\xi}}_i^T\right) \left(\frac{1}{n} \sum_{i=1}^n \hat{\boldsymbol{\xi}}_i \hat{\boldsymbol{\xi}}_i^T\right)^{-1} \rightarrow_p \mathbb{E}[\boldsymbol{\xi}_i \boldsymbol{\xi}_i^T]^{-1} \mathbb{E}[\varepsilon_i^2 \boldsymbol{\xi}_i \boldsymbol{\xi}_i^T] \mathbb{E}[\boldsymbol{\xi}_i \boldsymbol{\xi}_i^T]^{-1}. \quad (9)$$

Theorem 2 shows that two-step inference is valid when $\kappa = 0$. In this case, we have

$$\sqrt{n}(\hat{\boldsymbol{\psi}} - \boldsymbol{\psi}) \rightarrow_d N\left(\mathbf{0}, \mathbb{E}[\boldsymbol{\xi}_i \boldsymbol{\xi}_i^T]^{-1} \mathbb{E}[\varepsilon_i^2 \boldsymbol{\xi}_i \boldsymbol{\xi}_i^T] \mathbb{E}[\boldsymbol{\xi}_i \boldsymbol{\xi}_i^T]^{-1}\right).$$

Here $\hat{\boldsymbol{\psi}}$ has the same asymptotic distribution as the (infeasible) OLS estimator obtained by regressing Y_i on the true latent $\boldsymbol{\theta}_i$. Moreover, Eicker–Huber–White standard errors

⁹Bing et al. (2020), Wu et al. (2023), and Ke and Wang (2022) derive finite-sample guarantees for different estimators $\hat{\mathbf{B}}$ of \mathbf{B} . Each of their results implies the corresponding estimator $\hat{\mathbf{B}}$ converges at the optimal rate $(nC)^{-1/2}$ (up to log terms) where, for simplicity, the C_i are all of the same order C . Hence, all estimators $\hat{\mathbf{B}}$ converge faster than $n^{-1/2}$ when C grows with n , as we have here by (7).

computed using $\hat{\boldsymbol{\theta}}_i$ are consistent. The reason is that $\kappa = 0$ corresponds to a scenario where measurement error is of smaller order than sampling error. Hence, measurement error can effectively be ignored when performing inference on $\boldsymbol{\psi}$.

At an abstract level, this case is analogous to asymptotic theory for factor-augmented regressions. In that setting, latent factors \mathbf{F}_t at each date are imputed from a vector of N predictor variables \mathbf{x}_t , then the estimated factors $\hat{\mathbf{F}}_t$ are treated as covariates in a regression model. Bai and Ng (2006) show that treating the estimated factors $\hat{\mathbf{F}}_t$ as if they are the true latent factors \mathbf{F}_t leads to valid inference provided $\sqrt{T}/N \rightarrow 0$, where T is the time-series dimension and N is the cross-sectional dimension. Their \mathbf{F}_t is analogous to our $\boldsymbol{\theta}_i$, their T is analogous to our n , and their $1/N$ is analogous to our $\mathbb{E}[C_i^{-1}]$. Hence, their condition $\sqrt{T}/N \rightarrow 0$ is analogous to $\kappa = 0$.

An important insight developed in Theorem 2 is that $\kappa = 0$ is in fact *necessary* for the validity of two-step inference. If $\kappa > 0$, then the asymptotic distribution of $\hat{\boldsymbol{\psi}}$ has the correct variance (which is consistently estimated by the Eicker–Huber–White estimator computed using $\hat{\boldsymbol{\theta}}_i$) but its center is shifted away from the origin to $\kappa\mathbf{b}$ due to measurement error bias.¹⁰ Consequently, confidence intervals have the correct width but incorrect centering, and therefore have coverage below their nominal coverage. The bias, and hence the degree to which confidence intervals under-cover, is increasing in κ . It is worth noting that measurement error in $\hat{\boldsymbol{\theta}}_i$ will bias inference not just for $\boldsymbol{\gamma}$ but also for $\boldsymbol{\alpha}$ whenever $\mathbf{S}\boldsymbol{\theta}_i$ and \mathbf{q}_i are correlated. The implication is that, even when unstructured data is used to create controls variables (see Avivi (2024) for a recent example), the two-step strategy can bias inference.

3.3 Extensions and Complements

We close this section discussing how our findings extend to other estimation problems and frameworks.

3.3.1 VARs

A number of prominent studies including the seminal work of Baker et al. (2016) have included text-derived measures as variables in vector autoregressions (VARs). These studies again use a two-step strategy: variables of interest $\mathbf{S}\boldsymbol{\theta}_t$ (policy uncertainty, say) are first estimated from unstructured data, then their estimates $\mathbf{S}\hat{\boldsymbol{\theta}}_t$ are plugged into a VAR as regular data together with other variables \mathbf{q}_t (interest rate, industrial production

¹⁰This is the opposite of a generated regressors problem (Pagan 1984), where the asymptotic variance is inflated but there is no location shift. With generated regressors there is a common finite-dimensional parameter estimated in the first stage whereas here all n covariates $\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_n$ are estimated in the first stage. See Bai and Ng (2006) for further discussion in the context of factor-augmented regressions.

unemployment, say). Standard inference on VAR parameters and/or impulse response functions (IRFs) is performed, again treating the estimates $\mathbf{S}\hat{\boldsymbol{\theta}}_t$ as regular data.

More formally, consider the VAR

$$\boldsymbol{\xi}_t = \boldsymbol{\psi}_0 + \boldsymbol{\Psi}_1 \boldsymbol{\xi}_{t-1} + \dots + \boldsymbol{\Psi}_p \boldsymbol{\xi}_{t-p} + \varepsilon_t,$$

where $\boldsymbol{\xi}_t = (\mathbf{S}\boldsymbol{\theta}_t^T, \mathbf{q}_t^T)^T$. Given data $(\mathbf{q}_t, \mathbf{x}_t, C_t)_{t=1}^n$, suppose that the $\boldsymbol{\theta}_t$ are estimated for $t = 1, \dots, n$, then the VAR parameters $\boldsymbol{\psi} = (\boldsymbol{\psi}_0, \boldsymbol{\Psi}_1, \dots, \boldsymbol{\Psi}_p)$ are estimated by regressing $\hat{\boldsymbol{\xi}}_t = (\mathbf{S}\hat{\boldsymbol{\theta}}_t^T, \mathbf{q}_t^T)^T$ on its lagged values and a constant. Let $\hat{\boldsymbol{\psi}}$ denote the OLS estimator.

The theory developed above carries over to the VAR setting. Suppose the process $(\boldsymbol{\theta}_1, \mathbf{q}_1, C_1), (\boldsymbol{\theta}_2, \mathbf{q}_2, C_2), \dots$ is strictly stationary and satisfies suitable weak dependence conditions; each \mathbf{x}_t depends on $(C_t, \boldsymbol{\theta}_t)$ as in (5) but is independent conditional on $(C_t, \boldsymbol{\theta}_t)$ of all other variables at all leads and lags; and the ε_t are martingale differences and satisfy suitable weak dependence conditions.

Partition the VAR coefficient matrices into blocks corresponding to whether (a) $\mathbf{S}\boldsymbol{\theta}_t$ or (b) \mathbf{q}_t is the dependent variable. Under suitable modification of Assumptions 1 and 2, OLS estimators of block (b) parameters behave essentially as described in Theorems 1 and Theorem 2. That is, they are inconsistent in a fixed-population setting and \sqrt{n} -consistent and asymptotically normal with a location shift proportional to κ in a sequence-of-populations framework in which

$$\sqrt{n} \mathbb{E} \left[\frac{1}{C_t} \right] \rightarrow \kappa \in [0, \infty).$$

OLS estimators of block (a) parameters have a measurement error $\mathbf{S}\hat{\boldsymbol{\theta}}_t - \mathbf{S}\boldsymbol{\theta}_t$ in the dependent variable also. But this measurement error is asymptotically negligible in a sequence-of-populations framework. As such, OLS estimators of block (a) parameters are \sqrt{n} -consistent and asymptotically normal with a location shift proportional to κ . Further, the asymptotic variance of $\hat{\boldsymbol{\psi}}$ will be the same as if the VAR was estimated on the true latent $\boldsymbol{\theta}_t$ and can be consistently estimated even when κ is positive, as in Theorem 2.

These consequences carry over to delta-method inference on functionals of $\boldsymbol{\psi}$, such as IRFs. That is, IRF estimators computed from $\hat{\boldsymbol{\psi}}$ will be \sqrt{n} -consistent and asymptotically normal with a location shift proportional to κ . As such, delta-method confidence intervals for IRFs using the two-step strategy will have the correct width but incorrect centering, and will therefore under-cover.

3.3.2 Similarity Measures

So far we have focused on regressions of numeric outcomes on topic weights. However, numeric outcomes are often regressed on similarity measures formed from term frequen-

cies or topic weights. In this case, the two-step strategy first estimates similarity from unstructured data then regresses numerical outcomes on the estimated similarity measures. We show that this strategy leads to biased inference on regression parameters unless sampling error dominates measurement error.

Suppose that the data are a random sample $(Y_i, \mathbf{x}_{1,i}, \mathbf{x}_{2,i}, C_{1,i}, C_{2,i})_{i=1}^n$, where

$$\mathbf{x}_{t,i} \sim \text{Multinomial}(C_{t,i}, \mathbf{p}_{t,i}), \quad t = 1, 2.$$

Each feature count vector $\mathbf{x}_{t,i}$ is a noisy signal of the true (latent) frequency $\mathbf{p}_{t,i}$. We are interested in performing inference on the parameter γ_1 in the regression model

$$Y_i = \gamma_0 + \gamma_1 (\mathbf{p}_{1,i} \cdot \mathbf{p}_{2,i}) + \varepsilon_i. \quad (10)$$

To give context, consider a setting loosely based on Kelly et al. (2021) in which Y_i denotes citations of patent i after it is filed, $\mathbf{x}_{1,i}$ is a vector of feature counts for patent i , and $\mathbf{x}_{2,i}$ is a vector of feature counts of an existing stock of patents at the time patent i was filed. The counts $\mathbf{x}_{1,i}$ and $\mathbf{x}_{2,i}$ are noisy signals of the true information contents $\mathbf{p}_{1,i}$ and $\mathbf{p}_{2,i}$ of the new patent and existing stock. The dissimilarity between $\mathbf{p}_{1,i}$ and $\mathbf{p}_{2,i}$ measures the novelty of patent i . We assume $\mathbb{E}[\varepsilon_i | \mathbf{p}_{1,i} \cdot \mathbf{p}_{2,i}] = 0$ and $\text{Var}(\mathbf{p}_{1,i} \cdot \mathbf{p}_{2,i}) > 0$ so that OLS regression of Y_i on $\mathbf{p}_{1,i} \cdot \mathbf{p}_{2,i}$ would be consistent if $\mathbf{p}_{1,i}$ and $\mathbf{p}_{2,i}$ were observed.

As $\mathbf{p}_{1,i}$ and $\mathbf{p}_{2,i}$ are not observed, a pragmatic two-step strategy is to estimate γ_1 by regressing Y_i on $(\hat{\mathbf{p}}_{1,i} \cdot \hat{\mathbf{p}}_{2,i})$ where

$$\hat{\mathbf{p}}_{1,i} = \frac{\mathbf{x}_{1,i}}{C_{1,i}} \quad \text{and} \quad \hat{\mathbf{p}}_{2,i} = \frac{\mathbf{x}_{2,i}}{C_{2,i}}$$

are the term frequencies. Let $\hat{\gamma}_1$ denote the OLS estimator. To simplify derivations, we assume (i) Y_i , $\mathbf{x}_{1,i}$, and $\mathbf{x}_{2,i}$ are independent conditional on $(C_{1,i}, C_{2,i}, \mathbf{p}_{1,i}, \mathbf{p}_{2,i})$, and (ii) $C_{1,i}$ and $C_{2,i}$ are independent of each other and of $(\mathbf{p}_{1,i}, \mathbf{p}_{2,i}, Y_i)$. Consider a sequence of populations in which the conditional distribution of Y_i , $\mathbf{x}_{1,i}$, $\mathbf{x}_{2,i}$, $\mathbf{p}_{1,i}$, and $\mathbf{p}_{2,i}$ conditional on $(C_{1,i}, C_{2,i})$ is fixed, and the distribution of $(C_{1,i}, C_{2,i})$ grows with n so that

$$\sqrt{n} \mathbb{E} \left[\frac{1}{C_{t,i}} \right] \rightarrow \kappa_t \in [0, \infty) \quad (11)$$

for $t = 1, 2$.

Theorem 3. *Consider the sequence of populations just described. Then*

$$\sqrt{n}(\hat{\gamma}_1 - \gamma_1) \rightarrow_d N \left(\kappa_1 b_1 + \kappa_2 b_2, \frac{\mathbb{E}[\varepsilon_i^2 (\mathbf{p}_{1,i} \cdot \mathbf{p}_{2,i} - \mathbb{E}[\mathbf{p}_{1,i} \cdot \mathbf{p}_{2,i}])^2]}{\text{Var}(\mathbf{p}_{1,i} \cdot \mathbf{p}_{2,i})^2} \right),$$

where

$$b_1 = - \left(\frac{\mathbb{E}[\mathbf{p}_{2,i}^T (\text{diag}(\mathbf{p}_{1,i}) - \mathbf{p}_{1,i} \mathbf{p}_{1,i}^T) \mathbf{p}_{2,i}]}{\text{Var}(\mathbf{p}_{1,i} \cdot \mathbf{p}_{2,i})} \right) \gamma_1, \quad b_2 = - \left(\frac{\mathbb{E}[\mathbf{p}_{1,i}^T (\text{diag}(\mathbf{p}_{2,i}) - \mathbf{p}_{2,i} \mathbf{p}_{2,i}^T) \mathbf{p}_{1,i}]}{\text{Var}(\mathbf{p}_{1,i} \cdot \mathbf{p}_{2,i})} \right) \gamma_1.$$

As with our earlier results, Theorem 3 shows that $\hat{\mathbf{p}}_{1,i}$ and $\hat{\mathbf{p}}_{2,i}$ can be treated as if they are the true $\mathbf{p}_{1,i}$ and $\mathbf{p}_{2,i}$ for inference on γ_1 provided the amount of unstructured data per observation is “large” in the sense that both $\kappa_1 = 0$ and $\kappa_2 = 0$. For instance, in the patent example we would require both the amount of unstructured data per patent to be large (so that $\kappa_1 = 0$) and the amount of unstructured data for the existing stock to be large (so that $\kappa_2 = 0$). Otherwise, the location of two-step confidence intervals is shifted towards the origin, which leads to biased inference.

We focused on the simplest case of regression on dot-product similarity between term frequencies to simplify exposition. But our findings will extend to regression on other measures including cosine similarity, TF-IDF measures, transforms (e.g., logs) of similarity measures, and similarity measures formed from topic weights.

4 One-Step Strategy

Having shown the bias inherent to the popular two-step approach for unstructured data analysis, we now develop an alternative one-step approach that allows for valid inference. The starting point is to note that the topic model (5) provides a likelihood for \mathbf{x}_i conditional on C_i and $\boldsymbol{\theta}_i$. We next build a joint likelihood for $(\mathbf{x}_i, Y_i, \boldsymbol{\theta}_i)$ by specifying a parametric distribution for the regression errors in (4). Finally, in the spirit of correlated random effects estimation in panel data models, we also specify a distribution for the topic shares $\boldsymbol{\theta}_i$ conditional on covariates \mathbf{g}_i . The covariates \mathbf{g}_i (J in total) may or may not be the same as \mathbf{q}_i . In practice, one must specify particular distributions to complete the specification of the likelihood, but for now we keep the discussion general to highlight the broad applicability of the approach.

Together, these components combine to give a likelihood $l(\mathbf{x}_i, Y_i, \boldsymbol{\theta}_i | C_i, \mathbf{g}_i, \mathbf{q}_i)$ for \mathbf{x}_i , Y_i , and $\boldsymbol{\theta}_i$ conditional on C_i and covariates \mathbf{g}_i and \mathbf{q}_i . As $\boldsymbol{\theta}_i$ is latent, we can integrate it out to obtain a likelihood $l(\mathbf{x}_i, Y_i | C_i, \mathbf{g}_i, \mathbf{q}_i)$ depending only on observable variables, which can then be used for maximum likelihood estimation of model parameters $\boldsymbol{\delta}$, consisting of \mathbf{B} , $\boldsymbol{\gamma}$, $\boldsymbol{\alpha}$, and any other parameters in the regression error and topic share distributions. However, there are two challenges. First, the integration has no closed-form solution and so must be performed numerically. Moreover, this numerical integration is high-dimensional and must be done observation-by-observation. As such, standard likelihood-based estimation is not computationally feasible. In the remainder of this section, we discuss how we overcome this challenge and present a specific model—the Supervised

Topic Model with Covariates—that we use for the empirical results.

4.1 Inference Approach for One-Step Strategy

The inference approach we take, while frequentist, is instead based on Bayesian computation. The integration step is performed implicitly as part of the sampling procedure. Similar approaches are taken to deal with latent states in Bayesian estimation of DSGE models (Herbst and Schorfheide 2016). In this approach, we introduce a prior for the model parameters $\boldsymbol{\delta}$ and treat the latent $\boldsymbol{\theta}_i$ as “parameters” drawn from a distribution that potentially depends on covariates \mathbf{g}_i (as discussed above). We sample from the posterior distribution of $(\boldsymbol{\delta}, (\boldsymbol{\theta}_i)_{i=1}^n)$ conditional on the observed data $(\mathbf{x}_i, Y_i, C_i, \mathbf{g}_i, \mathbf{q}_i)_{i=1}^n$. The marginal draws for $\boldsymbol{\delta}$ represent draws from the posterior distribution for $\boldsymbol{\delta}$ based on the *integrated* likelihood $l(\mathbf{x}_i, Y_i | C_i, \mathbf{g}_i, \mathbf{q}_i)$.

It is important to emphasize that while our approach uses Bayesian computation, one does in fact perform valid *frequentist* inference on model parameters $\boldsymbol{\delta}$ using this method. The maximum likelihood estimator $\hat{\boldsymbol{\delta}}$ of $\boldsymbol{\delta}$ is asymptotically normal under standard regularity conditions (e.g., Theorem 5.41 of van der Vaart 1998). By the Bernstein–von Mises Theorem (see Theorem 10.1 of van der Vaart 1998 and discussion), the posterior mean $\bar{\boldsymbol{\delta}}$ of $\boldsymbol{\delta}$ is first-order asymptotically equivalent to the MLE $\hat{\boldsymbol{\delta}}$. Moreover, the posterior distribution of $\boldsymbol{\delta}$ is asymptotically normal with mean $\bar{\boldsymbol{\delta}}$ and variance (when appropriately scaled with n) equal to the asymptotic variance of the MLE. As such, Bayesian credible sets for $\boldsymbol{\delta}$ —or any of its components such as $\boldsymbol{\gamma}$ —are valid frequentist confidence sets with the desired asymptotic coverage. This approach is also *efficient* for inference on $\boldsymbol{\delta}$ and its components, as it is asymptotically equivalent to likelihood-based inference.

4.1.1 Hamiltonian Monte Carlo

Our problem is to sample from the posterior distribution $q(\boldsymbol{\zeta} | (\mathbf{x}_i, Y_i, C_i, \mathbf{g}_i, \mathbf{q}_i)_{i=1}^n)$ where $\boldsymbol{\zeta} = (\boldsymbol{\delta}, (\boldsymbol{\theta}_i)_{i=1}^n)$. To do so, we use Hamiltonian Monte Carlo (HMC), a modern Markov chain Monte Carlo (MCMC) algorithm that is particularly well-suited to high-dimensional models.¹¹ MCMC algorithms define a stochastic process, i.e., a Markov chain, whose ergodic distribution coincides with the posterior distribution one wishes to sample from. Samples from this Markov chain can be used to form estimates and credible sets/confidence intervals. Efficient MCMC algorithms have low autocorrelation across samples which improves the accuracy of the resulting estimates.

A popular and simple MCMC method is the Metropolis-Hastings (MH) algorithm. Note the posterior is proportional to $q_n(\boldsymbol{\zeta}) := q(\boldsymbol{\zeta}, (\mathbf{x}_i, Y_i)_{i=1}^n | (C_i, \mathbf{g}_i, \mathbf{q}_i)_{i=1}^n)$, which is

¹¹More in-depth overviews of HMC are provided in Neal (2012), Hoffman and Gelman (2014), and Betancourt (2018). We are not aware of the application of HMC to topic models in the literature.

formed by multiplying the likelihood by the prior. The MH algorithm generates samples from the posterior in two steps: (1) propose a new state ζ' from the current state ζ using a pre-specified proposal distribution; then (2) accept the new proposal with a probability that increases in the ratio $q_n(\zeta')/q_n(\zeta)$. A challenge in practice is that the proposal distribution must be chosen carefully to avoid slow convergence. Taking small steps in a random direction can have a high acceptance probability but also high autocorrelation across samples and slow convergence. Taking a large step in a random direction can drastically reduce $q_n(\zeta')$ and hence the acceptance probability.

The HMC algorithm addresses this problem by utilizing the geometry of q_n to propose distant states that nonetheless have high chance of acceptance. This is achieved by proposing a new state ζ' by following Hamiltonian dynamics for a certain number of steps, starting from the initial state ζ . This process is determined by the curvature of q_n , and so determining the path to follow requires evaluating the gradient of q_n with respect to the parameters ζ . The specific variant of HMC that we use is the No-U-Turn Sampler (Hoffman and Gelman 2014, NUTS). The intuitive idea of NUTS is to follow the Hamiltonian dynamics for a random number of steps, and to stop when the path starts to double back on itself. This is not only more efficient than following the dynamics for a fixed number of steps, but also avoids the need to specify the number of steps in advance.

4.1.2 Implementation with probabilistic programming

From an implementation perspective, an advantage of HMC is that it is amenable to probabilistic programming. This allows one to define a data generating process for a statistical model in computer code, after which sampling is performed “automatically” in the background by following a generic set of algorithmic procedures adapted to the given model. In practice, modern probabilistic programming libraries use automatic differentiation to compute the gradients of highly flexible families of densities. Furthermore, the density and gradient computations are typically parallelizable as they are additive with respect to the data points.¹² This facilitates the use of the same specialized hardware normally used for machine learning tasks.

NUTS is implemented in many probabilistic programming libraries, the most popular of which is Stan. For this paper, we instead use NumPyro (Phan et al. 2019), which utilizes a state-of-the-art automatic differentiation engine Jax (Bradbury et al. 2018) and allows users to easily deploy these computations to specialized hardware such as Graphical Processing Units (GPUs) and Tensor Processing Units (TPUs), resulting in a dramatic improvement in computation time. Furthermore, NumPyro is a Python library, not a standalone program, which means that it is easy to integrate with other libraries and

¹²More precisely, the logarithm of q_n is additive with respect to the data points, and the gradient of the logarithm of q_n is the sum of the gradients of the log-likelihood and the logarithm of the prior.

benefits from the host of functionalities that Python provides. This said, our goal is not to advocate for any particular library, but to demonstrate that software and hardware have evolved to a point that allows Bayesian computation to be performed at scale without the need to manually derive sampling equations.

4.2 Supervised Topic Model with Covariates

We now specify distributional assumptions to fully formulate a likelihood for empirical analysis. We assume the regression errors in (4) are normally distributed and that the distribution of $\boldsymbol{\theta}_i$ conditional on \mathbf{g}_i is logistic normal. These assumptions are made for illustrative purposes and applied researchers may modify them as desired. The resulting model, which we call the *Supervised Topic Model with Covariates* (STMC), is formalized in Model 1.

$$\begin{aligned} \boldsymbol{\theta}_i &\sim \text{LogisticNormal}(\boldsymbol{\Phi}\mathbf{g}_i, I_K\sigma_\theta^2) && \text{(Upstream Topic Model)} \\ \mathbf{x}_i &\sim \text{Multinomial}(C_i, \mathbf{B}^T\boldsymbol{\theta}_i) \\ Y_i &\sim \text{Normal}(\boldsymbol{\gamma}^T\mathbf{S}\boldsymbol{\theta}_i + \boldsymbol{\alpha}^T\mathbf{q}_i, \sigma_Y^2) && \text{(Downstream Regression Model)} \end{aligned}$$

Model 1: Supervised Topic Model with Covariates

The matrix $\boldsymbol{\Phi}$ is a $K \times J$ matrix of regression coefficients. The k th row of $\boldsymbol{\Phi}$, denoted ϕ_k , captures how variation in covariates maps to variation in the prevalence of the k th topic across observations. Hence, a number of research questions can be addressed by performing inference on $\boldsymbol{\Phi}$. Model 1 also introduces scale parameters σ_θ and σ_Y . While we have modeled the error terms in the downstream regression and upstream logistic normal as homoskedastic to simplify presentation, this can easily be relaxed.

Example: Monetary Policy Speeches (Continued). To return to the example of Section 3, the downstream regression model could capture how policymakers’ attention predicts policy actions controlling for economic conditions. But policymakers’ attention can itself be a function of speaker characteristics such as demographic variables, or past experience of economic conditions (Malmendier et al. 2021). Such variables would enter \mathbf{g}_i but arguably not directly affect policy decisions beyond their effect on attention; i.e., they would not enter \mathbf{q}_i .

To our knowledge, STMC is new in the literature. Roberts et al. (2014) presents a model in which a logistic normal distribution over $\boldsymbol{\theta}_i$ is parameterized by covariates but without a downstream regression. Blei and McAuliffe (2010) and Ahrens et al. (2021) present models in which linear combinations of topic shares explain a normally distributed

response variable, but do not allow covariates to enter the distribution over θ_i . As such, we view STMC as of independent interest in the literature on topic modeling, although its primary purpose is to provide an example in which dimensionality reduction and linear regression are part of the same joint model and one cares about doing valid inference on model parameters.

Following the literature on topic modeling, we specify the following standard prior distributions for model parameters:

$$\begin{aligned}
\beta_k &\sim \text{Dirichlet}(\eta) \quad \forall k \\
\phi_{j,k} &\sim \text{Normal}(0, \sigma_\phi^2) \quad \forall j, k \\
\gamma_k &\sim \text{Normal}(0, \sigma_\gamma^2) \quad \forall k \\
\alpha_m &\sim \text{Normal}(0, \sigma_\alpha^2) \quad \forall m \\
\sigma_Y &\sim \text{Gamma}(s_0, s_1)
\end{aligned}
\tag{Priors}$$

In total, the model has seven hyperparameters: the three σ^2 terms in (Priors) as well as σ_θ^2 in (Upstream Topic Model); the symmetric Dirichlet parameter η in (Priors); the two Gamma distribution parameters in (Priors).

Appendix C displays the NumPyro code needed to draw samples from the posterior distribution of STMC. The core code is only several dozen lines long, and individual elements can be quickly modified to specify alternative distributions or models. An interested researcher should be able to modify it as needed to accommodate different data; to test robustness of the conclusions to specifying alternative distributions for the data; or to test robustness with respect to choice of priors. The key is to avoid having to re-derive complex inference algorithms every time the model is adjusted, and this is precisely the main advantage of automatic inference methods.

5 Empirical Results

In theory, the one-step strategy should outperform the two-step strategy, but establishing the empirical relevance of the bias that the latter produces is clearly important. While our computational approach makes the one-step strategy straightforward to implement, easier still would be for applied researchers to continue to use off-the-shelf packages for information retrieval and then to import the outputs into familiar regression software. This section establishes that there is indeed a quantitatively meaningful difference in regression parameter estimates produced by the two methods, both in simulated and actual data. Moreover, the differences we observe are consistent with the key theoretical results established above. This highlights the broad relevance of the one-step strategy for the empirical literature.

In all exercises, we perform inference using Hamiltonian Monte Carlo applied to the Supervised Topic Model with Covariates with hyperparameters detailed in Appendix B. We choose $K = 2$, which implies that each observation’s topic share vector can be written $\boldsymbol{\theta}_i = (\theta_i, 1 - \theta_i)$. For the *one-step strategy*, we sample from the posterior distribution implied by the full structure of STMC. For the *two-step strategy*, we first sample from (Upstream Topic Model) and include only a constant in \mathbf{g}_i and thus ignore any dependence that exists between covariates and $\boldsymbol{\theta}_i$ in the upstream model. Including the constant allows for $\boldsymbol{\theta}_i$ to have an asymmetric prior. We use the sampled values of θ_i to compute an estimate $\hat{\theta}_i$ of the posterior mean. We then estimate the following regression models using HMC:

$$\log \left(\frac{\hat{\theta}_i}{1 - \hat{\theta}_i} \right) = \phi_0 + \boldsymbol{\phi}_1^T \mathbf{g}_i + u_i, \quad (12)$$

$$Y_i = \gamma_0 + \gamma_1 \hat{\theta}_i + \boldsymbol{\alpha}^T \mathbf{q}_i + \varepsilon_i, \quad (13)$$

where the error terms are normal. The prior distributions over the regression coefficients are the same in both strategies. This procedure is designed to emulate the typical approach in the empirical literature while ensuring that any observed differences between the two strategies are not driven by different inference methods or implicit modeling choices.

Finally, our focus here is on inference rather than identification. Ke et al. (2021) highlight that the parameters of topic models are generally set- rather than point-identified. To restore point identification, a common assumption in the machine learning literature is the existence of “anchor words” (Arora et al. 2012) which we adopt as explained below.¹³

5.1 Simulation

We start with a simulation exercise that compares the one- and two-step strategies in terms of (i) the evolution of the bias in regression coefficients across different values of κ and (ii) the coverage of confidence intervals. We simulate the data according to the data generating process described in Model (1).¹⁴ We conduct three sets of simulations. Within each set, the amount of unstructured data per observation is the same for all observations and equal to $C_i = C \in \{10, 25, 200\}$. Together with the total number of observations, $n = 10000$, this implies $\kappa \in \{10, 4, 0.5\}$, for the three sets of simulations, respectively.

¹³An alternative approach would be to dispense with the anchor words assumption, thereby allowing for the possibility of partial identification, and use an identification-robust method for constructing confidence sets based on the HMC draws as in Chen et al. (2018).

¹⁴We impose the anchor word assumptions in the simulation in the following way. We first draw $\boldsymbol{\beta}_1$ and $\boldsymbol{\beta}_2$ from symmetric $(V - 50)$ -dimensional Dirichlet priors. Then we insert 50 zeros into both $\boldsymbol{\beta}_1$ and $\boldsymbol{\beta}_2$ in such that there is no feature v where $\beta_{v,1} = \beta_{v,2} = 0$. Data is then simulated from these modified topic-feature distributions.

Table 1: Coverage Rates of Confidence Intervals

κ	(a) Coverage for γ				(b) Coverage for ϕ		
	2-Step	Bias-corrected	1-Step	Infeasible	2-Step	1-Step	Infeasible
10	0.575	0.095	0.955	0.955	0.000	0.920	0.975
4	0.635	0.915	0.965	0.955	0.000	0.955	0.975
0.5	0.910	0.935	0.960	0.955	0.025	0.965	0.975

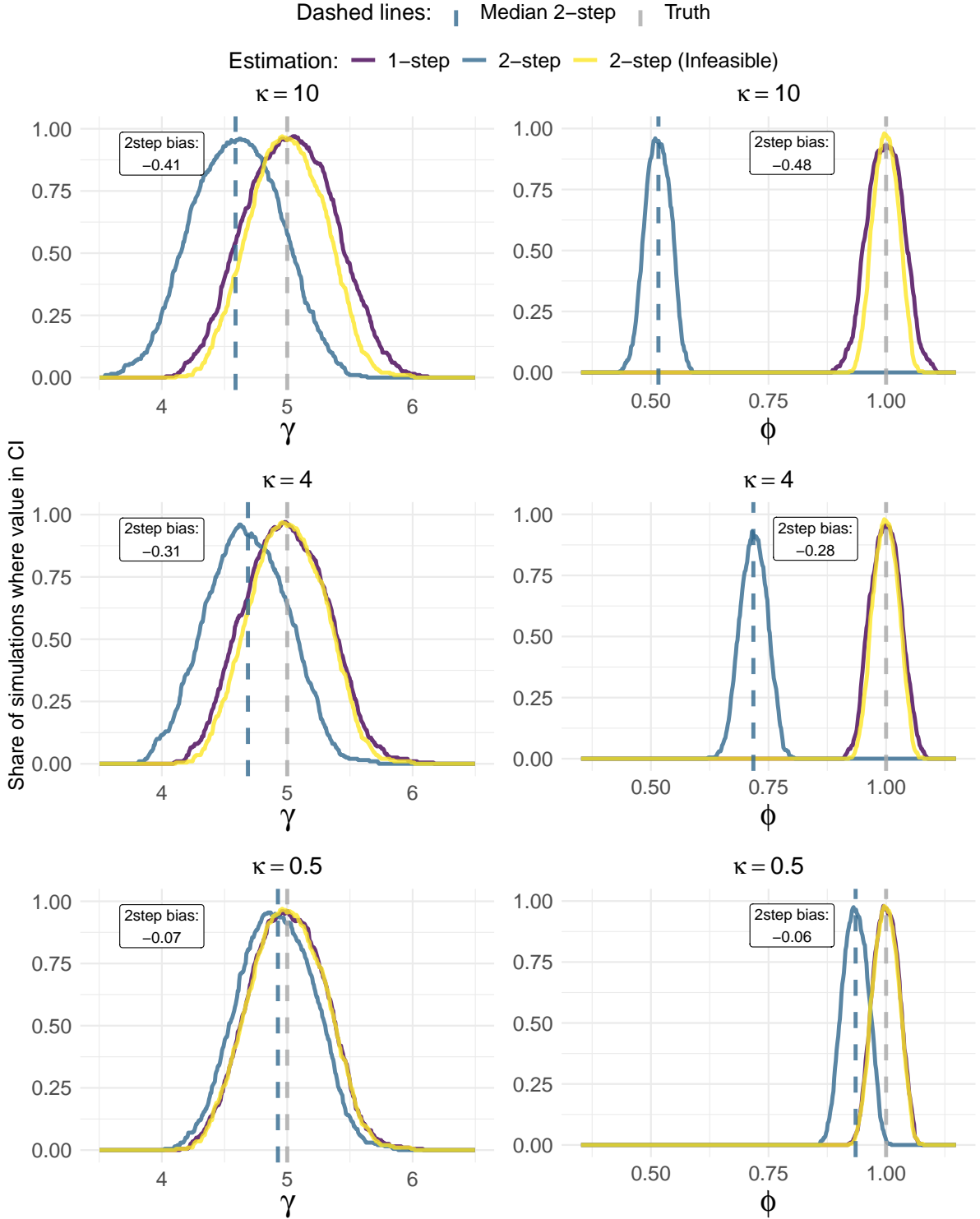
Note: This table reports the coverage rates of 95% confidence intervals for γ_1 and ϕ_1 across different values of κ . The values are reported for the two-step and one-step strategies, as well as for the infeasible estimator that uses the true θ_i . The bias-corrected estimates are obtained by subtracting the bias estimated using the formula from Theorem 1 from the two-step estimates. The coverage rates are computed as the share of simulations in which the true value of the parameter is included in the 95% confidence interval.

We conduct 200 simulations per set. Further details are included in Appendix B.

We focus on the estimation of two coefficients: (1) γ_1 , the effect of the increase in θ_i on Y_i ; and (2) ϕ_1 , the effect of a numerical covariate in (12). Our general theoretical results in Section 3 are directly applicable to two-step inference on γ_1 . Proposition 2 also shows that one should expect a bias for ϕ_1 that is also increasing in κ , and especially prominent when the distribution of θ_i has non-trivial mass at extreme values. To illustrate that the difference between the one-step and two-step strategies is due to mis-measurement of θ_i , we also estimate the regression coefficients using the true (known) values θ_i as an input instead of $\hat{\theta}_i$. This approach is, of course, not feasible in practice, but it allows us to isolate the effect of mis-measurement of θ_i on the regression coefficients.

Figure 1 presents the results. Each panel shows the coverage rates of confidence intervals for different parameter values: the share of simulations in which the values of the parameters are included in the 95% confidence intervals. The grey vertical dashed lines show the true value of the parameter. The blue vertical dashed line represents the median (across simulations) of mean posterior estimates for the two-step strategy. The two top panels show the results for the set of simulations where the amount of unstructured data per observation is the smallest and so $\kappa = 10$ is relatively large. The theory presented earlier suggests that in this case we should expect the two-step strategy to perform badly. This is indeed the case. The median (across simulations) estimate of γ_1 in top left, and ϕ_1 in top right, are both substantially biased towards zero. Further, as predicted by theory, the width of the CIs using the two-step strategy is similar to the infeasible estimator that uses the true θ_i . This, together with the bias, means that the CIs based on the two-step strategy under-cover. As reported in Table 1, for γ_1 the true value is included in the 95% CI in only 115/200 (57.5%) of simulations. For ϕ_1 this looks even worse: the true value is never included in the CIs.

On the other hand, the one-step strategy performs very well. The estimates appear



Note: Each mountain plot presents the share of simulations in which the value of γ_1 (respectively ϕ_1) on the x -axis is included in the 95% confidence interval. The grey vertical dashed lines show the true value of the parameter. The blue vertical dashed line represents the median (across simulations) of mean posterior estimates from the two-step strategy. The bias reported is the difference between the truth and this median value.

Figure 1: Evolution of Bias in Regression Coefficients across κ Values

unbiased, and the CIs have close to nominal coverage. The coverage is 95.5 and 92% respectively for γ_1 and ϕ_1 . The difference between the lengths of CIs using the one-step strategy and those using the infeasible estimator is small but noticeable. In the former, θ_i is recognized as latent and the uncertainty in θ_i is accounted for when performing inference on γ_1 and ϕ_1 . The resulting CIs are approximately 20% wider than those obtained with the infeasible estimator that uses the true θ_i .

Moving down from the top panels, we can see the evolution of bias and coverage as the amount of unstructured data per observation C increases and κ decreases. As predicted by theory, the bias in the two-step strategy becomes smaller as κ decreases. Increasing C from 10 in the top panels to 25 in the middle, substantially reduces the absolute value of median bias in the two-step estimate of γ_1 , while the width of a typical CI virtually does not change. This results in a noticeable increase in the coverage rate from 57.5 to 63.5%, but the coverage remains far from nominal. A similar pattern is observed for ϕ_1 . Meanwhile, the one-step strategy continues to perform well and one-step CIs are now hardly distinguishable from (infeasible) CIs based on the true θ_i . Finally, in the bottom panels, where $C = 200$ and $\kappa = 0.5$, the pattern continues. The bias in the two-step strategy is now very small for γ_1 and, consequently, coverage rate increases to 91%. A small, though still noticeable bias remains for ϕ_1 .

In principle, one can bias-correct the two-step CIs to restore valid inference. This is feasible because the model for unstructured data we use is sufficiently tractable that the bias of the OLS estimator of γ_1 can be characterized analytically, as in Theorem 1. By plugging in estimates for \mathbf{B} and Θ , we can correct for the asymptotic bias. Notably though, the formula for the asymptotic bias is based on the first-order approximation and so for large values of κ it is likely to perform poorly. This is indeed what we find. For small $\kappa = .5$ the median asymptotic bias for γ is -0.03 , compared with empirical median bias of -0.07 . Consequently, as reported in Table 1, applying the bias correction to 2-step estimate improves coverage rate from 91 to 93.5%. For $\kappa = 4$ the median estimated asymptotic bias is slightly larger than the empirical median bias (-0.36 vs. -0.31), but applying it dramatically improves the coverage rate from 63.5 to 91.5%. On the other hand, for $\kappa = 10$, the bias is poorly estimated and so the coverage rate actually decreases (from 57.5 to 9.5%) if the bias correction is applied¹⁵. This suggests that the bias correction is not a panacea and that the one-step strategy is a more reliable approach in practice. Moreover, more complex models for unstructured data may not admit analytical expressions for bias. The one-step strategy remains feasible in these cases provided the model for unstructured data is generative (i.e., has a likelihood).

Overall, the simulations confirm the three main insights from Theorem 2: (1) there is a first order bias in the two-step strategy, which is driven by the mis-measurement

¹⁵For more details on the bias correction, see Appendix B.1

of θ_i ; (2) the bias is larger when κ is larger; (3) the width of the confidence intervals is not substantially affected by mis-measurement of θ_i . The simulations also show that analytical bias-correction may be useful but only if κ is small, i.e. in the cases where the two-step strategy should already be performing relatively well. On the other hand, the one-step strategy performs well across all values of κ and is thus a more reliable approach in practice. The one-step strategy is not only theoretically sound, but also leads to substantially less biased inference in practice.

Finally, a word on the computational performance is in order. We have found that Numpyro’s HMC implementation of the STMC model is fast—each simulation took approximately 4 minutes when estimated on a single mid-range professional GPU, the Nvidia V100. As such, we think that the one-step strategy is feasible for most researchers, and that the computational cost is not a major concern.

5.2 CEO Behavior

To show that modeling joint dependence and estimating jointly matters in practice, we revisit the study of Bandiera et al. (2020), which collects and analyzes data on CEO time use in a sample of manufacturing firms in several countries. The goal of that paper is to describe salient differences in executive time use, and to relate those differences to firm and CEO characteristics as well as firm outcomes.

The estimation sample consists of 916 CEOs, each of whom participated in a survey that recorded features of time use in each 15-minute interval of a given week, e.g. Monday 8am-8:15am, Monday 8:15am-8:30am, and so forth. The recorded categories are (1) the type of activity (meeting, public event, etc.); (2) duration of activity (15m, 30m, etc.); (3) whether the activity is planned or unplanned; (4) the number of participants in the activity; (5) the functions of the participants in the activity (HR, finance, suppliers, etc.). In total there are 654 unique combinations of these categories observed in the data. We let $x_{i,j}$ denote the number of times feature j appears in the time use diary of CEO i . The average value of C_i is 88.4, with a minimum of 2 and a maximum of 222. Bandiera et al. (2020) uses LDA with $K = 2$ dimensions to organize the time use data. The authors refer to the separate distributions over time use combinations β_1 and β_2 as *pure behaviors*. The share of CEO i ’s time devoted to pure behavior 1, θ_i , is referred to as the *CEO index*.

The authors use the following inference procedure. First, estimate LDA on the time use data using the collapsed Gibbs sampler of Griffiths and Steyvers (2004), then form an estimate $\hat{\theta}_i$ based on the posterior means. They then use $\hat{\theta}_i$ as an input into productivity regressions where Y_i is the log of firm i sales, and \mathbf{q}_i is a vector of firm observables. Further, they separately analyze which CEO and firm characteristics are associated with behaviors by regressing $\hat{\theta}_i$ on a vector of characteristics \mathbf{g}_i .

We re-examine these questions using the Supervised Topic Model with Covariates. To explain CEO behavior, in \mathbf{g}_i we include log employment (a measure of firm size) and an indicator for whether the CEO has an MBA degree. To explain sales, in \mathbf{q}_i we include log employment and fixed effects for year and country. As before, we use HMC for inference and the same priors for both strategies.¹⁶ The priors used are the same as in the simulation exercise, except that we set the Dirichlet concentration parameter $\eta = 0.1$ to follow the original paper.

As demonstrated both theoretically and through the simulation exercise, the key quantity that governs the relative importance of sampling error and measurement error is κ . In the context of the CEO behavior data, the empirical analog of κ is the product of the square root of the number of observations (CEOs) and the average value of the inverse of the number of activities per CEO. This value is 0.44, which is close to the lowest value of κ in the simulation exercise. This suggests that the two-step approach should perform relatively well in this application. To further test our theory, we also estimate the model using data where we first sample 10% of the activities for each CEO, without replacement, with a minimum of one. This scenario could represent a researcher observing only half of a workday for each CEO, instead of a full five-day workweek. Such sampling increases the analogue of κ to 4.26, which is near the middle value of κ in the simulation exercise, indicating that we should expect the two-step approach to perform poorly under these conditions.

Turning to results, in Table 2 we report the relative probability of observing certain activities in Pure Behavior 1 relative to Pure Behavior 2. The table shows that estimated pure behaviors obtained with one-step and two-step strategies are very similar. What is more, they are also similar to those obtained with LDA and reported in the original paper. The table suggests that interacting with C-Suite executives, spending time communicating, and holding multi-function meetings are much more likely under Pure Behavior 1. Conversely, spending time on plant visits and interacting solely with suppliers are more likely under Pure Behavior 2. Based on these observations, the original authors label the CEOs with high values of $\hat{\theta}_i$ as *leaders* and those with low values as *managers*.

In terms of the regression coefficient estimates, we find patterns that are consistent with theory and the simulation results. In Table 3, we report the estimates of the regression coefficients under the two-step and one-step strategies. In Panel (a), we show the estimates for the downstream coefficient γ_1 , and in Panel (b), we show the estimates for the upstream coefficients ϕ_1 . In both panels, columns (1) and (2) report the estimates obtained using one- and two-step strategies, respectively, for the full sample. The coefficient on the CEO index in the downstream model is equal to 0.4 and 0.402, respectively,

¹⁶We impose the anchor word assumption by zeroing out from β_1 (β_2) the activity that is relatively least likely in Pure Behavior 1 (2).

Table 2: Comparison of Pure Behaviors

Activity	1-step	2-step	Bandiera et al (2020)
Plant Visits	0.1	0.09	0.11
Suppliers	0.61	0.74	0.32
Production	0.38	0.33	0.46
Just Outsiders	0.74	1.21	0.58
Communication	1.44	1.23	1.49
Multi-Function	1.35	1.12	1.9
Insiders and Outsiders	1.8	1.83	1.9
C-suite	29.78	16.76	33.9

Note: This table reports the relative probability of observing certain activities in Pure Behavior 1 relative to Pure Behavior 2. The value of 1 indicates that this activity is equally likely under both Pure Behaviors. Values higher than 1 mean that this type of activity is more likely to be performed under Pure Behavior 1. The values are reported in columns (1) and (2) are computed by first obtaining mean posterior probabilities of each activity in the given types. In column (3) we present values reported in Bandiera et al. (2020).

in the two strategies; the CIs have a similar length and exclude 0. Thus, both strategies suggest that a larger share of time spent on Pure Behavior 1 is associated with higher firm productivity. In the upstream model, we see larger differences between the two strategies as in the simulations. While having an MBA and managing larger firms are both associated with a higher CEO index, the point estimates differ substantially. As suggested in the simulations, there appears to be a downward bias in the two-step strategy: for instance, the coefficient on the MBA dummy is equal to 0.307 in the two-step strategy, compared to 0.606 in the one-step strategy. The CIs are marginally wider in the one-step strategy (0.297 vs. 0.261), but as the theory predicts, the difference is not substantial. Note there is no overlap in the CIs for these coefficients: the one-step CIs lie entirely to the right of the two-step CIs.

The differences between the strategies are substantially more pronounced when we consider the estimates obtained using the 10% subsample of unstructured data. Under the one-step strategy, the empirical conclusions are largely the same as when using the full data. For example, the point estimate on γ_1 changes from 0.402 to 0.439. While the confidence intervals are 54% wider than when using the full data (reflecting the increased uncertainty in estimated θ_i), there is still a strong estimated relationship between CEO behavior and firm performance. This is not so with the two-step strategy: the point estimate of γ_1 is now halved to 0.211, and the CI includes 0. Likewise, in the upstream model, the estimate of the coefficient on the MBA indicator remains large and statistically significant in the one-step strategy, but is reduced by 62% and is no longer statistically

Table 3: Regression Coefficient Estimates under Alternative Model Specifications

	Dependent variable: Log(sales)			
	(1) 2-Step	(2) 1-Step	(3) 2-Step	(4) 1-Step
CEO Index	0.4 (0.219, 0.572)	0.402 (0.240, 0.603)	0.211 (-0.028, 0.449)	0.439 (0.153, 0.711)
Log Employment	1.212 (1.159, 1.268)	1.198 (1.154, 1.248)	1.239 (1.186, 1.29)	1.199 (1.148, 1.26)
Controls	X	X	X	X
Activities' Sample	Full	Full	10%	10%

(a) Downstream Model: CEO Index and Firm Productivity

	Dependent variable: Un-normalized CEO index			
	(1) 2-Step	(2) 1-Step	(3) 2-Step	(4) 1-Step
MBA	0.307 (0.176, 0.437)	0.606 (0.446, 0.743)	0.118 (-0.012, 0.249)	0.323 (0.107, 0.486)
Log Employment	0.356 (0.306, 0.406)	0.492 (0.432, 0.548)	0.154 (0.104, 0.204)	0.443 (0.376, 0.507)
Controls	X	X	X	X
Activities' Sample	Full	Full	10%	10%

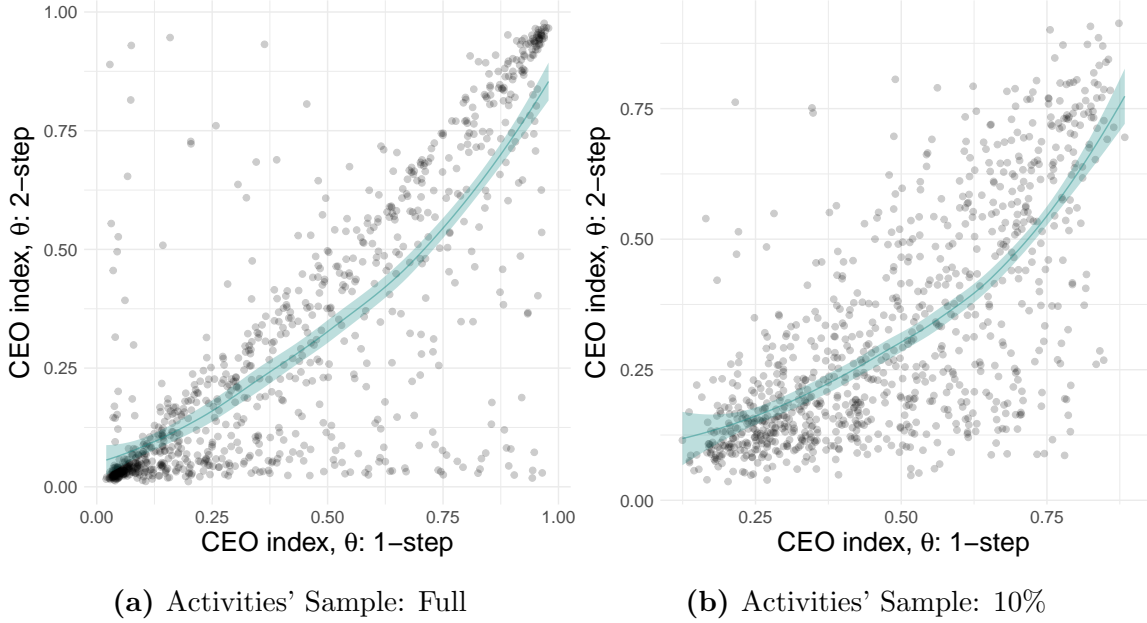
(b) Upstream Model: MBA and CEO Index

Note: In parentheses we report symmetric (equal-tailed) 95% confidence interval.

significant in the two-step strategy. This is consistent with the theory and simulation results, which suggest that the two-step strategy should perform particularly poorly in this scenario.

What explains the differences in estimates across strategies? To answer this question, we plot the estimated CEO indices in Figure 2. Panel (a) plots the estimated CEO indices obtained using the full sample, while Panel (b) plots the estimated CEO indices obtained using the 10% subsample. The blue line represents the local polynomial fit (with confidence intervals). The figure shows that when the full sample is used, both strategies find a large number of CEOs with $\hat{\theta}_i$ close to 0 and 1, and a strong correlation between the two estimates. However, the correlation is much weaker for the 10% subsample, suggesting that there is a large scope for mis-measurement of θ_i . Interestingly, Proposition 2 suggests that the bias in the two-step estimate of ϕ_1 can be severe when θ_i has mass near 0 and 1, as appears to be the case in this dataset. This provides an explanation for why the two-step strategy produces smaller estimates of ϕ_1 even in the full dataset.

Taken together, both the simulation results and the analysis of CEO behavior data highlight the importance of having a large amount of unstructured data per observation. Without it, the coefficients estimated using the two-step strategy can be badly biased,



Note: Each point represents the mean posterior estimate of a single CEO's index, $\hat{\theta}_i$. The blue line is the local polynomial fit (with confidence intervals) obtained with 'ggplots's' 'geom_smooth' with default parameters.

Figure 2: Scatterplots of Estimated CEO Indices $\hat{\theta}_i$

which can lead to incorrect empirical conclusions. The good statistical and computational performance of the one-step strategy makes it attractive to guard against this risk.

6 Conclusion

The leading strategy for analyzing unstructured data uses two steps. First, quantitative representations of unstructured data are extracted in an upstream information retrieval step. Second, the derived quantitative representations are plugged into downstream econometric models, with the representations treated as regular numerical data for the purposes of estimation and inference. This paper highlights, both theoretically and empirically, a previously unrecognized problem with this popular two-step strategy: measurement error introduced in the first step leads to biased estimates and invalid inference for downstream regression coefficients. The degree of bias, and therefore the degree to which it distorts inference, depends on the relative importance of measurement error and sampling error, but it can be material in applications. To guard against it, we propose a robust inference method based on joint maximum likelihood estimation of the IR and regression models. Joint estimation is straightforward using HMC and modern probabilistic programming languages. This strategy outperforms the two-step strategy in simulations and generates quantitatively important differences in a leading application.

Implementing the one-step strategy requires formulating a likelihood. Latest-generation machine-learning- and AI-based approaches to information retrieval increasingly use neural networks with no obvious structure that yields a likelihood. While implementing the one-step strategy may not be possible in these settings, the measurement error problem does not thereby disappear: it simply becomes harder to characterize. Such approaches are often given statistical foundations following their adoption and, as this process plays out, the scope for the one-step strategy will expand accordingly.¹⁷

Finally, we note there are limits on the scalability of HMC, even when fully optimized. When one confronts a vast amount of data, alternative approaches for approximating the joint distribution in the one-step strategy must be used. One popular choice in computer science is variational inference (VI; Jordan et al. 1998, Wainwright and Jordan 2008) which has recently seen applications in economics (Bonhomme 2021, Olenski and Sacher 2022, Mele and Zhu 2023). VI is no more complicated to implement because it too can be formulated within probabilistic programming languages that rely on automatic differentiation (Hoffman et al. 2013). However, VI has fewer statistical guarantees than HMC because it uses an approximate likelihood in place of the true likelihood. In ongoing work, we are studying how to best perform scalable inference in the one-step strategy with massive data.

¹⁷One illustrative example is the popular *word2vec* model for producing word embeddings. The original model (Mikolov et al. 2013b,a) had no statistical interpretation but yielded word representations that nevertheless captured semantic relationships well. Word2vec has subsequently been adopted by economists as part of the two-step strategy, for example to measure occupation-level exposure to technological change (Kogan et al. 2019) and emotionality in political speech (Gennaro and Ash 2022). In parallel, a literature has developed likelihood-based interpretations of embeddings (Arora et al. 2016, Dieng et al. 2020, Ruiz et al. 2020) which could in principle be adapted for use in the one-step strategy.

References

- Adams, R. B., Raganathan, V., and Tumarkin, R. (2021). Death by committee? An analysis of corporate board (sub-) committees. *Journal of Financial Economics*, 141(3):1119–1146.
- Ahrens, M., Ashwin, J., Calliess, J.-P., and Nguyen, V. (2021). Bayesian Topic Regression for Causal Inference. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 8162–8188, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Allon, G., Chen, D., Jiang, Z., and Zhang, D. (2023). Machine Learning and Prediction Errors in Causal Inference. *SSRN Electronic Journal*.
- Arora, S., Ge, R., Kannan, R., and Moitra, A. (2012). Computing a nonnegative matrix factorization – provably. In *Proceedings of the Forty-Fourth Annual ACM Symposium on Theory of Computing, STOC '12*, pages 145–162, New York, NY, USA. Association for Computing Machinery.
- Arora, S., Li, Y., Liang, Y., Ma, T., and Risteski, A. (2016). A Latent Variable Model Approach to PMI-based Word Embeddings. *Transactions of the Association for Computational Linguistics*, 4:385–399.
- Ash, E. and Hansen, S. (2023). Text Algorithms in Economics. *Annual Review of Economics*, 15(1):659–688.
- Ash, E., Morelli, M., and Vannoni, M. (2022). More Laws, More Growth? Evidence from U.S. States.
- Avivi, H. (2024). Are Patent Examiners Gender Neutral? Unpublished manuscript.
- Bai, J. and Ng, S. (2006). Confidence Intervals for Diffusion Index Forecasts and Inference for Factor-Augmented Regressions. *Econometrica*, 74(4):1133–1150.
- Baker, S. R., Bloom, N., and Davis, S. J. (2016). Measuring Economic Policy Uncertainty*. *The Quarterly Journal of Economics*, 131(4):1593–1636.
- Bandiera, O., Prat, A., Hansen, S., and Sadun, R. (2020). CEO Behavior and Firm Performance. *Journal of Political Economy*, 128(4):1325–1369.
- Bernanke, B. S., Boivin, J., and Eliasziw, P. (2005). Measuring the Effects of Monetary Policy: A Factor-Augmented Vector Autoregressive (FAVAR) Approach. *The Quarterly Journal of Economics*, 120(1):387–422.
- Betancourt, M. (2018). A Conceptual Introduction to Hamiltonian Monte Carlo. *arXiv:1701.02434 [stat]*.
- Bing, X., Bunea, F., and Wegkamp, M. (2020). Optimal estimation of sparse topic models. *Journal of Machine Learning Research*, 21:1–45.
- Bingham, E., Chen, J. P., Jankowiak, M., Obermeyer, F., Pradhan, N., Karaletsos, T.,

- Singh, R., Szerlip, P., Horsfall, P., and Goodman, N. D. (2018). Pyro: Deep Universal Probabilistic Programming. *arXiv:1810.09538 [cs, stat]*.
- Blei, D. M. and McAuliffe, J. D. (2010). Supervised Topic Models. *arXiv:1003.0783 [stat]*.
- Blei, D. M., Ng, A. Y., and Jordan, M. I. (2003). Latent dirichlet allocation. *The Journal of Machine Learning Research*, 3(null):993–1022.
- Bonhomme, S. (2021). Teams: Heterogeneity, Sorting, and Complementarity. *SSRN Electronic Journal*.
- Bradbury, J., Frostig, R., Hawkins, P., Johnson, M. J., Leary, C., Maclaurin, D., Necula, G., Paszke, A., VanderPlas, J., Wanderman-Milne, S., and Zhang, Q. (2018). JAX: Composable transformations of Python+NumPy programs.
- Bybee, L., Kelly, B. T., Manela, A., and Xiu, D. (2020). The Structure of Economic News. Technical Report w26648, National Bureau of Economic Research.
- Chen, X., Christensen, T. M., and Tamer, E. (2018). Monte Carlo Confidence Sets for Identified Sets. *Econometrica*, 86(6):1965–2018.
- Chesher, A. (1991). The effect of measurement error. *Biometrika*, 78(3):451–462.
- Compiani, G., Morozov, I., and Seiler, S. (2023). Demand Estimation with Text and Image Data. Technical Report 10695, CESifo.
- Dieng, A. B., Ruiz, F. J. R., and Blei, D. M. (2020). Topic Modeling in Embedding Spaces. *Transactions of the Association for Computational Linguistics*, 8:439–453.
- Draca, M. and Schwarz, C. (2021). How Polarized are Citizens? Measuring Ideology from the Ground-Up. SSRN Scholarly Paper ID 3154431, Social Science Research Network, Rochester, NY.
- Egami, N., Hinck, M., Stewart, B., and Wei, H. (2023). Using Imperfect Surrogates for Downstream Inference: Design-based Supervised Learning for Social Science Applications of Large Language Models. *Advances in Neural Information Processing Systems*, 36:68589–68601.
- Einav, L., Finkelstein, A., and Mahoney, N. (2022). Producing Health: Measuring Value Added of Nursing Homes.
- Fong, C. and Tyler, M. (2021). Machine Learning Predictions as Regression Covariates. *Political Analysis*, 29(4):467–484.
- Gabaix, X., Koijen, R. S. J., and Yogo, M. (2023). Asset Embeddings. *SSRN Electronic Journal*.
- Gennaro, G. and Ash, E. (2022). Emotion and Reason in Political Language. *The Economic Journal*, 132(643):1037–1059.
- Gentzkow, M., Kelly, B., and Taddy, M. (2019a). Text as Data. *Journal of Economic Literature*, 57(3):535–574.

- Gentzkow, M., Shapiro, J. M., and Taddy, M. (2019b). Measuring Group Differences in High-Dimensional Choices: Method and Application to Congressional Speech. *Econometrica*, 87(4):1307–1340.
- Gorodnichenko, Y., Pham, T., and Talavera, O. (2023). The Voice of Monetary Policy. *American Economic Review*, 113(2):548–584.
- Griffiths, T. L. and Steyvers, M. (2004). Finding scientific topics. *Proceedings of the National Academy of Sciences*, 101(suppl 1):5228–5235.
- Hahn, J. and Kuersteiner, G. (2002). Asymptotically Unbiased Inference for a Dynamic Panel Model with Fixed Effects when Both n and T Are Large. *Econometrica*, 70(4):1639–1657.
- Hansen, S., McMahon, M., and Prat, A. (2018). Transparency and Deliberation Within the FOMC: A Computational Linguistics Approach*. *The Quarterly Journal of Economics*, 133(2):801–870.
- Herbst, E. P. and Schorfheide, F. (2016). *Bayesian Estimation of DSGE Models*. Princeton University Press, Princeton.
- Hoberg, G. and Phillips, G. (2016). Text-Based Network Industries and Endogenous Product Differentiation. *Journal of Political Economy*, 124(5):1423–1465.
- Hoffman, M. D., Blei, D. M., Wang, C., and Paisley, J. (2013). Stochastic Variational Inference. *Journal of Machine Learning Research*, 14(4):1303–1347.
- Hoffman, M. D. and Gelman, A. (2014). The No-U-Turn Sampler: Adaptively Setting Path Lengths in Hamiltonian Monte Carlo. *Journal of Machine Learning Research*, 15(47):1593–1623.
- Hofmann, T. (1999). Probabilistic latent semantic indexing. In *Proceedings of the 22nd Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 50–57, Berkeley California USA. ACM.
- Jordan, M. I., Ghahramani, Z., Jaakkola, T. S., and Saul, L. K. (1998). An Introduction to Variational Methods for Graphical Models. In Jordan, M. I., editor, *Learning in Graphical Models*, NATO ASI Series, pages 105–161. Springer Netherlands, Dordrecht.
- Ke, S., Olea, J. L. M., and Nesbit, J. (2021). Robust Machine Learning Algorithms for Text Analysis. Unpublished manuscript.
- Ke, Z. T. and Wang, M. (2022). Using SVD for Topic Modeling. *Journal of the American Statistical Association*, pages 1–16.
- Kelly, B., Papanikolaou, D., Seru, A., and Taddy, M. (2021). Measuring Technological Innovation over the Long Run. *American Economic Review: Insights*, 3(3):303–320.
- Kogan, L., Papanikolaou, D., Schmidt, L., and Seegmiller, B. (2019). Technology, Vintage-Specific Human Capital, and Labor Displacement: Evidence from Linking Patents with Occupations.

- Larsen, V. H. and Thorsrud, L. A. (2019). The value of news for economic developments. *Journal of Econometrics*, 210(1):203–218.
- MacKay, D. J. C. (2003). *Information Theory, Inference and Learning Algorithms*. Cambridge University Press, Cambridge, UK ; New York, illustrated edition edition.
- Magnolfi, L., McClure, J., and Sorensen, A. (2022). Embeddings and Distance-based Demand for Differentiated Products. In *Proceedings of the 23rd ACM Conference on Economics and Computation*, EC '22, page 607, New York, NY, USA. Association for Computing Machinery.
- Malmendier, U., Nagel, S., and Yan, Z. (2021). The making of hawks and doves. *Journal of Monetary Economics*, 117:19–42.
- Mardia, J., Jiao, J., Tánčzos, E., Nowak, R. D., and Weissman, T. (2019). Concentration Inequalities for the Empirical Distribution.
- Mele, A. and Zhu, L. (2023). Approximate Variational Estimation for a Model of Network Formation. *The Review of Economics and Statistics*, 105(1):113–124.
- Mikolov, T., Chen, K., Corrado, G., and Dean, J. (2013a). Efficient Estimation of Word Representations in Vector Space. *arXiv:1301.3781 [cs]*.
- Mikolov, T., Sutskever, I., Chen, K., Corrado, G., and Dean, J. (2013b). Distributed Representations of Words and Phrases and their Compositionality. *arXiv:1310.4546 [cs, stat]*.
- Mueller, H. and Rauh, C. (2018). Reading Between the Lines: Prediction of Political Violence Using Newspaper Text. *American Political Science Review*, 112(2):358–375.
- Munro, E. and Ng, S. (2022). Latent Dirichlet Analysis of Categorical Survey Responses. *Journal of Business & Economic Statistics*, 40(1):256–271.
- Neal, R. M. (2012). MCMC using Hamiltonian dynamics. *arXiv:1206.1901 [physics, stat]*.
- Nimczik, J. S. (2017). Job Mobility Networks and Endogenous Labor Markets. Technical Report 168147, Verein für Socialpolitik / German Economic Association.
- Olencki, A. and Sacher, S. (2022). Estimating Nursing Home Quality with Selection.
- Olivella, S., Pratt, T., and Imai, K. (2021). Dynamic Stochastic Blockmodel Regression for Network Data: Application to International Militarized Conflicts. *arXiv:2103.00702 [cs, stat]*.
- Pagan, A. (1984). Econometric Issues in the Analysis of Regressions with Generated Regressors. *International Economic Review*, 25(1):221–247.
- Phan, D., Pradhan, N., and Jankowiak, M. (2019). Composable Effects for Flexible and Accelerated Probabilistic Programming in NumPyro. *arXiv:1912.11554 [cs, stat]*.
- Phillips, P. C. B. (1987). Towards a unified asymptotic theory for autoregression. *Biometrika*, 74(3):535–547.

- Roberts, M. E., Stewart, B. M., Tingley, D., Lucas, C., Leder-Luis, J., Gadarian, S. K., Albertson, B., and Rand, D. G. (2014). Structural Topic Models for Open-Ended Survey Responses. *American Journal of Political Science*, 58(4):1064–1082.
- Ruiz, F. J. R., Athey, S., and Blei, D. M. (2020). SHOPPER: A probabilistic model of consumer choice with substitutes and complements. *The Annals of Applied Statistics*, 14(1):1–27.
- Staiger, D. and Stock, J. H. (1997). Instrumental Variables Regression with Weak Instruments. *Econometrica*, 65(3):557–586.
- Stock, J. H. and Watson, M. W. (2002). Forecasting Using Principal Components from a Large Number of Predictors. *Journal of the American Statistical Association*, 97(460):1167–1179.
- Thorsrud, L. A. (2020). Words are the New Numbers: A Newsy Coincident Index of the Business Cycle. *Journal of Business & Economic Statistics*, 38(2):393–409.
- Vafa, K., Athey, S., and Blei, D. M. (2023). Decomposing Changes in the Gender Wage Gap over Worker Careers. In *NBER Summer Institute*, Boston, MA.
- van der Vaart, A. W. (1998). *Asymptotic Statistics*. Cambridge University Press, Cambridge, UK.
- Wainwright, M. J. and Jordan, M. I. (2008). Graphical Models, Exponential Families, and Variational Inference. *Foundations and Trends® in Machine Learning*, 1(1–2):1–305.
- Wu, R., Zhang, L., and Tony Cai, T. (2023). Sparse Topic Modeling: Computational Efficiency, Near-Optimal Algorithms, and Statistical Inference. *Journal of the American Statistical Association*, 118(543):1849–1861.
- Zhang, J., Xue, W., Yu, Y., and Tan, Y. (2023). Debiasing Machine-Learning- or AI-Generated Regressors in Partial Linear Models. *SSRN Electronic Journal*.

A Proofs

Notation Let $\|\cdot\|$ denote the Euclidean norm when applied to vectors and the spectral norm when applied to matrices. Let $\|\cdot\|_F$ denote the Frobenius norm.

A.1 Proofs for Section 2

Proof of Proposition 1. We start by writing

$$\begin{aligned}\sqrt{n}(\hat{\gamma}_1 - \gamma_1) &= \frac{\frac{1}{\sqrt{n}} \sum_{i=1}^n (Y_i - \gamma_1(\hat{\theta}_i - \bar{\theta}))(\hat{\theta}_i - \bar{\theta})}{\frac{1}{n} \sum_{i=1}^n (\hat{\theta}_i - \bar{\theta})^2} \\ &= -\gamma_1 \frac{\frac{1}{\sqrt{n}} \sum_{i=1}^n (\hat{\theta}_i - \theta_i)(\hat{\theta}_i - \bar{\theta})}{\frac{1}{n} \sum_{i=1}^n (\hat{\theta}_i - \bar{\theta})^2} + \frac{\frac{1}{\sqrt{n}} \sum_{i=1}^n \varepsilon_i(\hat{\theta}_i - \bar{\theta})}{\frac{1}{n} \sum_{i=1}^n (\hat{\theta}_i - \bar{\theta})^2} =: T_{1,n} + T_{2,n},\end{aligned}$$

where $\bar{\theta} = \frac{1}{n} \sum_{i=1}^n \hat{\theta}_i$. Note by Chebyshev's inequality that for integers $k_1, k_2 \geq 0$ and any $t > 0$, we have

$$\Pr \left(\left| \frac{1}{n} \sum_{i=1}^n \hat{\theta}_i^{k_1} \theta_i^{k_2} - \mathbb{E}[\hat{\theta}_i^{k_1} \theta_i^{k_2}] \right| > t \right) \leq \frac{\mathbb{E}[\hat{\theta}_i^{2k_1} \theta_i^{2k_2}]}{t^2 n} \leq \frac{1}{t^2 n}. \quad (14)$$

Consider the denominator term in $T_{1,n}$ and $T_{2,n}$. By inequality (14), we have

$$\left| \frac{1}{n} \sum_{i=1}^n (\hat{\theta}_i - \bar{\theta})^2 - \text{Var}(\hat{\theta}_i) \right| \rightarrow_p 0,$$

where, by the law of total variance and independence of C_i and θ_i ,

$$\text{Var}(\hat{\theta}_i) = \text{Var}(\theta_i) + \mathbb{E} \left[\frac{1}{C_i} \right] \mathbb{E}[\theta_i(1 - \theta_i)] \rightarrow \text{Var}(\theta_i)$$

because $\mathbb{E}[C_i^{-1}] \rightarrow 0$.

For the numerator in $T_{1,n}$, we similarly have

$$\left| \frac{1}{\sqrt{n}} \sum_{i=1}^n (\hat{\theta}_i - \theta_i)(\hat{\theta}_i - \bar{\theta}) - \frac{1}{\sqrt{n}} \sum_{i=1}^n (\hat{\theta}_i - \theta_i)(\hat{\theta}_i - \mathbb{E}[\hat{\theta}_i]) \right| \rightarrow_p 0.$$

Because $\mathbb{E}[\hat{\theta}_i | \theta_i] = \theta_i$ and $\text{Var}(\hat{\theta}_i | \theta_i, C_i) = C_i^{-1} \theta_i(1 - \theta_i)$, we have

$$\begin{aligned}\mathbb{E} \left[\sqrt{n}(\hat{\theta}_i - \theta_i)(\hat{\theta}_i - \mathbb{E}[\hat{\theta}_i]) \right] &= \mathbb{E} \left[\sqrt{n}(\hat{\theta}_i - \theta_i)^2 \right] \\ &= \sqrt{n} \mathbb{E} \left[\frac{1}{C_i} \right] \mathbb{E}[\theta_i(1 - \theta_i)] \rightarrow \kappa \mathbb{E}[\theta_i(1 - \theta_i)].\end{aligned}$$

A second application of inequality (14) gives

$$\begin{aligned} \Pr \left(\left| \frac{1}{n} \sum_{i=1}^n \sqrt{n}(\hat{\theta}_i - \theta_i)(\hat{\theta}_i - \mathbb{E}[\hat{\theta}_i]) - \mathbb{E}[\sqrt{n}(\hat{\theta}_i - \theta_i)(\hat{\theta}_i - \mathbb{E}[\hat{\theta}_i])] \right| > t \right) \\ \leq \frac{\mathbb{E}[(\hat{\theta}_i - \theta_i)^2(\hat{\theta}_i - \mathbb{E}[\hat{\theta}_i])^2]}{t^2} \leq \frac{\mathbb{E}[(\hat{\theta}_i - \theta_i)^2]}{t^2} = \mathbb{E} \left[\frac{1}{C_i} \right] \frac{\mathbb{E}[\theta_i(1 - \theta_i)]}{t^2} \rightarrow 0. \end{aligned}$$

Hence,

$$T_{1,n} \rightarrow_p -\kappa \gamma_1 \frac{\mathbb{E}[\theta_i(1 - \theta_i)]}{\text{Var}(\theta_i)}.$$

For $T_{2,n}$, we have

$$\left| \frac{1}{\sqrt{n}} \sum_{i=1}^n \varepsilon_i(\hat{\theta}_i - \bar{\theta}) - \frac{1}{\sqrt{n}} \sum_{i=1}^n \varepsilon_i(\hat{\theta}_i - \theta_i) - \frac{1}{\sqrt{n}} \sum_{i=1}^n \varepsilon_i(\theta_i - \mathbb{E}[\theta_i]) \right| \rightarrow_p 0$$

because $(\bar{\theta} - \mathbb{E}[\theta_i]) \times \frac{1}{\sqrt{n}} \sum_{i=1}^n \varepsilon_i \rightarrow_p 0$. Note that $\mathbb{E}[\varepsilon_i(\hat{\theta}_i - \theta_i)] = 0$ because Y_i and (X_i, C_i) are independent conditional on θ_i and both ε_i and $\hat{\theta}_i - \theta_i$ have conditional (on θ_i) mean zero. Hence by Chebyshev's inequality, for any $t > 0$ we have

$$\Pr \left(\left| \frac{1}{n} \sum_{i=1}^n \sqrt{n} \varepsilon_i(\hat{\theta}_i - \theta_i) \right| > t \right) \leq \frac{\mathbb{E}[\varepsilon_i^2(\hat{\theta}_i - \theta_i)^2]}{t^2} = \mathbb{E} \left[\frac{1}{C_i} \right] \frac{\mathbb{E}[\varepsilon_i^2 \theta_i(1 - \theta_i)]}{t^2} \rightarrow 0,$$

because ε_i and (X_i, C_i) are independent conditional on θ_i , C_i and θ_i are independent, and $\mathbb{E}[C_i^{-1}] \rightarrow 0$. Finally, $\frac{1}{\sqrt{n}} \sum_{i=1}^n \varepsilon_i(\theta_i - \mathbb{E}[\theta_i])$ is asymptotically $N(0, \mathbb{E}[\varepsilon_i^2(\theta_i - \mathbb{E}[\theta_i])^2])$ by the central limit theorem. \blacksquare

The next assumption is used to derive Proposition 2.

Assumption 3. (i) $\text{Var}(Z_i) > 0$, $\mathbb{E}[Z_i^2] < \infty$, and $\mathbb{E}[u_i^2(Z_i - \mathbb{E}[Z_i])^2] < \infty$.

(ii) $\Pr(\theta_i \in [\delta, 1 - \delta]) = 1$ for some $\delta > 0$.

(iii) $C_i \gtrsim (\log n)^{1+\varepsilon}$ almost surely for some $\varepsilon > 0$.

Part (i) is standard and ensures the OLS estimator of ϕ_1 without measurement error is well defined with finite asymptotic variance. Part (ii) is made to simplify technical arguments and can be relaxed, e.g., by controlling the rate at which the distribution of θ_i behaves at the boundary of its support. Finally, part (iii) is the same as Assumption 2(vi) and is also made to simplify technical derivations and can be relaxed.

Proof of Proposition 2. To simplify notation, let $Y_i = \log\left(\frac{\theta_i}{1-\theta_i}\right)$ and $\hat{Y}_i = \log\left(\frac{\hat{\theta}_i}{1-\hat{\theta}_i}\right)$.

We have

$$\sqrt{n}(\hat{\phi}_1 - \phi_1) = \frac{\frac{1}{\sqrt{n}} \sum_{i=1}^n u_i(Z_i - \bar{Z})}{\frac{1}{n} \sum_{i=1}^n (Z_i - \bar{Z})^2} + \frac{\frac{1}{\sqrt{n}} \sum_{i=1}^n (\hat{Y}_i - Y_i)(Z_i - \bar{Z})}{\frac{1}{n} \sum_{i=1}^n (Z_i - \bar{Z})^2} =: T_{1,n} + T_{2,n},$$

where $\bar{Z} = \frac{1}{n} \sum_{i=1}^n Z_i$. It follows by standard arguments that under Assumption 3(i), we have

$$T_{1,n} \rightarrow_d N\left(0, \frac{\mathbb{E}[u_i^2(Z_i - \mathbb{E}[Z_i])^2]}{\text{Var}(Z_i)^2}\right)$$

and $\frac{1}{n} \sum_{i=1}^n (Z_i - \bar{Z})^2 \rightarrow_p \text{Var}(Z_i)$.

It remains to characterize the numerator of $T_{2,n}$. To this end, first note that with δ as in Assumption 3(ii), we have

$$\begin{aligned} \Pr\left(\min_{1 \leq i \leq n} \hat{\theta}_i < \delta/2 \mid \{(C_i, \theta_i)\}_{i=1}^n\right) &\leq \sum_{i=1}^n \Pr\left(\hat{\theta}_i < \delta/2 \mid C_i, \theta_i\right) \\ &\leq \sum_{i=1}^n \Pr\left(\hat{\theta}_i < \theta_i/2 \mid C_i, \theta_i\right) \quad (\text{almost surely}) \\ &\leq \sum_{i=1}^n e^{-\frac{1}{8}C_i\theta_i} \quad (\text{almost surely}) \\ &\leq ne^{-\frac{1}{8}\delta c(\log n)^{1+\epsilon}} \quad (\text{almost surely}), \end{aligned}$$

where the first inequality is by the union bound, the second is by Assumption 3(ii), the third is by Chernoff's inequality for Binomial random variables, and the fourth is because $C_i \geq c(\log n)^{1+\epsilon}$ for some $c > 0$ and $\theta_i \geq \delta$ both hold for all i with probability one by Assumption 3(ii)-(iii). Therefore,

$$\Pr\left(\min_{1 \leq i \leq n} \hat{\theta}_i < \delta/2\right) \leq ne^{-\frac{1}{8}\delta c(\log n)^{1+\epsilon}} \rightarrow 0. \quad (15)$$

We may similarly deduce that

$$\Pr\left(\max_{1 \leq i \leq n} \hat{\theta}_i > 1 - \delta/2\right) \rightarrow 0, \quad (16)$$

and that

$$\max_{1 \leq i \leq n} |\hat{\theta}_i - \theta_i| \rightarrow_p 0. \quad (17)$$

In view of Assumption 3(ii), condition (17) also implies that $\max_{1 \leq i \leq n} |\hat{Y}_i - Y_i| \rightarrow_p 0$ because $x \mapsto \log(\frac{x}{1-x})$ is uniformly continuous on $[\delta, 1 - \delta]$. But then note that this

implies

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n (\hat{Y}_i - Y_i)(\mathbb{E}[Z_i] - \bar{Z}) \leq \max_{1 \leq i \leq n} |\hat{Y}_i - Y_i| |\sqrt{n}(\mathbb{E}[Z_i] - \bar{Z})| \rightarrow_p 0$$

by Assumption 3(i). It therefore remains to show

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n (\hat{Y}_i - Y_i) (Z_i - \mathbb{E}[Z_i]) \rightarrow_p \kappa \text{Cov} \left(\frac{2\theta_i - 1}{2\theta_i(1 - \theta_i)}, Z_i \right).$$

By Taylor's theorem, we have

$$\hat{Y}_i - Y_i = \frac{\hat{\theta}_i - \theta_i}{\theta_i(1 - \theta_i)} + \frac{(2\theta_i - 1)(\hat{\theta}_i - \theta_i)^2}{2\theta_i^2(1 - \theta_i)^2} + \frac{(3\tilde{\theta}_i^2 - 3\tilde{\theta}_i + 1)(\hat{\theta}_i - \theta_i)^3}{3\tilde{\theta}_i^3(1 - \tilde{\theta}_i)^3}$$

where $\tilde{\theta}_i$ is between θ_i and $\hat{\theta}_i$. Note that Assumption 3(ii) implies that $\theta_i(1 - \theta_i) \geq \delta^2$. We also have by (15) and (16) that $\tilde{\theta}_i(1 - \tilde{\theta}_i) \geq \delta^2/4$ with probability approaching one (wpa1). Thus, all terms on the right-hand side are well defined wpa1. We control the covariance of Z_i with these terms using $\mathbb{E}[\hat{\theta}_i|\theta_i] = \theta_i$ and $\text{Var}(\hat{\theta}_i|\theta_i, C_i) = C_i^{-1}\theta_i(1 - \theta_i)$ and the fact that (X_i, C_i) and Z_i are independent conditional on θ_i as follows:

First, by Chebyshev's inequality, we have for $t > 0$ that

$$\begin{aligned} \Pr \left(\left| \frac{1}{\sqrt{n}} \sum_{i=1}^n \frac{\hat{\theta}_i - \theta_i}{\theta_i(1 - \theta_i)} (Z_i - \mathbb{E}[Z_i]) \right| > t \right) &\leq \frac{1}{t^2} \mathbb{E} \left[\frac{(\hat{\theta}_i - \theta_i)^2}{\theta_i^2(1 - \theta_i)^2} (Z_i - \mathbb{E}[Z_i])^2 \right] \\ &\leq \frac{1}{t^2} \mathbb{E} \left[\frac{1}{C_i} \right] \mathbb{E} \left[\frac{(Z_i - \mathbb{E}[Z_i])^2}{\theta_i(1 - \theta_i)} \right] \rightarrow 0 \end{aligned}$$

by (3), independence of (X_i, C_i) and Z_i conditional on θ_i , and independence of C_i and (θ_i, Z_i) . Second, we similarly have

$$\begin{aligned} \sqrt{n} \mathbb{E} \left[\frac{(2\theta_i - 1)(\hat{\theta}_i - \theta_i)^2}{2\theta_i^2(1 - \theta_i)^2} (Z_i - \mathbb{E}[Z_i]) \right] &= \sqrt{n} \mathbb{E} \left[\frac{1}{C_i} \right] \mathbb{E} \left[\frac{(2\theta_i - 1)}{2\theta_i(1 - \theta_i)} (Z_i - \mathbb{E}[Z_i]) \right] \\ &\rightarrow \kappa \text{Cov} \left(\frac{2\theta_i - 1}{2\theta_i(1 - \theta_i)}, Z_i \right). \end{aligned}$$

Moreover, letting $W_i = \frac{(2\theta_i - 1)(\hat{\theta}_i - \theta_i)^2}{2\theta_i^2(1 - \theta_i)^2} (Z_i - \mathbb{E}[Z_i])$ and noting $|2\theta_i - 1| \leq 1$ and $|\hat{\theta}_i - \theta_i| \leq 1$ because $\theta_i, \hat{\theta}_i \in [0, 1]$, we have by Chebyshev's inequality that for $t > 0$,

$$\begin{aligned} \Pr \left(\left| \frac{1}{\sqrt{n}} \sum_{i=1}^n W_i - \mathbb{E}[W_i] \right| > t \right) &\leq \frac{1}{t^2} \mathbb{E}[W_i^2] \leq \frac{1}{4t^2} \mathbb{E} \left[\frac{(\hat{\theta}_i - \theta_i)^2}{\theta_i^4(1 - \theta_i)^4} (Z_i - \mathbb{E}[Z_i])^2 \right] \\ &\leq \frac{1}{4t^2} \mathbb{E} \left[\frac{1}{C_i} \right] \mathbb{E} \left[\frac{(Z_i - \mathbb{E}[Z_i])^2}{\theta_i^3(1 - \theta_i)^3} \right] \rightarrow 0. \end{aligned}$$

Finally, because $\tilde{\theta}_i \in [\delta/2, 1 - \delta/2]$ holds for all $1 \leq i \leq n$ wpa1, there is a positive constant D such that

$$\begin{aligned} \left| \frac{1}{\sqrt{n}} \sum_{i=1}^n \frac{(3\tilde{\theta}_i^2 - 3\tilde{\theta}_i + 1)(\hat{\theta}_i - \theta_i)^3}{3\tilde{\theta}_i^3(1 - \tilde{\theta}_i)^3} (Z_i - \mathbb{E}[Z_i]) \right| \\ \leq D \max_{1 \leq i \leq n} |\hat{\theta}_i - \theta_i| \left(\frac{1}{\sqrt{n}} \sum_{i=1}^n (\hat{\theta}_i - \theta_i)^2 |Z_i - \mathbb{E}[Z_i]| \right) \end{aligned}$$

holds wpa1. Hence, in view of (17), it suffices to show that the right-hand side term is bounded in probability. To this end, note by Markov's inequality that for $t > 0$,

$$\begin{aligned} \Pr \left(\frac{1}{\sqrt{n}} \sum_{i=1}^n (\hat{\theta}_i - \theta_i)^2 |Z_i - \mathbb{E}[Z_i]| > t \right) &\leq \frac{1}{t} \sqrt{n} \mathbb{E} \left[(\hat{\theta}_i - \theta_i)^2 |Z_i - \mathbb{E}[Z_i]| \right] \\ &= \frac{1}{t} \sqrt{n} \mathbb{E} \left[\frac{1}{C_i} \right] \mathbb{E} [\theta_i(1 - \theta_i) |Z_i - \mathbb{E}[Z_i]|] \\ &\rightarrow \frac{1}{t} \kappa \mathbb{E} [\theta_i(1 - \theta_i) |Z_i - \mathbb{E}[Z_i]|], \end{aligned}$$

as required. ■

A.2 Proofs for Section 3

The next two lemmas apply in both fixed-populations and sequences-of-populations.

Lemma 1. *Suppose that (5) holds. Then*

$$\mathbb{E} [\hat{\mathbf{p}}_i \hat{\mathbf{p}}_i^T | C_i, \boldsymbol{\theta}_i] = \mathbf{B}^T \boldsymbol{\theta}_i \boldsymbol{\theta}_i^T \mathbf{B} + \frac{1}{C_i} (\text{diag}(\mathbf{B}^T \boldsymbol{\theta}_i) - \mathbf{B}^T \boldsymbol{\theta}_i \boldsymbol{\theta}_i^T \mathbf{B}),$$

and

$$\mathbb{E} [(\hat{\mathbf{p}}_i - \mathbf{p}_i)(\hat{\mathbf{p}}_i - \mathbf{p}_i)^T | C_i, \boldsymbol{\theta}_i] = \frac{1}{C_i} (\text{diag}(\mathbf{B}^T \boldsymbol{\theta}_i) - \mathbf{B}^T \boldsymbol{\theta}_i \boldsymbol{\theta}_i^T \mathbf{B}).$$

Proof of Lemma 1. First note by (5) that

$$\begin{aligned} \mathbb{E} [\hat{\mathbf{p}}_i \hat{\mathbf{p}}_i^T | C_i, \boldsymbol{\theta}_i] &= \frac{1}{C_i^2} \mathbb{E} [\mathbf{x}_i \mathbf{x}_i^T | C_i, \boldsymbol{\theta}_i] \\ &= \frac{1}{C_i^2} \left(\mathbb{E} [\mathbf{x}_i | C_i, \boldsymbol{\theta}_i] \mathbb{E} [\mathbf{x}_i | C_i, \boldsymbol{\theta}_i]^T + \text{Var} [\mathbf{x}_i | C_i, \boldsymbol{\theta}_i] \right) \\ &= \mathbf{B}^T \boldsymbol{\theta}_i \boldsymbol{\theta}_i^T \mathbf{B} + \frac{1}{C_i} (\text{diag}(\mathbf{B}^T \boldsymbol{\theta}_i) - \mathbf{B}^T \boldsymbol{\theta}_i \boldsymbol{\theta}_i^T \mathbf{B}), \end{aligned}$$

where the last line follows from the mean and variance of the multinomial distribution.

The second result now follows because $\mathbb{E} [\hat{\mathbf{p}}_i | C_i, \boldsymbol{\theta}_i] = \mathbf{p}_i = \mathbf{B}^T \boldsymbol{\theta}_i$. ■

Lemma 2. *Let Assumption 1(i)-(iii) hold. Then*

$$\left\| \frac{1}{n} \sum_{i=1}^n \hat{\boldsymbol{\theta}}_i \hat{\boldsymbol{\theta}}_i^T - \mathbb{E} [\boldsymbol{\theta}_i \boldsymbol{\theta}_i^T] - \mathbb{E} \left[\frac{1}{C_i} \right] \left((\mathbf{B}\mathbf{B}^T)^{-1} \mathbf{B} \text{diag}(\mathbf{B}^T \mathbb{E}[\boldsymbol{\theta}_i]) \mathbf{B}^T (\mathbf{B}\mathbf{B}^T)^{-1} - \mathbb{E} [\boldsymbol{\theta}_i \boldsymbol{\theta}_i^T] \right) \right\| \rightarrow_p 0.$$

Proof of Lemma 2. In view of Assumption 1(iii), we have

$$\left\| \frac{1}{n} \sum_{i=1}^n \hat{\boldsymbol{\theta}}_i \hat{\boldsymbol{\theta}}_i^T - (\hat{\mathbf{B}}\hat{\mathbf{B}}^T)^{-1} \hat{\mathbf{B}} \left(\frac{1}{n} \sum_{i=1}^n \hat{\mathbf{p}}_i \hat{\mathbf{p}}_i^T \right) \hat{\mathbf{B}}^T (\hat{\mathbf{B}}\hat{\mathbf{B}}^T)^{-1} \right\| \rightarrow_p 0$$

where $(\hat{\mathbf{B}}\hat{\mathbf{B}}^T)^{-1}$ exists with probability approaching one by Assumption 1(i)-(ii). Each element of $\hat{\mathbf{p}}_i \hat{\mathbf{p}}_i^T$ is bounded between 0 and 1, so we may deduce by Chebyshev's inequality that

$$\left\| \frac{1}{n} \sum_{i=1}^n \hat{\mathbf{p}}_i \hat{\mathbf{p}}_i^T - \mathbb{E} [\hat{\mathbf{p}}_i \hat{\mathbf{p}}_i^T] \right\| \rightarrow_p 0.$$

Hence, by Assumption 1(ii) and Slutsky's theorem, we have

$$\left\| \frac{1}{n} \sum_{i=1}^n \hat{\boldsymbol{\theta}}_i \hat{\boldsymbol{\theta}}_i^T - (\mathbf{B}\mathbf{B}^T)^{-1} \mathbf{B} \mathbb{E} [\hat{\mathbf{p}}_i \hat{\mathbf{p}}_i^T] \mathbf{B}^T (\mathbf{B}\mathbf{B}^T)^{-1} \right\| \rightarrow_p 0.$$

The result follows by Lemma 1 and independence of C_i and $\boldsymbol{\theta}_i$. ■

Lemma 3. *Let Assumption 1(iv) hold. Then*

$$\frac{1}{n} \sum_{i=1}^n \mathbf{q}_i \mathbf{q}_i^T \rightarrow_p \mathbb{E} [\mathbf{q}_i \mathbf{q}_i^T], \quad \left(\frac{1}{n} \sum_{i=1}^n \mathbf{q}_i \mathbf{q}_i^T \right)^{-1} \rightarrow_p \mathbb{E} [\mathbf{q}_i \mathbf{q}_i^T]^{-1}.$$

Proof of Lemma 3. The first result follows by the law of large numbers (LLN). The second result then follows because the rank condition on $\mathbb{E} [\boldsymbol{\xi}_i \boldsymbol{\xi}_i^T]$ implies $\mathbb{E} [\mathbf{q}_i \mathbf{q}_i^T]$ has full rank. ■

Lemma 4. *Let Assumption 1 hold. Then*

$$\frac{1}{n} \sum_{i=1}^n \hat{\boldsymbol{\theta}}_i \mathbf{q}_i^T \rightarrow_p \mathbb{E} [\boldsymbol{\theta}_i \mathbf{q}_i^T].$$

Proof of Lemma 4. In view of Assumption 1(iii)-(iv), we have

$$\begin{aligned} & \left\| \frac{1}{n} \sum_{i=1}^n \hat{\boldsymbol{\theta}}_i \mathbf{q}_i^T - (\hat{\mathbf{B}}\hat{\mathbf{B}}^T)^{-1} \hat{\mathbf{B}} \left(\frac{1}{n} \sum_{i=1}^n \hat{\mathbf{p}}_i \mathbf{q}_i^T \right) \right\| \\ & \leq \left(\max_{1 \leq i \leq n} \|\hat{\boldsymbol{\theta}}_i - (\hat{\mathbf{B}}\hat{\mathbf{B}}^T)^{-1} \hat{\mathbf{B}} \hat{\mathbf{p}}_i\| \right) \times \frac{1}{n} \sum_{i=1}^n \|\mathbf{q}_i\| \rightarrow_p 0. \end{aligned}$$

Moreover, $(\hat{\mathbf{B}}\hat{\mathbf{B}}^T)^{-1}\hat{\mathbf{B}} \rightarrow_p (\mathbf{B}\mathbf{B}^T)^{-1}\mathbf{B}$ by Assumption 1(i)-(ii). Let $\mathbf{X}_i = \mathbf{q}_i\hat{\mathbf{p}}_i^T - \mathbb{E}[\mathbf{q}_i\hat{\mathbf{p}}_i^T]$ and let D denote the dimension of \mathbf{q}_i . Then

$$\begin{aligned} \mathbb{E} \left[\left\| \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i \right\|_F^2 \right] &= \frac{1}{n} \sum_{j=1}^D \sum_{k=1}^V \mathbb{E} \left[(\mathbf{X}_i)_{j,k}^2 \right] \leq \frac{1}{n} \sum_{j=1}^D \sum_{k=1}^V \mathbb{E} \left[(\mathbf{q}_{i,j})^2 (\hat{\mathbf{p}}_{i,k})^2 \right] \\ &\leq \frac{1}{n} \mathbb{E} [\|\mathbf{q}_i\|^2] \rightarrow 0, \end{aligned}$$

by Assumption 1(iv) and the fact that $\hat{\mathbf{p}}_i$ is in the simplex. Hence, $\frac{1}{n} \sum_{i=1}^n \mathbf{X}_i \rightarrow_p 0$. The result follows by Slutsky's theorem, noting that $\mathbb{E}[\mathbf{q}_i\hat{\mathbf{p}}_i^T] = \mathbb{E}[\mathbf{q}_i\boldsymbol{\theta}_i^T]\mathbf{B}$. \blacksquare

In what follows, we let $\mathbf{0}$ denote a conformable matrix of zeros.

Lemma 5. *Let Assumption 1 hold. Then*

$$\left\| \frac{1}{n} \sum_{i=1}^n \hat{\boldsymbol{\xi}}_i \hat{\boldsymbol{\xi}}_i^T - \mathbb{E} [\boldsymbol{\xi}_i \boldsymbol{\xi}_i^T] - \begin{bmatrix} \mathbb{E} \left[\frac{1}{C_i} \right] \mathbf{S} \left((\mathbf{B}\mathbf{B}^T)^{-1} \mathbf{B} \text{diag}(\mathbf{B}^T \mathbb{E}[\boldsymbol{\theta}_i]) \mathbf{B}^T (\mathbf{B}\mathbf{B}^T)^{-1} - \mathbb{E} [\boldsymbol{\theta}_i \boldsymbol{\theta}_i^T] \right) \mathbf{S}^T & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \right\| \rightarrow_p 0.$$

Proof of Lemma 5. Note that

$$\frac{1}{n} \sum_{i=1}^n \hat{\boldsymbol{\xi}}_i \hat{\boldsymbol{\xi}}_i^T = \begin{bmatrix} \mathbf{S} \left(\frac{1}{n} \sum_{i=1}^n \hat{\boldsymbol{\theta}}_i \hat{\boldsymbol{\theta}}_i^T \right) \mathbf{S}^T & \mathbf{S} \left(\frac{1}{n} \sum_{i=1}^n \hat{\boldsymbol{\theta}}_i \mathbf{q}_i^T \right) \\ \left(\frac{1}{n} \sum_{i=1}^n \mathbf{q}_i \hat{\boldsymbol{\theta}}_i^T \right) \mathbf{S}^T & \frac{1}{n} \sum_{i=1}^n \mathbf{q}_i \mathbf{q}_i^T \end{bmatrix}.$$

The result follows by Lemmas 2, 3, and 4. \blacksquare

Proof of Theorem 1. First consider the denominator. By Lemma 5, we have

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \hat{\boldsymbol{\xi}}_i \hat{\boldsymbol{\xi}}_i^T &\rightarrow_p \mathbb{E} [\boldsymbol{\xi}_i \boldsymbol{\xi}_i^T] \\ &+ \begin{bmatrix} \mathbb{E} \left[\frac{1}{C_i} \right] \mathbf{S} \left((\mathbf{B}\mathbf{B}^T)^{-1} \mathbf{B} \text{diag}(\mathbf{B}^T \mathbb{E}[\boldsymbol{\theta}_i]) \mathbf{B}^T (\mathbf{B}\mathbf{B}^T)^{-1} - \mathbb{E} [\boldsymbol{\theta}_i \boldsymbol{\theta}_i^T] \right) \mathbf{S}^T & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}. \end{aligned} \quad (18)$$

For the numerator term, first note that

$$\frac{1}{n} \sum_{i=1}^n \hat{\boldsymbol{\xi}}_i Y_i = \begin{bmatrix} \mathbf{S} \left(\frac{1}{n} \sum_{i=1}^n \hat{\boldsymbol{\theta}}_i Y_i \right) \\ \frac{1}{n} \sum_{i=1}^n \mathbf{q}_i Y_i \end{bmatrix}. \quad (19)$$

For the upper block on the right-hand side, we use Assumption 1(iii) to deduce

$$\left\| \frac{1}{n} \sum_{i=1}^n \hat{\boldsymbol{\theta}}_i Y_i - (\hat{\mathbf{B}}\hat{\mathbf{B}}^T)^{-1} \hat{\mathbf{B}} \left(\frac{1}{n} \sum_{i=1}^n \hat{\mathbf{p}}_i Y_i \right) \right\| \rightarrow_p 0.$$

We have $(\hat{\mathbf{B}}\hat{\mathbf{B}}^T)^{-1} \hat{\mathbf{B}} \rightarrow_p (\mathbf{B}\mathbf{B}^T)^{-1} \mathbf{B}$ by Assumption 1(i)-(ii). Moreover, by the LLN and independence of \mathbf{x}_i and Y_i conditional on $(C_i, \boldsymbol{\theta}_i)$, we have

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \hat{\mathbf{p}}_i Y_i &\rightarrow_p \mathbb{E} [\hat{\mathbf{p}}_i Y_i] = \mathbb{E} \left[\mathbb{E} \left[\frac{\mathbf{x}_i}{C_i} \middle| C_i, \boldsymbol{\theta}_i \right] \mathbb{E} [Y_i | C_i, \boldsymbol{\theta}_i] \right] \\ &= \mathbb{E} [\mathbf{B}^T \boldsymbol{\theta}_i \mathbb{E} [Y_i | C_i, \boldsymbol{\theta}_i]] \\ &= \mathbf{B}^T \mathbb{E} [\boldsymbol{\theta}_i \mathbb{E} [Y_i | \boldsymbol{\theta}_i, \mathbf{q}_i]] \\ &= \mathbf{B}^T \mathbb{E} [\boldsymbol{\theta}_i \boldsymbol{\theta}_i^T] \mathbf{S}^T \boldsymbol{\gamma} + \mathbf{B}^T \mathbb{E} [\boldsymbol{\theta}_i \mathbf{q}_i^T] \boldsymbol{\alpha}, \end{aligned}$$

Hence,

$$\mathbf{S} \left(\frac{1}{n} \sum_{i=1}^n \hat{\boldsymbol{\theta}}_i Y_i \right) \rightarrow_p \mathbf{S} \mathbb{E} [\boldsymbol{\theta}_i \boldsymbol{\theta}_i^T] \mathbf{S}^T \boldsymbol{\gamma} + \mathbf{S} \mathbb{E} [\boldsymbol{\theta}_i \mathbf{q}_i^T] \boldsymbol{\alpha}.$$

Similarly, for the lower block on the right-hand side of (19), we have by the LLN that

$$\frac{1}{n} \sum_{i=1}^n \mathbf{q}_i Y_i \rightarrow_p \mathbb{E} [\mathbf{q}_i Y_i] = \mathbb{E} [\mathbf{q}_i \boldsymbol{\theta}_i^T] \mathbf{S}^T \boldsymbol{\gamma} + \mathbb{E} [\mathbf{q}_i \mathbf{q}_i^T] \boldsymbol{\alpha}.$$

Combining the above two displays, we have

$$\frac{1}{n} \sum_{i=1}^n \hat{\boldsymbol{\xi}}_i Y_i \rightarrow_p \mathbb{E} [\boldsymbol{\xi}_i \boldsymbol{\xi}_i^T] \boldsymbol{\psi}. \quad (20)$$

The first result follows from (18) and (20). Note that the matrix on the right-hand side of (18) is bounded below (in Loewner order) by $\mathbb{E} [\boldsymbol{\xi}_i \boldsymbol{\xi}_i^T]$ so its inverse is well defined by Assumption 1(iv). The second result follows because $(\mathbf{A} + \boldsymbol{\Delta})^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1} \boldsymbol{\Delta} \mathbf{A}^{-1} + O(\|\boldsymbol{\Delta}\|^2)$ for \mathbf{A} invertible and $\boldsymbol{\Delta}$ small. \blacksquare

The next three lemmas are used in the proof of Theorem 2. They are derived in the sequence-of-populations asymptotic framework.

Lemma 6. *Let Assumption 2(i)-(iii) hold. Then*

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n \hat{\boldsymbol{\theta}}_i (\boldsymbol{\theta}_i - \hat{\boldsymbol{\theta}}_i)^T \rightarrow_p -\kappa \left((\mathbf{B}\mathbf{B}^T)^{-1} \mathbf{B} \text{diag}(\mathbf{B}^T \mathbb{E}[\boldsymbol{\theta}_i]) \mathbf{B}^T (\mathbf{B}\mathbf{B}^T)^{-1} - \mathbb{E}[\boldsymbol{\theta}_i \boldsymbol{\theta}_i^T] \right).$$

Proof of Lemma 6. First note that $\left\| \frac{1}{\sqrt{n}} \sum_{i=1}^n \hat{\boldsymbol{\theta}}_i (\boldsymbol{\theta}_i - \hat{\boldsymbol{\theta}}_i)^T - T_{1,n} - T_{2,n} \right\| \rightarrow_p 0$ by As-

sumption 2(iii), where

$$T_{1,n} = (\hat{\mathbf{B}}\hat{\mathbf{B}}^T)^{-1}\hat{\mathbf{B}} \left(\left(\frac{1}{n} \sum_{i=1}^n \hat{\mathbf{p}}_i \mathbf{p}_i^T \right) \sqrt{n} \left(\mathbf{B}^T (\mathbf{B}\mathbf{B}^T)^{-1} - \hat{\mathbf{B}}^T (\hat{\mathbf{B}}\hat{\mathbf{B}}^T)^{-1} \right) \right)$$

$$T_{2,n} = (\hat{\mathbf{B}}\hat{\mathbf{B}}^T)^{-1}\hat{\mathbf{B}} \left(\frac{1}{\sqrt{n}} \sum_{i=1}^n \hat{\mathbf{p}}_i (\mathbf{p}_i - \hat{\mathbf{p}}_i)^T \right) \hat{\mathbf{B}}^T (\hat{\mathbf{B}}\hat{\mathbf{B}}^T)^{-1}.$$

Assumption 2(i)-(ii) implies $\sqrt{n}(\mathbf{B}^T(\mathbf{B}\mathbf{B}^T)^{-1} - \hat{\mathbf{B}}^T(\hat{\mathbf{B}}\hat{\mathbf{B}}^T)^{-1}) \rightarrow_p \mathbf{0}$. As $\|\frac{1}{n} \sum_{i=1}^n \hat{\mathbf{p}}_i \mathbf{p}_i^T\| \leq 1$, it follows that $T_{1,n} \rightarrow_p \mathbf{0}$. For term $T_{2,n}$, note by Lemma 1 that

$$\begin{aligned} \mathbb{E} \left[\hat{\mathbf{p}}_i (\hat{\mathbf{p}}_i - \mathbf{p}_i)^T \right] &= \mathbb{E} \left[(\hat{\mathbf{p}}_i - \mathbf{p}_i) (\hat{\mathbf{p}}_i - \mathbf{p}_i)^T \right] \\ &= \mathbb{E} \left[\frac{1}{C_i} \right] \left(\text{diag}(\mathbf{B}^T \mathbb{E}[\boldsymbol{\theta}_i]) - \mathbf{B}^T \mathbb{E}[\boldsymbol{\theta}_i \boldsymbol{\theta}_i^T] \mathbf{B} \right). \end{aligned} \quad (21)$$

Let $\mathbf{X}_i = \hat{\mathbf{p}}_i (\hat{\mathbf{p}}_i - \mathbf{p}_i)^T - \mathbb{E} \left[\hat{\mathbf{p}}_i (\hat{\mathbf{p}}_i - \mathbf{p}_i)^T \right]$. Then

$$\begin{aligned} \mathbb{E} \left[\left\| \frac{1}{\sqrt{n}} \sum_{i=1}^n \mathbf{X}_i \right\|_F^2 \right] &= \sum_{j=1}^V \sum_{k=1}^V \mathbb{E} \left[(\mathbf{X}_i)_{j,k}^2 \right] \leq \sum_{j=1}^V \sum_{k=1}^V \mathbb{E} \left[(\hat{\mathbf{p}}_{i,j})^2 (\hat{\mathbf{p}}_{i,k} - \mathbf{p}_{i,k})^2 \right] \\ &\leq \sum_{k=1}^V \mathbb{E} \left[(\hat{\mathbf{p}}_{i,k} - \mathbf{p}_{i,k})^2 \right] \rightarrow 0, \end{aligned}$$

where the second inequality is because $\hat{\mathbf{p}}_i$ is in the simplex and the convergence to zero holds in view of (7) and (21). It follows that

$$\left\| \frac{1}{\sqrt{n}} \sum_{i=1}^n \hat{\mathbf{p}}_i (\mathbf{p}_i - \hat{\mathbf{p}}_i)^T - \sqrt{n} \mathbb{E} \left[\hat{\mathbf{p}}_i (\mathbf{p}_i - \hat{\mathbf{p}}_i)^T \right] \right\| \rightarrow_p 0.$$

We conclude that $T_{2,n} \rightarrow_p -\kappa \left((\mathbf{B}\mathbf{B}^T)^{-1} \mathbf{B} \text{diag}(\mathbf{B}^T \mathbb{E}[\boldsymbol{\theta}_i]) \mathbf{B}^T (\mathbf{B}\mathbf{B}^T)^{-1} - \mathbb{E}[\boldsymbol{\theta}_i \boldsymbol{\theta}_i^T] \right)$ by (7), (21), and Assumption 2(i)-(ii) \blacksquare

Lemma 7. *Let Assumption 2(i)-(iv) hold. Then*

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n \mathbf{q}_i (\boldsymbol{\theta}_i - \hat{\boldsymbol{\theta}}_i)^T \rightarrow_p \mathbf{0}.$$

Proof of Lemma 7. First note that $\left\| \frac{1}{\sqrt{n}} \sum_{i=1}^n \mathbf{q}_i (\boldsymbol{\theta}_i - \hat{\boldsymbol{\theta}}_i)^T - T_{1,n} - T_{2,n} \right\| \rightarrow_p 0$ by As-

sumption 2(iii)-(iv), where

$$T_{1,n} = \left(\left(\frac{1}{n} \sum_{i=1}^n \mathbf{q}_i \mathbf{p}_i^T \right) \sqrt{n} \left(\mathbf{B}^T (\mathbf{B} \mathbf{B}^T)^{-1} - \hat{\mathbf{B}}^T (\hat{\mathbf{B}} \hat{\mathbf{B}}^T)^{-1} \right) \right)$$

$$T_{2,n} = \left(\frac{1}{\sqrt{n}} \sum_{i=1}^n \mathbf{q}_i (\mathbf{p}_i - \hat{\mathbf{p}}_i)^T \right) \hat{\mathbf{B}}^T (\hat{\mathbf{B}} \hat{\mathbf{B}}^T)^{-1}.$$

Assumption 2(i)-(ii) implies that $\sqrt{n}(\mathbf{B}^T(\mathbf{B}\mathbf{B}^T)^{-1} - \hat{\mathbf{B}}^T(\hat{\mathbf{B}}\hat{\mathbf{B}}^T)^{-1}) \rightarrow_p \mathbf{0}$. Moreover, $\|\frac{1}{n} \sum_{i=1}^n \mathbf{q}_i \mathbf{p}_i^T\| \leq \frac{1}{n} \sum_{i=1}^n \|\mathbf{q}_i\|$, which is bounded in probability by Assumption 2(iv). It follows that $T_{1,n} \rightarrow_p \mathbf{0}$. For $T_{2,n}$, note that $\mathbb{E}[\mathbf{q}_i(\hat{\mathbf{p}}_i - \mathbf{p}_i)^T] = \mathbf{0}$ by independence of \mathbf{q}_i and \mathbf{x}_i conditional on $(C_i, \boldsymbol{\theta}_i)$. Let $\mathbf{X}_i = \mathbf{q}_i(\hat{\mathbf{p}}_i - \mathbf{p}_i)^T$ and D the dimension of \mathbf{q}_i . Then

$$\begin{aligned} \mathbb{E} \left[\left\| \frac{1}{\sqrt{n}} \sum_{i=1}^n \mathbf{X}_i \right\|_F^2 \right] &= \sum_{j=1}^D \sum_{k=1}^V \mathbb{E} \left[(\mathbf{X}_i)_{j,k}^2 \right] = \sum_{j=1}^D \sum_{k=1}^V \mathbb{E} \left[(\mathbf{q}_{i,j})^2 (\hat{\mathbf{p}}_{i,k} - \mathbf{p}_{i,k})^2 \right] \\ &\leq \sum_{j=1}^D \sum_{k=1}^V \mathbb{E} \left[(\mathbf{q}_{i,j})^4 \right]^{1/2} \mathbb{E} \left[(\hat{\mathbf{p}}_{i,k} - \mathbf{p}_{i,k})^4 \right]^{1/2} \\ &\leq \text{constant} \times \sum_{k=1}^V \mathbb{E} \left[(\hat{\mathbf{p}}_{i,k} - \mathbf{p}_{i,k})^2 \right]^{1/2} \rightarrow 0, \end{aligned}$$

where the first inequality is by Cauchy-Schwarz, the second is by Assumption 2(iv) and the fact that $|\hat{\mathbf{p}}_{i,k} - \mathbf{p}_{i,k}| \leq 1$, and convergence to zero is by (7) and (21). It follows that $\frac{1}{\sqrt{n}} \sum_{i=1}^n \mathbf{X}_i \rightarrow_p \mathbf{0}$. We conclude by Assumption 2(i)-(ii) that $T_{2,n} \rightarrow_p \mathbf{0}$. ■

Lemma 8. *Let Assumption 2(vi) hold. Then*

$$\max_{1 \leq i \leq n} \|\hat{\mathbf{p}}_i - \mathbf{p}_i\| \rightarrow_p 0.$$

Proof of Lemma 8. Let $\|\cdot\|_1$ be the ℓ^1 norm. As $\hat{\mathbf{p}}_i | (C_i, \boldsymbol{\theta}_i) \sim C_i^{-1} \text{Multinomial}(C_i, \mathbf{p}_i)$, for all $t > 0$ we have

$$\Pr \left(\max_{1 \leq i \leq n} \|\hat{\mathbf{p}}_i - \mathbf{p}_i\|_1 > t \mid \{(C_i, \boldsymbol{\theta}_i)\}_{i=1}^n \right) \leq \sum_{i=1}^n (2^V - 2) \exp \left\{ -\frac{C_i t^2}{2K} \right\}$$

by the union bound and Lemma 1 of Mardia et al. (2019). Then by Assumption 2(vi),

$$\Pr \left(\max_{1 \leq i \leq n} \|\hat{\mathbf{p}}_i - \mathbf{p}_i\|_1 > t \right) \leq n(2^V - 2) \exp \left\{ -\frac{c(\log n)^{1+\epsilon} t^2}{2K} \right\},$$

where $c, \epsilon > 0$. Hence, $\max_{1 \leq i \leq n} \|\hat{\mathbf{p}}_i - \mathbf{p}_i\|_1 \rightarrow_p 0$. The result now follows because the ℓ^1 norm is weakly greater than the Euclidean norm. ■

Proof of Theorem 2. First consider the denominator term. By Lemma 5 and condition (7), we have $\frac{1}{n} \sum_{i=1}^n \hat{\boldsymbol{\xi}}_i \hat{\boldsymbol{\xi}}_i^T \rightarrow_p \mathbb{E}[\boldsymbol{\xi}_i \boldsymbol{\xi}_i^T]$. Hence, by Assumption 2(iv),

$$\left(\frac{1}{n} \sum_{i=1}^n \hat{\boldsymbol{\xi}}_i \hat{\boldsymbol{\xi}}_i^T \right)^{-1} \rightarrow_p \mathbb{E}[\boldsymbol{\xi}_i \boldsymbol{\xi}_i^T]^{-1}. \quad (22)$$

Now consider the numerator term. We have

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n \hat{\boldsymbol{\xi}}_i (Y_i - \hat{\boldsymbol{\xi}}_i^T \boldsymbol{\psi}) = \frac{1}{\sqrt{n}} \sum_{i=1}^n \hat{\boldsymbol{\xi}}_i (\boldsymbol{\xi}_i - \hat{\boldsymbol{\xi}}_i)^T \boldsymbol{\psi} + \frac{1}{\sqrt{n}} \sum_{i=1}^n \hat{\boldsymbol{\xi}}_i \varepsilon_i =: T_{1,n} + T_{2,n}.$$

First consider term $T_{1,n}$. By definition,

$$T_{1,n} = \begin{bmatrix} \mathbf{S} \left(\frac{1}{\sqrt{n}} \sum_{i=1}^n \hat{\boldsymbol{\theta}}_i (\boldsymbol{\theta}_i - \hat{\boldsymbol{\theta}}_i)^T \right) \mathbf{S}^T \boldsymbol{\gamma} \\ \left(\frac{1}{\sqrt{n}} \sum_{i=1}^n \mathbf{q}_i (\boldsymbol{\theta}_i - \hat{\boldsymbol{\theta}}_i)^T \right) \mathbf{S}^T \boldsymbol{\gamma} \end{bmatrix}.$$

It follows by Lemmas 6 and 7 that

$$T_{1,n} \rightarrow_p \begin{bmatrix} -\kappa \mathbf{S} \left((\mathbf{B}\mathbf{B}^T)^{-1} \mathbf{B} \text{diag}(\mathbf{B}^T \mathbb{E}[\boldsymbol{\theta}_i]) \mathbf{B}^T (\mathbf{B}\mathbf{B}^T)^{-1} - \mathbb{E}[\boldsymbol{\theta}_i \boldsymbol{\theta}_i^T] \right) \mathbf{S}^T \boldsymbol{\gamma} \\ \mathbf{0} \end{bmatrix}. \quad (23)$$

Now consider term $T_{2,n}$. By construction,

$$T_{2,n} = \begin{bmatrix} \mathbf{S} \left(\frac{1}{\sqrt{n}} \sum_{i=1}^n \hat{\boldsymbol{\theta}}_i \varepsilon_i \right) \\ \frac{1}{\sqrt{n}} \sum_{i=1}^n \mathbf{q}_i \varepsilon_i \end{bmatrix}.$$

We may deduce by arguments similar to those in the proof of Lemma 6 that

$$\left\| \frac{1}{\sqrt{n}} \sum_{i=1}^n \hat{\boldsymbol{\theta}}_i \varepsilon_i - (\mathbf{B}\mathbf{B}^T)^{-1} \mathbf{B} \left(\frac{1}{\sqrt{n}} \sum_{i=1}^n \hat{\mathbf{p}}_i \varepsilon_i \right) \right\| \rightarrow_p 0.$$

under Assumption 2(i)-(iii). Moreover,

$$\begin{aligned} \mathbb{E} [\varepsilon_i^2 \|\hat{\mathbf{p}}_i - \mathbf{p}_i\|^2] &= \mathbb{E} [\mathbb{E} [\varepsilon_i^2 | C_i, \boldsymbol{\theta}_i] \mathbb{E} [\|\hat{\mathbf{p}}_i - \mathbf{p}_i\|^2 | C_i, \boldsymbol{\theta}_i]] \\ &= \mathbb{E} \left[\mathbb{E} [\varepsilon_i^2 | C_i, \boldsymbol{\theta}_i] \frac{1}{C_i} \text{tr} \{ \text{diag}(\mathbf{B}^T \boldsymbol{\theta}_i) - \mathbf{B}^T \boldsymbol{\theta}_i \boldsymbol{\theta}_i^T \mathbf{B} \} \right] \\ &= \mathbb{E} \left[\frac{1}{C_i} \right] \mathbb{E} [\varepsilon_i^2 (\text{diag}(\mathbf{B}^T \boldsymbol{\theta}_i) - \mathbf{B}^T \boldsymbol{\theta}_i \boldsymbol{\theta}_i^T \mathbf{B})] \rightarrow 0. \end{aligned}$$

In the above display, the first equality is by independence of (Y_i, \mathbf{q}_i) and \mathbf{x}_i conditional on $(C_i, \boldsymbol{\theta}_i)$, the second is by Lemma 1, and the third is by independence of C_i and $(Y_i, \boldsymbol{\theta}_i, \mathbf{q}_i)$.

Therefore,

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n (\hat{\mathbf{p}}_i - \mathbf{p}_i) \varepsilon_i \rightarrow_p 0$$

and so

$$\left\| T_{2,n} - \frac{1}{\sqrt{n}} \sum_{i=1}^n \boldsymbol{\xi}_i \varepsilon_i \right\| \rightarrow_p 0.$$

Note that $\mathbb{E} [\varepsilon_i^2 \boldsymbol{\xi}_i \boldsymbol{\xi}_i^T]$ is finite by Assumption 2(iv)-(v). It follows by the central limit theorem that

$$T_{2,n} \rightarrow_d N(\mathbf{0}, \mathbb{E} [\varepsilon_i^2 \boldsymbol{\xi}_i \boldsymbol{\xi}_i^T]). \quad (24)$$

Result (8) now follows by combining (22), (23), and (24).

For result (9), it remains to show

$$\frac{1}{n} \sum_{i=1}^n \hat{\varepsilon}_i^2 \hat{\boldsymbol{\xi}}_i \hat{\boldsymbol{\xi}}_i^T \rightarrow_p \mathbb{E} [\varepsilon_i^2 \boldsymbol{\xi}_i \boldsymbol{\xi}_i^T].$$

To this end, first write

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \hat{\varepsilon}_i^2 \hat{\boldsymbol{\xi}}_i \hat{\boldsymbol{\xi}}_i^T &= \frac{1}{n} \sum_{i=1}^n \varepsilon_i^2 \boldsymbol{\xi}_i \boldsymbol{\xi}_i^T + \frac{1}{n} \sum_{i=1}^n \varepsilon_i^2 (\hat{\boldsymbol{\xi}}_i \hat{\boldsymbol{\xi}}_i^T - \boldsymbol{\xi}_i \boldsymbol{\xi}_i^T) \\ &\quad + \frac{1}{n} \sum_{i=1}^n (\hat{\varepsilon}_i^2 - \varepsilon_i^2) \hat{\boldsymbol{\xi}}_i \hat{\boldsymbol{\xi}}_i^T =: T_{3,n} + T_{4,n} + T_{5,n}. \end{aligned}$$

Evidently, $T_{3,n} \rightarrow_p \mathbb{E} [\varepsilon_i^2 \boldsymbol{\xi}_i \boldsymbol{\xi}_i^T]$ by the LLN. Now consider $T_{4,n}$. By construction, we have

$$T_{4,n} = \begin{bmatrix} \frac{1}{n} \sum_{i=1}^n \varepsilon_i^2 (\hat{\boldsymbol{\theta}}_i \hat{\boldsymbol{\theta}}_i^T - \boldsymbol{\theta}_i \boldsymbol{\theta}_i^T) & \frac{1}{n} \sum_{i=1}^n \varepsilon_i^2 (\hat{\boldsymbol{\theta}}_i - \boldsymbol{\theta}_i) \mathbf{q}_i^T \\ \frac{1}{n} \sum_{i=1}^n \varepsilon_i^2 \mathbf{q}_i (\hat{\boldsymbol{\theta}}_i - \boldsymbol{\theta}_i)^T & \mathbf{0} \end{bmatrix}.$$

Consider the upper-left block. We may deduce by arguments similar to those in the proof of Lemma 6 that

$$\left\| \frac{1}{n} \sum_{i=1}^n \varepsilon_i^2 (\hat{\boldsymbol{\theta}}_i \hat{\boldsymbol{\theta}}_i^T - \boldsymbol{\theta}_i \boldsymbol{\theta}_i^T) - (\mathbf{B}\mathbf{B}^T)^{-1} \mathbf{B} \left(\frac{1}{n} \sum_{i=1}^n \varepsilon_i^2 (\hat{\mathbf{p}}_i \hat{\mathbf{p}}_i^T - \mathbf{p}_i \mathbf{p}_i^T) \right) \mathbf{B}^T (\mathbf{B}\mathbf{B}^T)^{-1} \right\| \rightarrow_p 0,$$

by Assumption 2(i)-(iii),(v). Since \mathbf{p}_i and $\hat{\mathbf{p}}_i$ both take values in the simplex, we have $\|\hat{\mathbf{p}}_i \hat{\mathbf{p}}_i^T - \mathbf{p}_i \mathbf{p}_i^T\| \leq 2\|\hat{\mathbf{p}}_i - \mathbf{p}_i\|$ and so

$$\left\| \frac{1}{n} \sum_{i=1}^n \varepsilon_i^2 (\hat{\mathbf{p}}_i \hat{\mathbf{p}}_i^T - \mathbf{p}_i \mathbf{p}_i^T) \right\| \leq 2 \left(\max_{1 \leq i \leq n} \|\hat{\mathbf{p}}_i - \mathbf{p}_i\| \right) \frac{1}{n} \sum_{i=1}^n \varepsilon_i^2 \rightarrow_p 0,$$

by Lemma 8 and Assumption 2(v). Now consider the lower-left (equivalently, upper-right)

block. Again by arguments similar to those in the proof of Lemma 7, we have

$$\left\| \frac{1}{n} \sum_{i=1}^n \varepsilon_i^2 \mathbf{q}_i (\hat{\boldsymbol{\theta}}_i - \boldsymbol{\theta}_i)^T - \left(\frac{1}{n} \sum_{i=1}^n \varepsilon_i^2 \mathbf{q}_i (\hat{\mathbf{p}}_i - \mathbf{p}_i)^T \right) \mathbf{B}^T (\mathbf{B} \mathbf{B}^T)^{-1} \right\| \rightarrow_p 0,$$

by Assumption 2(i)-(v). But note that

$$\left\| \frac{1}{n} \sum_{i=1}^n \varepsilon_i^2 \mathbf{q}_i (\hat{\mathbf{p}}_i - \mathbf{p}_i)^T \right\| \leq \left(\max_{1 \leq i \leq n} \|\hat{\mathbf{p}}_i - \mathbf{p}_i\| \right) \frac{1}{n} \sum_{i=1}^n \varepsilon_i^2 \|\mathbf{q}_i\| \rightarrow_p 0,$$

by Lemma 8 and Assumption 2(iv)-(v). Therefore, $T_{4,n} \rightarrow_p \mathbf{0}$.

Now consider $T_{5,n}$. We have

$$\hat{\varepsilon}_i - \varepsilon_i = (\mathbf{S} \hat{\boldsymbol{\theta}}_i)^T (\boldsymbol{\gamma} - \hat{\boldsymbol{\gamma}}) + (\mathbf{S} (\boldsymbol{\theta}_i - \hat{\boldsymbol{\theta}}_i))^T \boldsymbol{\gamma} + \mathbf{q}_i^T (\boldsymbol{\alpha} - \hat{\boldsymbol{\alpha}}),$$

where

$$\max_{1 \leq i \leq n} |(\mathbf{S} \hat{\boldsymbol{\theta}}_i)^T (\boldsymbol{\gamma} - \hat{\boldsymbol{\gamma}})| \leq \left(\max_{1 \leq i \leq n} \|\mathbf{S} (\hat{\boldsymbol{\theta}}_i - (\hat{\mathbf{B}} \hat{\mathbf{B}}^T)^{-1} \hat{\mathbf{B}} \hat{\mathbf{p}}_i)\| + \|\mathbf{S} (\hat{\mathbf{B}} \hat{\mathbf{B}}^T)^{-1} \hat{\mathbf{B}}\| \right) \|\hat{\boldsymbol{\gamma}} - \boldsymbol{\gamma}\| \rightarrow_p 0$$

by Assumption 2(iii), consistency of $\hat{\boldsymbol{\gamma}}$, and the fact that $\|(\hat{\mathbf{B}} \hat{\mathbf{B}}^T)^{-1} \hat{\mathbf{B}}\|$ is bounded in probability by Assumption 2(i)-(ii) and $\|\hat{\mathbf{p}}_i\| \leq 1$. Moreover,

$$\begin{aligned} \max_{1 \leq i \leq n} |(\mathbf{S} (\boldsymbol{\theta}_i - \hat{\boldsymbol{\theta}}_i))^T \boldsymbol{\gamma}| &\leq \|\mathbf{S}\| \left(\max_{1 \leq i \leq n} \|\boldsymbol{\theta}_i - (\hat{\mathbf{B}} \hat{\mathbf{B}}^T)^{-1} \hat{\mathbf{B}} \hat{\mathbf{p}}_i\| \right. \\ &\quad \left. + \left\| (\hat{\mathbf{B}} \hat{\mathbf{B}}^T)^{-1} \hat{\mathbf{B}} - (\mathbf{B} \mathbf{B}^T)^{-1} \mathbf{B} \right\| + \|(\mathbf{B} \mathbf{B}^T)^{-1} \mathbf{B}\| \max_{1 \leq i \leq n} \|\hat{\mathbf{p}}_i - \mathbf{p}_i\| \right) \|\boldsymbol{\gamma}\|. \end{aligned}$$

Consider the three terms in parentheses on the right-hand side of this display. The first two terms converge in probability to zero by Assumption 2(i)-(iii) and the third converges in probability to zero by Lemma 8. Finally, we have

$$\max_{1 \leq i \leq n} \|\mathbf{q}_i^T (\boldsymbol{\alpha} - \hat{\boldsymbol{\alpha}})\| \leq \left(\max_{1 \leq i \leq n} \|\mathbf{q}_i\| \right) \|\boldsymbol{\alpha} - \hat{\boldsymbol{\alpha}}\| \rightarrow_p 0$$

by \sqrt{n} -consistency of $\hat{\boldsymbol{\alpha}}$ and the fact that $n^{-1/4} \max_{1 \leq i \leq n} \|\mathbf{q}_i\| \rightarrow_p 0$ by Assumption 2(iv). Therefore, $\max_{1 \leq i \leq n} |\hat{\varepsilon}_i - \varepsilon_i| \rightarrow_p 0$.

Now, since

$$\hat{\varepsilon}_i^2 - \varepsilon_i^2 = 2(\hat{\varepsilon}_i - \varepsilon_i)\varepsilon_i + (\hat{\varepsilon}_i - \varepsilon_i)^2,$$

we have

$$T_{5,n} = \frac{2}{n} \sum_{i=1}^n (\hat{\varepsilon}_i - \varepsilon_i) \varepsilon_i \hat{\boldsymbol{\xi}}_i \hat{\boldsymbol{\xi}}_i^T + \frac{1}{n} \sum_{i=1}^n (\hat{\varepsilon}_i - \varepsilon_i)^2 \hat{\boldsymbol{\xi}}_i \hat{\boldsymbol{\xi}}_i^T,$$

and so

$$\begin{aligned} \|T_{5,n}\| &\leq 2 \left(\max_{1 \leq i \leq n} |\hat{\varepsilon}_i - \varepsilon_i| \right) \frac{1}{n} \sum_{i=1}^n |\varepsilon_i| \|\hat{\boldsymbol{\xi}}_i\|^2 + \left(\max_{1 \leq i \leq n} |\hat{\varepsilon}_i - \varepsilon_i|^2 \right) \frac{1}{n} \sum_{i=1}^n \|\hat{\boldsymbol{\xi}}_i\|^2 \\ &= \left(\max_{1 \leq i \leq n} |\hat{\varepsilon}_i - \varepsilon_i| \right) \text{tr} \left\{ \frac{2}{n} \sum_{i=1}^n |\varepsilon_i| \hat{\boldsymbol{\xi}}_i \hat{\boldsymbol{\xi}}_i^T \right\} + \left(\max_{1 \leq i \leq n} |\hat{\varepsilon}_i - \varepsilon_i|^2 \right) \text{tr} \left\{ \frac{1}{n} \sum_{i=1}^n \hat{\boldsymbol{\xi}}_i \hat{\boldsymbol{\xi}}_i^T \right\} \rightarrow_p 0, \end{aligned}$$

because $\frac{1}{n} \sum_{i=1}^n |\varepsilon_i| \hat{\boldsymbol{\xi}}_i \hat{\boldsymbol{\xi}}_i^T$ is bounded in probability by control of $T_{3,n}$ and $T_{4,n}$, which together imply $\frac{1}{n} \sum_{i=1}^n \varepsilon_i^2 \hat{\boldsymbol{\xi}}_i \hat{\boldsymbol{\xi}}_i^T$ is bounded in probability, and $\frac{1}{n} \sum_{i=1}^n \hat{\boldsymbol{\xi}}_i \hat{\boldsymbol{\xi}}_i^T$ is bounded in probability by Lemma 5. \blacksquare

Proof of Theorem 3. Let $S_i = \mathbf{p}_{1,i} \cdot \mathbf{p}_{2,i}$ and $\hat{S}_i = \hat{\mathbf{p}}_{1,i} \cdot \hat{\mathbf{p}}_{2,i}$. By standard OLS algebra,

$$\sqrt{n}(\hat{\gamma}_1 - \gamma_1) = \frac{1}{\frac{1}{n} \sum_{i=1}^n (\hat{S}_i - \bar{\hat{S}})^2} \left(\frac{1}{\sqrt{n}} \sum_{i=1}^n \varepsilon_i (\hat{S}_i - \bar{\hat{S}}) - \left(\frac{1}{\sqrt{n}} \sum_{i=1}^n (\hat{S}_i - S_i) (\hat{S}_i - \bar{\hat{S}}) \right) \gamma_1 \right),$$

where $\bar{\hat{S}} = \frac{1}{n} \sum_{i=1}^n \hat{S}_i$.

By Chebyshev's inequality, for all integers $k_1, k_2 \geq 0$ and all $t > 0$, we have

$$\Pr \left(\left| \frac{1}{n} \sum_{i=1}^n \hat{S}_i^{k_1} S_i^{k_2} - \mathbb{E} \left[\hat{S}_i^{k_1} S_i^{k_2} \right] \right| > t \right) \leq \frac{1}{nt^2} \mathbb{E} \left[\hat{S}_i^{2k_1} S_i^{2k_2} \right] \leq \frac{1}{nt^2}, \quad (25)$$

because $|\hat{S}_i| \leq \|\hat{\mathbf{p}}_{1,i}\| \|\hat{\mathbf{p}}_{2,i}\| \leq 1$ by virtue of the fact that $\|\hat{\mathbf{p}}_{t,i}\| \leq \|\hat{\mathbf{p}}_{t,i}\|_1 = 1$ for $t = 1, 2$, with $\|\cdot\|_1$ denoting the ℓ^1 norm, and similarly for S_i . Let $F_i = (\mathbf{p}_{1,i}, \mathbf{p}_{2,i}, C_{1,i}, C_{2,i})$. Note that

$$\mathbb{E} \left[\hat{S}_i \mid F_i \right] = S_i$$

because $\mathbf{x}_{1,i}$, and $\mathbf{x}_{2,i}$ are independent conditional on $(C_{1,i}, C_{2,i}, \mathbf{p}_{1,i}, \mathbf{p}_{2,i})$. Hence, $\mathbb{E}[\hat{S}_i] = \mathbb{E}[S_i]$ and so it follows by (25) that $\bar{\hat{S}} \rightarrow_p \mathbb{E}[S_i]$. Moreover, for any conformable non-stochastic matrix \mathbf{M} , we have for $t = 1, 2$ that

$$\mathbb{E} \left[\hat{\mathbf{p}}_{t,i}^T \mathbf{M} \hat{\mathbf{p}}_{t,i} \mid F_i \right] = \mathbf{p}_{t,i}^T \mathbf{M} \mathbf{p}_{t,i} + \frac{1}{C_{t,i}} \text{tr} \left\{ \mathbf{M} (\text{diag}(\mathbf{p}_{t,i}) - \mathbf{p}_{t,i} \mathbf{p}_{t,i}^T) \right\}.$$

Hence, by independence of $\mathbf{x}_{1,i}$ and $\mathbf{x}_{2,i}$ conditional on F_i , we have

$$\begin{aligned}
\mathbb{E} \left[\hat{S}_i^2 \middle| F_i \right] &= \mathbb{E} \left[\mathbb{E} \left[\hat{\mathbf{p}}_{1,i}^T (\hat{\mathbf{p}}_{2,i} \hat{\mathbf{p}}_{2,i}^T) \hat{\mathbf{p}}_{1,i} \middle| \hat{\mathbf{p}}_{2,i}, F_i \right] \middle| F_i \right] \\
&= \mathbb{E} \left[\mathbf{p}_{1,i}^T (\hat{\mathbf{p}}_{2,i} \hat{\mathbf{p}}_{2,i}^T) \mathbf{p}_{1,i} + \frac{1}{C_{1,i}} \operatorname{tr} \left\{ (\hat{\mathbf{p}}_{2,i} \hat{\mathbf{p}}_{2,i}^T) (\operatorname{diag}(\mathbf{p}_{1,i}) - \mathbf{p}_{1,i} \mathbf{p}_{1,i}^T) \right\} \middle| F_i \right] \\
&= \mathbb{E} \left[\hat{\mathbf{p}}_{2,i}^T (\mathbf{p}_{1,i} \mathbf{p}_{1,i}^T) \hat{\mathbf{p}}_{2,i} \middle| F_i \right] + \frac{1}{C_{1,i}} \mathbb{E} \left[\hat{\mathbf{p}}_{2,i}^T (\operatorname{diag}(\mathbf{p}_{1,i}) - \mathbf{p}_{1,i} \mathbf{p}_{1,i}^T) \hat{\mathbf{p}}_{2,i} \middle| F_i \right] \\
&= S_i^2 + \frac{1}{C_{2,i}} \mathbf{p}_{1,i}^T (\operatorname{diag}(\mathbf{p}_{2,i}) - \mathbf{p}_{2,i} \mathbf{p}_{2,i}^T) \mathbf{p}_{1,i} + \frac{1}{C_{1,i}} \mathbf{p}_{2,i}^T (\operatorname{diag}(\mathbf{p}_{1,i}) - \mathbf{p}_{1,i} \mathbf{p}_{1,i}^T) \mathbf{p}_{2,i} \\
&\quad + \frac{1}{C_{1,i} C_{2,i}} \operatorname{tr} \left\{ (\operatorname{diag}(\mathbf{p}_{2,i}) - \mathbf{p}_{2,i} \mathbf{p}_{2,i}^T) (\operatorname{diag}(\mathbf{p}_{1,i}) - \mathbf{p}_{1,i} \mathbf{p}_{1,i}^T) \right\}. \tag{26}
\end{aligned}$$

It follows by (11) and (25) that

$$\left| \frac{1}{n} \sum_{i=1}^n \hat{S}_i^2 - \mathbb{E}[S_i^2] \right| \rightarrow_p 0.$$

Hence, $\frac{1}{n} \sum_{i=1}^n (\hat{S}_i - \bar{S})^2 \rightarrow_p \operatorname{Var}(S_i)$.

For the numerator term, note that

$$\left| \frac{1}{\sqrt{n}} \sum_{i=1}^n \varepsilon_i (\hat{S}_i - \bar{S}) - \frac{1}{\sqrt{n}} \sum_{i=1}^n \varepsilon_i (S_i - \mathbb{E}[S_i]) \right| \rightarrow_p 0,$$

because: firstly,

$$\Pr \left(\left| \frac{1}{\sqrt{n}} \sum_{i=1}^n \varepsilon_i (\hat{S}_i - S_i) \right| > t \right) \leq \frac{1}{t^2} \mathbb{E} \left[\mathbb{E}[\varepsilon_i^2 | F_i] \mathbb{E}[\hat{S}_i^2 - S_i^2 | F_i] \right] \rightarrow 0$$

by (11) and (26), mutual independence of $C_{1,i}$, $C_{2,i}$, and $(\mathbf{p}_{1,i}, \mathbf{p}_{2,i}, Y_i)$, and using finite second moment of ε_i ; and, second,

$$\left| \frac{1}{\sqrt{n}} \sum_{i=1}^n \varepsilon_i (\bar{S} - \mathbb{E}[S_i]) \right| \leq \left| \bar{S} - \mathbb{E}[S_i] \right| \left| \frac{1}{\sqrt{n}} \sum_{i=1}^n \varepsilon_i \right| \rightarrow_p 0$$

by the CLT and consistency of \bar{S} . So by a second application of the CLT we deduce that

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n \varepsilon_i (\hat{S}_i - \bar{S}) \rightarrow_d N(0, \mathbb{E}[\varepsilon_i^2 (S_i - \mathbb{E}[S_i])^2]).$$

Finally to characterize the bias term, first note that

$$\left| \frac{1}{\sqrt{n}} \sum_{i=1}^n (\hat{S}_i - S_i)(\hat{S}_i - \bar{S}) - \frac{1}{\sqrt{n}} \sum_{i=1}^n (\hat{S}_i - S_i)\hat{S}_i \right| = |\bar{S}| \left| \frac{1}{\sqrt{n}} \sum_{i=1}^n (\hat{S}_i - S_i) \right| \rightarrow_p 0$$

by consistency of \bar{S} and because

$$\mathbb{E}[(\hat{S}_i - S_i)^2] \rightarrow 0$$

holds by (11) and (26) and mutual independence of $C_{1,i}$, $C_{2,i}$, and $(\mathbf{p}_{1,i}, \mathbf{p}_{2,i}, Y_i)$. Hence by Chebyshev's inequality, we have

$$\begin{aligned} \Pr \left(\left| \frac{1}{\sqrt{n}} \sum_{i=1}^n \left((\hat{S}_i - S_i)\hat{S}_i - \mathbb{E} \left[(\hat{S}_i - S_i)\hat{S}_i \right] \right) \right| > t \right) &\leq \frac{1}{t^2} \mathbb{E}[(\hat{S}_i - S_i)^2 \hat{S}_i^2] \\ &\leq \frac{1}{t^2} \mathbb{E}[(\hat{S}_i - S_i)^2] \rightarrow 0, \end{aligned}$$

where the second inequality is because $|\hat{S}_i| \leq 1$. We have therefore shown that

$$\left| \frac{1}{\sqrt{n}} \sum_{i=1}^n (\hat{S}_i - S_i)(\hat{S}_i - \bar{S}) - \mathbb{E}[\sqrt{n}(\hat{S}_i - S_i)\hat{S}_i] \right| \rightarrow_p 0.$$

Finally,

$$\begin{aligned} \mathbb{E}[\sqrt{n}(\hat{S}_i - S_i)\hat{S}_i] &\rightarrow \kappa_1 \mathbb{E} \left[\mathbf{p}_{2,i}^T (\text{diag}(\mathbf{p}_{1,i}) - \mathbf{p}_{1,i}\mathbf{p}_{1,i}^T) \mathbf{p}_{2,i} \right] \\ &\quad + \kappa_2 \mathbb{E} \left[\mathbf{p}_{1,i}^T (\text{diag}(\mathbf{p}_{2,i}) - \mathbf{p}_{2,i}\mathbf{p}_{2,i}^T) \mathbf{p}_{1,i} \right], \end{aligned}$$

by (11), (26), and mutual independence of $C_{1,i}$, $C_{2,i}$, and $(\mathbf{p}_{1,i}, \mathbf{p}_{2,i}, Y_i)$. ■

B Further Details on the Simulation Exercise

Table B.1 presents the parameters used for simulation exercise. Since we used $K = 2$ types and type shares must add to 1, only the differences in regression parameters, e.g. $\gamma_1 - \gamma_2$ are identified, therefore in the simulation and estimation we normalized γ_2 and ϕ_2 to 0. Second, since class ‘labels’ are not identified in estimation, it is necessary to adjust signs post estimation.

To investigate the impact of κ on the estimation of γ , we ran three sets of simulations which vary only by the total number of features drawn per observations. For simplicity, in each of the set of simulations we set C_i to be equal for all i . We set $C \in \{10, 25, 200\}$ which, given that N is fixed to 10000 corresponds to $\kappa \in \{10, 4, 0.5\}$.

To ensure the model is properly identified, in each simulation we set $A = 100$ features to be ‘anchor words’ meaning that $\beta_{j,0}$ or $\beta_{j,1}$ is set to 0.

We simulated data 200 times for each set and then estimated the model using 1-step approach, 2-step approach and the infeasible 2-step approach with known θ .

We construct 95% confidence intervals for γ_1 and ϕ_1 using the corresponding 95% posterior credible intervals for these parameters. This construction is justified in view of the discussion in Section 4.1.

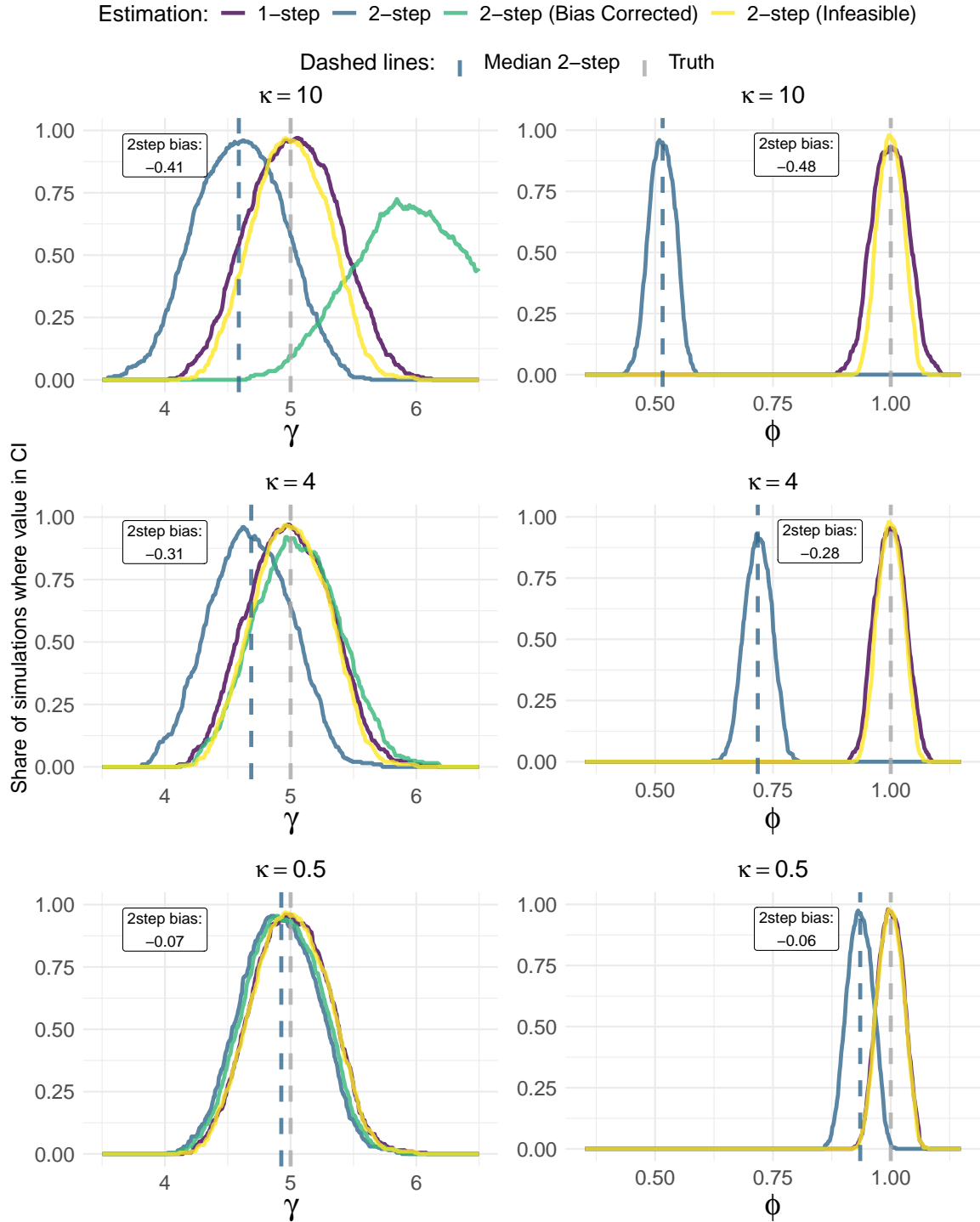
Table B.1: Parameters for the simulation exercise.

Parameter	Value	Description
(a) Data Simulation		
N	10000	Number of observations
V	300	Number of distinct features
C_i	$\{10, 25, 200\}$	Total number of features per document
K	2	Number of latent types
True ϕ	1	Effect of a covariates on un-normalized type shares
True γ	5	Effect of topic shares on numerical outcomes
True α	$(0, 1, 1, 1)$	Effect of additional covariates on numerical outcomes
g_i	$\sim N(0, \frac{\log(3)}{1.96})$	Covariate affecting type shares
$q_{i,m} \forall m \in (1, 2, 3)$	$\sim N(0, 3)$	Additional covariates affecting outcome
σ_Y^2	16	SD of the numeric outcome’s residual
σ_θ^2	1	SD of residual of the un-normalized type shares
η	0.2	Dirichlet concentration parameter
(b) Hyperparameters		
K	<i>as above</i>	Number of latent types
η	<i>as above</i>	Dirichlet concentration parameter
σ_θ^2	<i>as above</i>	SD of residual of the un-normalized type shares
$p(\phi_1)$	$N(0, 4)$	Prior for ϕ_1 , i.e. $\sigma_\phi^2 = 4$
$p(\gamma_1)$	$N(0, 100)$	Prior for γ_1 , i.e. $\sigma_\gamma^2 = 100$
$p(\alpha) \forall m \in (0, 1, 2, 3)$	$N(0, 100)$	Prior for α , i.e. $\sigma_\alpha^2 = 100$
$p(\sigma_Y)$	Gamma(1, 10)	Prior for σ_Y , i.e. $s_0 = 1$ and $s_1 = 10$

We performed the simulation on a ‘N1-highmem-2’ instance on the Google Cloud Platform. The instance has 2 vCPUs and 13 GB of memory. We also utilized a single Tesla V100 GPU. We run chose 500 warmup and 500 post-warmup iterations. A single simulation (consisting of drawing the data, and estimating the model in three ways) took approximately 6 minutes.

B.1 Additional Results

Figure B.1 below, replicates Figure 1 but adds the bias-corrected estimates for γ_1 . The bias-corrected estimates are obtained by subtracting the estimated bias computed using the formula from Theorem 1 from the two-step estimates. Since the estimates of the asymptotic bias are negative, the bias-corrected estimates are larger than the two-step estimates. Since the theorem relies on a first-order approximation, the bias-correction works best for small values of κ . As κ increases, the accuracy of the bias-corrected estimates diminishes, and for $\kappa = 10$, they perform worse than the uncorrected two-step estimates.



Note: Each mountain plot presents the share of simulations in which the value of γ_1 (respectively ϕ_1) on the x -axis is included in the 95% confidence interval. The grey vertical dashed lines show the true value of the parameter. The blue vertical dashed line represents the median (across simulations) of mean posterior estimates from the two-step strategy. The bias reported is the difference between the truth and this median value.

Figure B.1: Evolution of Bias in Regression Coefficients across κ Values, including Bias-Corrected Estimates

C Example Code

```
1
2 from numpyro import sample, plate
3 import numpyro.distributions as dist
4 import jax.numpy as jnp
5 from jax.nn import softmax
6
7 class SUPPMC:
8     def __init__(self, K, N, V, z, q, eta = .1, alpha = 1):
9         self.K = K # number of latent types
10        self.N = N # number of observations
11        self.V = V # number of distinct features
12        self.z = z # number of covariates affecting outcome
13        self.q = q # number of covariates affecting type shares
14        self.eta = eta
15        self.alpha = alpha
16
17    def model(self, C, Z, Q, Y=None, X=None):
18        # Supervised topic model with covariates
19
20        # Y : regression outcome
21        # X : feature count matrix
22        # C : total number of features per observation
23        # Z : covariates entering regression
24        # Q : covariates entering type shares
25        # K : number of types
26        # eta, alpha : Dirichlet hyperparameters
27
28        ##### Upstream Factor Model #####
29
30        with plate("topics", self.K):
31            beta = sample("beta", dist.Dirichlet(
32                self.eta * jnp.ones(self.V - self.num_anchors_per_class)))
33
34        phis = sample("phis", dist.Normal(0,2).expand([self.q, self.K-1]))
35
36        with plate_stack("docs", sizes = [self.N, self.K - 1]):
37            A = sample("A", dist.Normal(jnp.matmul(Q, phis) , self.alpha))
38
39        # document-topic distributions
40        theta = deterministic(
41            "theta",
42            softmax(jnp.hstack([A, jnp.zeros([self.D, 1])]), axis = -1)
43        )
44
45        distMultinomial = dist.Multinomial(
46            total_count=C,
47            probs = jnp.matmul(theta, beta)
48        )
49        with plate("hist", self.N):
50            X_bows = sample("obs_x", distMultinomial, obs = X)
51
52        ##### Downstream Regression Model #####
53
54        gammas = sample("gammas", dist.Normal(0, 10).expand([self.K-1]))
55        zetas = sample("zetas", dist.Normal(0, 10).expand([self.Z]))
56        sigma = sample("sigma", dist.Gamma(1, 10))
57
58        mean = jnp.matmul(theta[:,:(self.K-1)], gammas) + jnp.matmul(Z, zetas)
59
60        with plate("y", self.N):
61            Y = sample("obs_y", dist.Normal(mean, sigma), obs = Y)
```

Figure C.1: Numpyro’s code used to estimate Supervised Topic Model with Covariates