

# DYNAMICS OF HOMELESSNESS IN URBAN AMERICA<sup>1</sup>

BY CHRIS GLYNN AND EMILY B. FOX

*University of New Hampshire and University of Washington*

The relationship between housing costs and homelessness has important implications for the way that city and county governments respond to increasing homeless populations. Though many analyses in the public policy literature have examined inter-community variation in homelessness *rates* to identify causal mechanisms of homelessness [*J. Urban Aff.* **35** (2013) 607–625; *J. Urban Aff.* **25** (2003) 335–356; *Am. J. Publ. Health* **103** (2013) S340–S347], few studies have examined time-varying homeless *counts* within the same community [*J. Mod. Appl. Stat. Methods* **15** (2016) 15]. To examine trends in homeless population counts in the 25 largest U.S. metropolitan areas, we develop a dynamic Bayesian hierarchical model for time-varying homeless count data. Particular care is given to modeling uncertainty in the homeless count generating and measurement processes, and a critical distinction is made between the counted number of homeless and the true size of the homeless population. For each metro under study, we investigate the relationship between increases in the Zillow Rent Index and increases in the homeless population. Sensitivity of inference to potential improvements in the accuracy of point-in-time counts is explored, and evidence is presented that the inferred increase in the rate of homelessness from 2011–2016 depends on prior beliefs about the accuracy of homeless counts. A main finding of the study is that the relationship between homelessness and rental costs is strongest in New York, Los Angeles, Washington, D.C., and Seattle.

**1. Introduction.** Counts of people experiencing homelessness in cities such as Seattle, Los Angeles, and New York reveal alarming year-over-year increases in the raw numbers of enumerated individuals. In addition to rising counts of homeless, rental costs in these cities are significantly increasing as well. The relationship between housing costs and homelessness is a topic of great public importance and has received considerable attention [Fargo et al. (2013), Hanratty (2017), Stojanovic et al. (1999), Byrne et al. (2013), O’Flaherty (1995), Sclar (1990)].

Several challenges exist in quantifying the impact of increased rental costs on the size of the homeless population. The first challenge is that point-in-time homeless counts often occur on a single night in January and are thus subject to significant sampling variability. The second challenge is that the accuracy of the count itself is not the same from one year to the next. Differences in the number of volunteers, weather, and count methodologies lead to counts that are difficult to compare year-over-year.

---

Received November 2017; revised May 2018.

<sup>1</sup>Supported by Zillow, AFOSR Grant FA9550-16-1-0038, and NSF CAREER Award IIS-1350133.

*Key words and phrases.* Homelessness, housing costs, missing data, state-space.

These facts beg the question: are homeless populations across the county increasing? Or do the reported counts simply represent a higher fraction of the homeless population? Changing count accuracy over time directly impacts inferred trends in the size of the homeless population. In light of this fact, we investigate the impact of different count accuracy trajectories on the inferred change in homelessness rates from 2011–2016.

Inference on the relationship between trends in rental costs and trends in the homeless population is related to other trend analyses with data quality challenges [Tokdar et al. (2011), Coles and Sparks (2006), Cornulier et al. (2011), Kery and Royle (2010)]. Although we observe the number of *counted* homeless, we do not observe the true size of the homeless population. Plant-capture methods [Laska and Meisner (1993), Schwarz and Seber (1999)] have demonstrated that homeless counts systematically understate the size of the total homeless population [Hopper et al. (2008)]. One strategy to include the uncounted number of homeless in the analysis is to build a mechanism for the imperfect counting process into the statistical model, as McCandless et al. (2016) have done with plant-capture data from Edmonton, Canada.

In this paper, the total size of the homeless population is imputed, and uncertainty in the total homeless population and count accuracy is propagated to our assessment of the relationship between rental costs and homelessness. Our goal is to *jointly* model the collection of homeless count time series from the 25 largest metropolitan areas in the United States. In contrast to McCandless et al. (2016), who treat time-indexed counts as exchangeable, we directly model temporal dependence. We develop a Bayesian dynamic modeling framework to investigate the relationship between the number of homeless and rent costs subject to different prior beliefs about count accuracy over time.

The data in our analysis comes from three sources: the U.S. Census Bureau; the U.S. Department of Housing and Urban Development (HUD); and the housing website Zillow. The data include the total population, point-in-time homeless counts, and the Zillow Rent Index (ZRI) for continuums of care that service the 25 largest metro areas from 2011–2016. The continuum of care (CoC) is an administrative unit that seeks to integrate local homeless services, counts, and data under a single community organization [Office of Community Planning and Development, Dept. of Housing and Urban Development (2009)].

Numerous previous studies have utilized inter-community variation in homelessness *rates* to identify potential causal mechanisms of homelessness [Fargo et al. (2013), Byrne et al. (2013), Raphael (2010), Lee, Price-Spratlen and Kanan (2003), Early and Olsen (2002), Quigley and Raphael (2001), Quigley, Raphael and Smolensky (2001), Troutman, Jackson and Ekelund (1999), Hudson (1998), Grimes and Chressanthis (1997), Honig and Filer (1993), Burt (1992), Bohanon (1991), Appelbaum et al. (1991), Quigley (1990)]. Studies that model homelessness rates, defined as  $\frac{\text{total homeless}}{\text{total population}}$ , assume that both the numerator and denominator are observed without error. In practice, there is significant uncertainty in both

the numerator and denominator in any such homelessness rate calculation. To account for that uncertainty, we directly model time-varying *counts* within the same community. Working with time series of count data has two advantages. First, statistical models of counts more aptly characterize the sampling variability in the observed data; and second, focusing on within community variation over time avoids drawing conclusions from data generated across different municipal and state governments, climates, and social structures. A major contribution of our work is the development of a statistical framework that enables researchers, policymakers, and local continuum coordinators to address five specific questions for each metro:

(Q1) When adjusting for increases in count accuracy and total population growth, is the *rate* of homelessness increasing?

(Q2) If ZRI increases by  $x\%$ , what are the predicted increases in the counted and total number of people experiencing homelessness?

(Q3) How can we quantify uncertainty in the reported point-in-time counts?

(Q4) Given that  $C$  homeless are counted and count accuracy is imperfect, what is the expected range in the total number of people experiencing homelessness at a point in time?

(Q5) What is the one-year-ahead forecast of the total homeless population in 2017?

We identify New York, Los Angeles, Seattle, and Washington, D.C. as metros where (i) the inferred rate of homelessness significantly increased from 2011–2016 and (ii) there exists a strong relationship between housing costs and homelessness. We present evidence that the inferred change in the homelessness rate from 2011–2016 is sensitive to the trajectory of count accuracy. This point is emphasized to encourage researchers, policymakers, and continuum leaders to carefully quantify their beliefs and uncertainty about count accuracy.

The prior beliefs that we incorporate in this analysis are informed by existing literature and discussions with count coordinators, volunteers, and homelessness experts from around the country. Incorporating the expert opinions of count coordinators in every metro in the sample will lead to a more informed study. Our goal in this paper is to advance the statistical methodology utilized by researchers to analyze data on homelessness. We view this as a demonstration of a predictive modeling framework that will benefit from a partnership between private companies with relevant data, HUD, and local continuums of care.

Homelessness is a complex problem with many potential contributing factors. Though previous research has consistently identified the cost of rental housing as a significant predictor of homelessness [see, e.g., [Corinth \(2015\)](#), [Byrne et al. \(2013\)](#), [Quigley, Raphael and Smolensky \(2001\)](#)], we recognize that apartment vacancy rates, social safety net measures, changes in affordable/supportive housing policy, and local unemployment rates could contribute to increased homelessness as well. There are many paths to homelessness with complicated interactions among a host of factors; however, data scarcity at the metro level precludes robust estimation of

many correlated covariate effects. In the presence of limited data, it is essential to focus on a sparse predictive model. For this reason, we utilize changes in rental costs as our sole predictor and emphasize that we cannot draw causal conclusions from our analysis.

In Section 2, we discuss the data used in our analysis and necessary pre-processing steps to account for geographic mismatches between counties and continuums of care. Section 3 describes the Bayesian dynamic model that hierarchically shares information across all metros under study. Efficient information sharing, both locally in time and hierarchically across all metros, facilitates sharper inference on the relationship between rental costs and homelessness. Our hierarchical dynamic model allows us to estimate local relationships between homelessness and rental costs whereas the cross-sectional regression model in [Byrne et al. \(2013\)](#) estimates a single global effect. We discuss prior information and how that information translates to prior distributions for model parameters in Section 4. Model fitting with a custom Markov chain Monte Carlo algorithm is discussed in Section 5, and Section 6 presents results and addresses questions (Q1)–(Q5). Section 7 concludes with a discussion of our findings.

**2. Data.** The data in our study comes from three different sources: the U.S. Census Bureau, HUD, and the housing website Zillow. For the continuums of care in the 25 largest metros, we observe a collection of three time series that correspond to (i) the total number of people living in the metro, (ii) the *counted* number of homeless in the continuum(s) of that metro, and (iii) the ZRI for the metro.

For the total population data, we use county-level population estimates reported by the U.S. Census Bureau [[U.S. Census Bureau \(2016\)](#)]. The homeless counts are the number of individuals experiencing homelessness (both sheltered and unsheltered) at a point-in-time as reported by HUD [[U.S. Department of Housing and Urban Development \(2016\)](#)]. While the first two series of interest are fairly self-explanatory, the use of ZRI warrants further discussion.

According to Zillow's description in [Bun \(2012\)](#), ZRI is computed so that it does not depend on the set of homes that are currently for rent. Using proprietary statistical models trained on market rent data, Zillow estimates the market rent for every home, regardless of whether that home is currently for rent. The rent estimates of individual homes account for home attributes, physical characteristics, prior sales, tax assessments, and geographic location. The ZRI is computed to be the median of the market rent estimates for the housing stock in each metro.

One potential concern is that Zillow's estimated rent levels are skewed high as a consequence of training data from a biased sample of high quality homes. To mitigate the impact of this potential bias, we focus on percent differences in ZRI from one year to the next,  $\Delta\text{ZRI}$ , and do not directly utilize the level of ZRI itself.

One of the challenges in working with the HUD point-in-time data is that the jurisdiction of the HUD-defined continuums of care do not always agree with the boundaries of cities or counties. Often, each county will have a single continuum;

however, cases exist where this is not true. In some counties, there may be more than one continuum (e.g., Cook County, IL and Fulton County, GA have two). Other times, there may be multiple counties in a single continuum (e.g., the Denver, CO continuum spans seven different counties, and the New York City continuum spans five). Table 1 maps the 25 metros under study to the underlying HUD continuum(s) of care. In each metro, if the continuum does not match up with a single county, we construct a synthetic unit of analysis by following one of two approaches. If a county includes multiple continuums, we aggregate homeless totals

TABLE 1

*HUD continuums of care and state counties that correspond to the 25 largest metropolitan areas under study. In cases where more than one continuum of care is in a county, we aggregate homeless counts to form a synthetic continuum for that county. When a single continuum spans multiple counties, we construct a population-weighted ZRI measure and aggregate total population figures across the multiple counties*

	<b>Metro area</b>	<b>HUD continuum of care</b>	<b>Counties</b>
1	New York, NY	NY-600	New York, Bronx, Queens, Kings, Richmond
2	Los Angeles-Long Beach-Anaheim, CA	CA-600, CA-606, CA-607, CA-612	Los Angeles
3	Chicago, IL	IL-510, IL-511	Cook
4	Dallas-Fort Worth, TX	TX-600	Dallas
5	Philadelphia, PA	PA-500	Philadelphia
6	Houston, TX	TX-700	Harris, Fort Bend
7	Washington, DC	DC-500	District of Columbia
8	Miami-Fort Lauderdale, FL	FL-600	Miami-Dade
9	Atlanta, GA	GA-500, GA-502	Fulton
10	Boston, MA	MA-500	Suffolk
11	San Francisco, CA	CA-501	San Francisco
12	Detroit, MI	MI-501, MI-502	Wayne
13	Riverside, CA	CA-608	Riverside
14	Phoenix, AZ	AZ-502	Maricopa
15	Seattle, WA	WA-500	King
16	Minneapolis-St Paul, MN	MN-500	Hennepin
17	San Diego, CA	CA-601	San Diego
18	St. Louis, MO	MO-500, MO-501	St. Louis
19	Tampa, FL	FL-501	Hillsborough
20	Baltimore, MD	MD-501, MD-505	Baltimore
21	Denver, CO	CO-503	Adams, Arapahoe, Boulder, Broomfield, Denver, Douglas, Jefferson
22	Pittsburgh, PA	PA-600	Allegheny
23	Portland, OR	OR-501	Multnomah
24	Charlotte, NC	NC-505	Mecklenburg
25	Sacramento, CA	CA-503	Sacramento

reported by each continuum in the county. If a continuum includes multiple counties, we aggregate population totals reported by the different counties that make up a single continuum and construct a population-weighted ZRI metric.

We also focus on year-over-year changes in the metro-specific ZRI rather than the absolute level of the ZRI itself. This standardizes the analysis of rental markets across metros. The result is a data set of synthetic continuums that properly record the counted number of homeless, total population, and changes in rent levels in each metro. Because the ZRI is only available after October 2010, the time series for these three quantities are observed from 2011–2016 at an annual frequency. The homeless count data and ZRI are recorded each January, but the intercensal total population estimates from the Census Bureau are dated July 1. There is a six month temporal mismatch in both the homeless count and ZRI and the total metro population series. To account for this mismatch, we linearly interpolate the census population estimates to align the data so that the average of year  $t$  and  $t - 1$  is an estimate of the population in January for metro  $i$ . Figure 1 presents these three time series for the All Home King County continuum in Seattle, WA.

The count of homeless in Seattle/King County has dramatically increased since 2014 [Figure 1(b)]; however, the total population [Figure 1(a)] also significantly increased over that same time period. The King County, WA data demonstrate the need for modeling the homelessness rate to control for increases in the total population. The ZRI for King County, shown in Figure 1(c), has similarly increased.

In order to properly calculate the homelessness rate, it is necessary to account for time-variation in the count accuracy. We define the count accuracy to be the probability that a person who is homeless will be accounted for in the homeless count. If the count accuracy improves over time, more homeless are likely to be

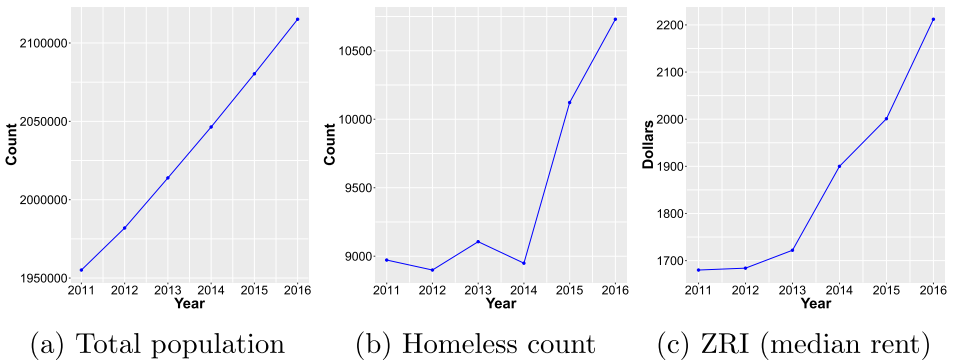


FIG. 1. Data from the All Home King County (WA) continuum of care from 2011–2016. Left: the total population in King County has rapidly increased in recent years. Increased population creates increased demand for rental housing and community services. Middle: The number of homeless counted in King County has dramatically increased since 2014. Right: The median rent, as measured by the ZRI, demonstrates the same basic pattern of increases as the count of people experiencing homelessness.

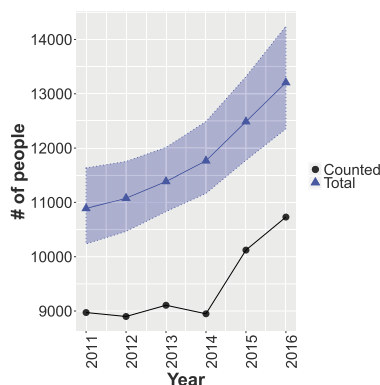


FIG. 2. Counted number of homeless (points) and imputed mean (triangles) and 95% predicted interval (shaded region) of the total number of homeless in King County. In this illustration, we assume the expected count accuracy is constant over time when inferring the distribution of the total number of homeless.

counted. In a scenario where the homeless count has improved, an increase in the number of homeless counted does not necessarily imply that the total size of the homeless population has increased. The count may simply represent a higher fraction of the total homeless population.

To account for the homeless not included in the HUD-reported count data, we impute the total homeless population in each metro from 2011–2016 and examine the impact of different trajectories in the count accuracy on the inferred homelessness rate. Figure 2 illustrates the distribution of the unobserved total number of homeless over time in King County, WA if we assume that the count accuracy does not improve with time and approximately 80% of homeless are included in the count.

It is reasonable to assume time variation in the count accuracy: in many continuums, the count accuracy may incrementally improve each year; in some continuums, the count accuracy could degrade over time due to lack of funding; in others, the accuracy may jump at a single year. A primary objective of our study is to assess the impact of different trajectories in the count accuracy on the relationship between homelessness rates and changes in ZRI. The model and prior distribution for different trajectories of the count accuracy will be discussed further in Sections 3.3 and 4.1.

**3. Model.** In this section, we develop a joint statistical model for collections of population-level and subpopulation counts. For each metro, we model (i) the number of homeless counted, (ii) the true number of homeless, and (iii) the total number of people living in the metro. Of the three quantities, only two are observed: the homeless counted and the total number of people. The true number of people experiencing homelessness is not observed, and we treat it as missing data.

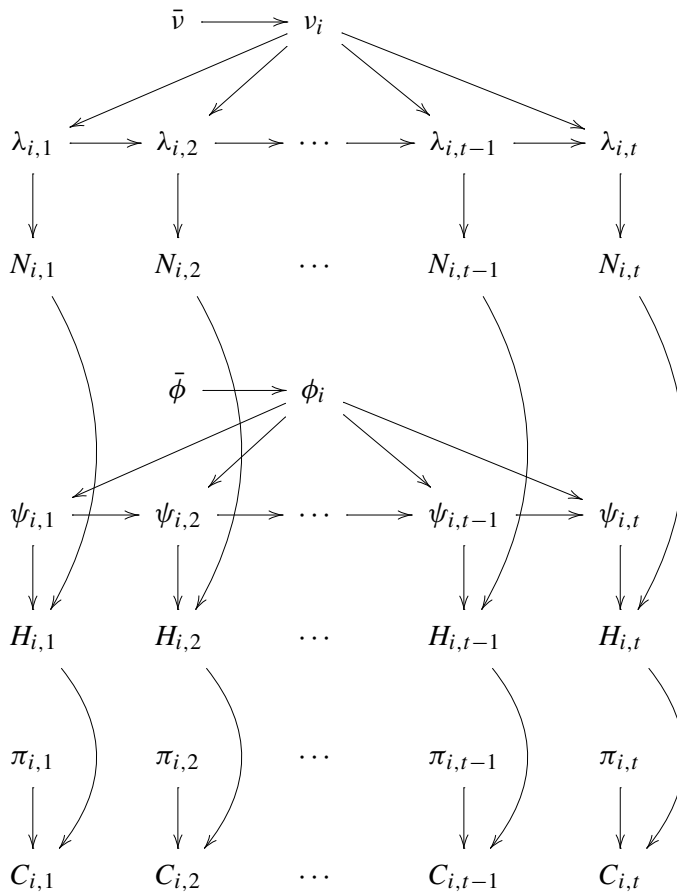


FIG. 3. Graphical model of the continuum-level homeless population. In metro  $i$  in year  $t$ , the total population is modeled by  $N_{i,t}$ , the homeless population is modeled by  $H_{i,t}$ , and the number of homeless counted is modeled by  $C_{i,t}$ . The dynamical processes and associated parameters are outlined in Sections 3.1–3.3.

The total population of a metro (as reported by the Census) is modeled as a noisy observation of the true total population.

Figure 3 is the graphical representation of our dynamic Bayesian hierarchical model. The random variable  $N_{i,t}$  is the total number of people that live in metro  $i$  in year  $t$ . The  $N_{i,1:T}$  variables depend on the dynamic process governing population growth,  $\lambda_{i,1:T}$ . The expected increase in population from one year to the next in metro  $i$  is modeled by parameter  $v_i$ , and the global population growth is modeled by parameter  $\bar{v}$ . Section 3.1 develops the total population model in greater detail.

The total number of homeless,  $H_{i,t}$ , depends on  $N_{i,t}$  and the probability of being homeless,  $p_{i,t}$ . The log odds of homelessness,  $\psi_{i,t} = \log\left(\frac{p_{i,t}}{1-p_{i,t}}\right)$ , is modeled by a



dynamic process that depends on changes in ZRI. The relationship between change in ZRI and the log odds of homelessness is modeled hierarchically by parameter  $\phi_i$  with global mean  $\bar{\phi}$ . The full model for  $H_{i,t}$  and the dynamics of  $\psi_{i,t}$  are discussed in Section 3.2. The counted number of homeless,  $C_{i,t}$ , depends on  $H_{i,t}$  and the probability that a homeless person is counted,  $\pi_{i,t}$ . We call  $\pi_{i,t}$  the count accuracy and discuss the count data generating process in Section 3.3.

3.1. *Total population model.* Significant interest lies in homelessness rates that facilitate comparison across different metros. The total population of the metro, or the denominator in a rate calculation, is uncertain. The intercensal population estimates reported annually by the U.S. Census Bureau are noisy. To properly quantify the uncertainty in the estimated homelessness rate of each metro, it is necessary to account for the uncertainty in the total population size.

A secondary reason for modeling the total population is that it facilitates forecasting. In order to forecast the size of the homeless population in future years, it is necessary to know the size of the future total population. A dynamic model for the total population enables a model-based forecast of the homeless population.

The total population size for metro  $i$  in year  $t$ ,  $N_{i,t}$ , is modeled as a time-indexed Poisson random variable that allows for growth and decay in the population of each metro.

$$(3.1) \quad N_{i,t} \sim \text{Poisson}(\lambda_{i,t}),$$

$$(3.2) \quad \lambda_{i,t} = \bar{\lambda}_i \theta_{i,t}.$$

The Poisson rate,  $\lambda_{i,t}$ , is the product of a static scale factor,  $\bar{\lambda}_i$ , and a latent time-varying component  $\theta_{i,t}$ . The dynamics of  $\lambda_{i,t}$  are driven by a dynamic process on the unit interval,  $\theta_{i,t} \in (0, 1)$ . Modeling  $\lambda_{i,t}$  as the product of the scaling factor and the dynamic term  $\theta_{i,t}$  provides an intuitive and computationally tractable dynamic model for Poisson counts. An auxiliary Poisson–Binomial thinning step for efficient computation is discussed in Section 5.1. The unit-interval-constrained dynamic process  $\theta_{i,1:T}$  is constructed with the logistic transformation and a real-valued stochastic process  $\eta_{i,1:T}$ .

$$(3.3) \quad \theta_{i,t} = \frac{e^{\eta_{i,t}}}{1 + e^{\eta_{i,t}}},$$

$$(3.4) \quad \eta_{i,t} = \eta_{i,t-1} + v_i + v_{i,t}, \quad v_{i,t} \sim N(0, \sigma_{\eta_i}^2).$$

The nonstationary  $\eta_{i,1:T}$  process is a random walk with a metro-specific drift term,  $v_i$ . The drift component is aimed at modeling population dynamics in cities like Seattle and Detroit. In Seattle, the population is rapidly growing which would correspond to a positive drift ( $v_i > 0$ ). On the other hand, the population in Detroit has recently decreased, which would correspond to negative drift ( $v_i < 0$ ). To borrow

information across metros, we model the drift components hierarchically. The parameter  $\bar{v}$  may be interpreted as the expected drift in population across all metros.

$$(3.5) \quad v_i = \bar{v} + \varepsilon_i, \quad \varepsilon_i \sim N(0, \sigma_{v_i}^2),$$

$$(3.6) \quad \bar{v} \sim N(0, \sigma_{\bar{v}}^2).$$

Because  $\eta_{i,1:T}$  is nonstationary, the Poisson marginals  $p(N_{i,t}|\lambda_{i,t}), \dots, p(N_{i,T}|\lambda_{i,T})$  are not identically distributed. Consequently,  $N_{i,1:T}|\lambda_{i,1:T}$  is a nonstationary process for the total population counts. While the PoINAR method of Aldor-Noiman et al. (2016) would be suitable for stationary population modeling, the expected total populations in these metro areas are clearly changing over time. A second modeling alternative would be to transform the large counts with a natural logarithm and model the transformed response with a Gaussian dynamic model. Despite the computational simplicity of such a model, we prefer to directly model the count data with discrete, time-varying distributions to more aptly characterize the uncertainty in the observed data.

*3.2. Homeless population model.* In a metro with  $N_{i,t}$  total residents, some small fraction of the residents will be homeless. For this reason, it is natural to model the total number of people experiencing homelessness,  $H_{i,t}$ , with a time-indexed binomial distribution. The binomial parameter  $p_{i,t}$  is the unobserved probability that a person in metro  $i$  is homeless in year  $t$  (i.e., the homelessness rate).

$$(3.7) \quad H_{i,t}|N_{i,t}, p_{i,t} \sim \text{Binomial}(N_{i,t}, p_{i,t}).$$

One of our primary objectives is to include ZRI as a covariate in the dynamic model for the homeless probability  $p_{i,t}$ . We achieve this by modeling the log odds of homelessness.

$$(3.8) \quad p_{i,t} = \frac{e^{\psi_{i,t}}}{1 + e^{\psi_{i,t}}},$$

$$(3.9) \quad \psi_{i,t} = \psi_{i,t-1} + \phi_i \Delta \text{ZRI}_{i,t} + w_{i,t}, \quad w_{i,t} \sim N(0, \sigma_w^2).$$

The dynamic process that controls the homelessness rate,  $\psi_{i,1:T}$ , linearly depends on the year-over-year rate of change in the ZRI,  $\Delta \text{ZRI}_{i,t}$ .

$$(3.10) \quad \Delta \text{ZRI}_{i,t} = \frac{\text{ZRI}_{i,t} - \text{ZRI}_{i,t-1}}{\text{ZRI}_{i,t-1}}.$$

The regression coefficient  $\phi_i$  models the relationship between change in rent levels and change in homelessness rates. As a concrete example, if ZRI increases by 1% in continuum  $i$  from one year to the next, the expected log odds of homelessness will increase by  $0.01\phi_i$ .

The connection between increased rental costs and homelessness rates is well established in the homelessness literature [Hanratty (2017), Fargo et al. (2013), Byrne et al. (2013), Stojanovic et al. (1999), O'Flaherty (1995), Sclar (1990)]. To model this positive relationship, the prior distribution for regression coefficient  $\phi_i$  favors values greater than zero but no truncation at zero is forced. The data inform the degree to which increasing rent levels are associated with higher homelessness rates. The parameter  $\phi_i$  is modeled hierarchically across metros to borrow strength and provide a more robust estimation of the relationship between rent increases and homelessness.

$$(3.11) \quad \phi_i \sim N(\bar{\phi}, \sigma_{\phi_i}^2),$$

$$(3.12) \quad \bar{\phi} \sim N(m_{\bar{\phi}}, \sigma_{\bar{\phi}}^2).$$

As noted at the beginning of Section 3,  $H_{i,t}$ , is not observed. Only the imperfect homeless count,  $C_{i,t}$ , is observed. In our study, we treat  $H_{i,t}$  as missing data and impute it to estimate each  $\phi_i$ . By modeling the relationship between  $\Delta ZRI_{i,t}$  and the imputed  $H_{i,t}$ , we obtain a more reliable quantification of the uncertainty in the posterior distribution for  $\phi_i$  and  $\bar{\phi}$ .

**3.3. Homeless count model.** Plant-capture studies and postcount surveys have demonstrated that homeless counts systematically understate the number of people experiencing homelessness [Hopper et al. (2008)]. While single night counts are imperfect, it is not clear that there exist feasible alternatives. Logistics, expenses, and privacy concerns preclude volunteers and continuums from counting every person without a home.

We model the imperfection in the homeless counts with a binomial thinning step. Of the true number of homeless,  $H_{i,t}$ , only  $C_{i,t}$  of them are counted. It is  $C_{i,t}$  that we observe.

$$(3.13) \quad C_{i,t} \sim \text{Binomial}(H_{i,t}, \pi_{i,t}),$$

$$(3.14) \quad \pi_{i,t} \sim \text{Beta}(a_{i,t}, b_{i,t}).$$

The parameter  $\pi_{i,t}$  is the probability that a person experiencing homelessness is counted, and it is modeled with a beta distribution. In other words,  $\pi_{i,t}$  is the accuracy of the count. If  $\pi_{i,t} = 1$ , the count in metro  $i$  in year  $t$  is perfectly accurate and every homeless person is counted.

Observe that we model  $\pi_{i,t}$  and  $\pi_{i,t-1}$  as independent random variables. While the count accuracy is surely time-varying and may exhibit trends, we believe there is no clear dependence of  $\pi_{i,t}$  on  $\pi_{i,t-1}$ . Factors driving count accuracy such as weather and volunteer turnout are unrelated across years. Even the count methodology utilized by a continuum may change from one year to the next. As an example, in 2017, the All Home King County continuum of care overhauled its count

methodology to enhance the accuracy [Beekman (2016)]. Rather than sending volunteers to known areas where homeless congregate, as in previous years, volunteers covered each census tract in the county. In addition, volunteers were lead by guides who were either currently or recently homeless themselves. As a result, the number of homeless counted in January 2017 was significantly higher than the number counted in January 2016. Due to changes in methodology, it is not necessarily accurate to conclude that the size of the homeless population,  $H_{i,t}$ , dramatically increased. By providing a mechanism for changes in count accuracy in each metro from one year to the next, it is possible to more reliably assess the local relationship between increased rental costs and the homeless population.

The count accuracy itself is an unknown quantity, and since we do not observe  $H_{i,t}$ , it is not possible to learn  $\pi_{i,t}$ . Instead of trying to learn  $\pi_{i,t}$ , we marginalize it out so that  $C_{i,t}|H_{i,t} \sim \text{Beta-Binomial}(H_{i,t}, a_{i,t}, b_{i,t})$ . Despite the lack of an underlying dynamic model for the count accuracy, we examine the impact of different time trends in  $E[\pi_{i,t}]$  on posterior inference for  $\phi_i$ . The trends are achieved through specification of the  $a_{i,t}$  and  $b_{i,t}$  parameters. Given the sequences of expected values and variances for  $\pi_{i,t}$ —which we assume are provided by the agencies conducting the counts—the hyperparameters  $a_{i,t}$  and  $b_{i,t}$  may be computed from (3.15) and (3.16).

$$(3.15) \quad a_{i,t} = E[\pi_{i,t}] \left( \frac{(1 - E[\pi_{i,t}])E[\pi_{i,t}]}{\text{Var}(\pi_{i,t})} - 1 \right),$$

$$(3.16) \quad b_{i,t} = \frac{\text{Var}(\pi_{i,t})}{E[\pi_{i,t}]^2} \left( \frac{a_{i,t}^2}{E[\pi_{i,t}]} + a_{i,t} \right).$$

By modeling each  $\pi_{i,t}$  with an independent beta distribution, it is possible to easily achieve different types of accuracy trajectories. We discuss three trajectories of specific interest in Section 4.1.

**4. Prior distributions.** With limited data for each metro, it is critically important to elicit well-informed prior distributions. Information from data that predates 2011, existing literature, and the expert opinion of homeless count coordinators have been combined to elicit prior distributions for four components of the model: (i) the count accuracy that is formalized through  $\pi_{i,1:T}$  (Section 4.1); (ii) the relationship between homelessness and rising rental costs, which is modeled with the regression coefficient  $\phi_i$  (Section 4.2); (iii) the dynamic process  $\eta_{i,1:T}$  that governs the total population (Section 4.3); and (iv) the dynamic process  $\psi_{i,1:T}$  that governs the total homeless subpopulation dynamics (Section 4.4).

4.1. *Priors for count accuracy.* We base our prior distribution for count accuracy on a study by Hopper et al. (2008), who report evidence that 60–70% of unsheltered individuals in New York were visible and included in the city’s 2005 count. They discuss one plant-capture study where only 59% of participants were

counted. The number of homeless used in our study includes *sheltered* homeless as well. Hopper et al. (2008) note that counts of sheltered homeless are more reliable than the counts of unsheltered homeless. To elicit our prior for count accuracy, we compute a weighted average of accuracy for sheltered and unsheltered populations, respectively. We use homeless counts from 2010,  $C_{i,0}$ , to compute this weighted average in year  $t = 0$ .

$$(4.1) \quad E[\pi_{i,0}] = (0.95) \frac{C_{i,0}^{\text{sheltered}}}{C_{i,0}} + (0.6) \frac{C_{i,0}^{\text{unsheltered}}}{C_{i,0}}.$$

Our prior expectation is that the probability that a sheltered homeless person is included in the homeless count is 0.95, which allows for a small discrepancy between the true number and counted number of sheltered homeless due to administrative and other count errors that may occur. From the Hopper study, we believe the probability that an unsheltered homeless person is included in the homeless count to be approximately 0.6. Because each metro has a different proportion of sheltered and unsheltered homeless, each metro is assigned a unique baseline prior distribution for count accuracy based on the 2010 data. We develop prior distributions for  $\pi_{i,1:T}$  that exhibit different expected trajectories: constant, linear, and step functions in time. In each of these cases,  $\text{Var}(\pi_{i,t}) = 0.0015$  is chosen so that reasonable prior mass covers the  $E[\pi_{i,t}] \pm 0.10$  interval.

The constant case corresponds to a count that utilizes relatively consistent procedures and resources from one year to the next. In this case, the mean and variance of the accuracy are constant over time (for all  $t$ ,  $E[\pi_{i,t}] = E[\pi_{i,0}]$ ). With  $E[\pi_{i,t}]$  and  $\text{Var}(\pi_{i,t})$ , calculation of  $a_{i,t}$  and  $b_{i,t}$  follows directly from (3.15) and (3.16). The prior for  $\pi_{i,1:T}$  with constant count accuracy in King County, WA is presented in Figure 4(a).

The linear case corresponds to a count where the accuracy incrementally improves by a fixed amount (called  $\delta_i$ ) until it reaches one, as shown in (4.2). We assume that  $\delta_i$  is known and is ideally specified by the agency conducting the count. Alternatively,  $\delta_i$  can be adjusted to examine sensitivity of inference to different accuracy scenarios. This is the approach adopted here. As an example, to consider an increase of  $\bar{\delta}$  in the accuracy of the unsheltered homeless count,  $\delta_i$  is computed as in (4.3). We assume that sheltered homeless are always counted with 95% accuracy, and the improvement in accuracy of  $\bar{\delta}$  applies only to the unsheltered count.

$$(4.2) \quad E[\pi_{i,t+1}] = \min(E[\pi_{i,t}] + \delta_i, 1),$$

$$(4.3) \quad \delta_i = \bar{\delta} \left( \frac{C_{i,0}^{\text{sheltered}}}{C_{i,0}} \right).$$

In the step scenario, the accuracy dramatically increases at a specific point in time due to improved count methodology. This is observed in practice with the All

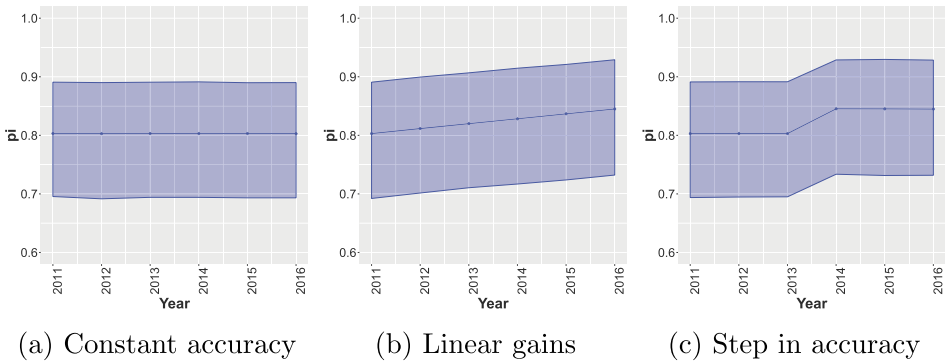


FIG. 4. Different prior beliefs about the trajectory of  $\pi_{i,1:T}$  in King County, WA. The solid lines are the expected count accuracy over time, and the shaded interval corresponds to the 99% prior uncertainty interval. Left: constant count accuracy. Middle: incremental (linear) increases in count accuracy. Right: Step in count accuracy.

Home King County continuum as discussed in Section 3.3. In this case, we assume that the year of change for metro  $i$ ,  $\tau_i$ , is known. For  $t < \tau_i$ ,  $E[\pi_{i,t}] = E[\pi_{i,0}]$ . For  $t \geq \tau_i$ ,  $E[\pi_{i,t}] = E[\pi_{i,0}] + \delta_i$ . A step in  $E[\pi_{i,t}]$  occurs at time  $\tau_i$ . Figure 4(c) illustrates a hypothetical step in count accuracy for King County in 2014. It is possible that there could be multiple steps for each metro and that there exists sequences  $\tau_i^1, \tau_i^2, \dots$  and  $\delta_i^1, \delta_i^2, \dots$  where steps of different size occur in different years. Because we do not know  $\tau_i$  for each metro, we do not investigate the step scenario further; however, with consultation of each local coordinator, this may be a very promising area of future work.

4.2. Priors for  $\phi_i$  and  $\bar{\phi}$ . We use previous work by Byrne et al. (2013) to form the basis of our prior distribution for  $\bar{\phi}$ . Byrne et al. (2013) found that in metropolitan continuums, when median rent increased by \$100, the expected homelessness rate increased by 6.34%. The average log odds of homelessness across all continuums in 2010 was  $\bar{f}_0 := -5.5$ . The average ZRI across continuums in 2010 was \$1534. So a \$100 increase in median rent would translate to a percent change in ZRI of  $\frac{100}{1534} \approx 6.5\%$ . This leads to calculation of the expectation of  $\bar{\phi}$  based on  $\bar{f}_0$ , the 6.5% increase in ZRI, and the expected increase in the homelessness rate of 6.34%:

$$(4.4) \quad \frac{1 + \exp\{-\bar{f}_0\}}{1 + \exp\{-\bar{f}_0 - \frac{\$100}{\$1534}m_{\bar{\phi}}\}} = 1.0634.$$

We calculate that  $m_{\bar{\phi}} = 0.94$ . Because of differences in methodology and data, we use  $\sigma_{\bar{\phi}}^2 = 0.005$  so that there is reasonable prior uncertainty about  $\bar{\phi}$ . We let  $\sigma_{\phi_i}^2 = 0.05$  so that there is modest shrinkage of each local effect toward the global



FIG. 5. Implied prior and posterior distribution of % change in homelessness rate with increases in ZRI for an arbitrary metro  $i$ . The triangle-marked (dashed) line is the prior (posterior) mean, and the shaded regions with solid (dotted) boundaries mark the 95% prior (posterior) credible interval.

mean  $\bar{\phi}$ . To examine the prior uncertainty in the relationship between increases in ZRI and increases in homelessness implied by our choices of  $m_{\bar{\phi}}$ ,  $\sigma_{\bar{\phi}}^2$ , and  $\sigma_{\phi_i}^2$ , we simulate from the marginal prior distribution for percent changes in the homelessness rate (see Figure 5).

Although we inform our prior using the results from Byrne et al. (2013), whose methodology we are trying to advance, notice a few things in Figure 5. One is that our prior is diffuse, and it becomes increasingly diffuse with larger percent increases in ZRI. Second, the inferred posterior concentrates on different values than the prior, indicating that we are indeed learning from data. Third, very little posterior mass is placed on values less than zero, providing evidence for the positive relationship between rising rents and homelessness. The conclusion is that using the Byrne et al. result is a useful way to center our prior.

4.3. *Prior for  $\eta_{i,1:T}$ .* The sampling distribution for the total population,  $N_{i,t}$ , depends on  $\eta_{i,t}$  through the Poisson rate,  $\lambda_{i,t}$  [refer to (3.1)–(3.3)]. Recall that  $\lambda_{i,t}$  is the product of a scaling factor,  $\bar{\lambda}_i$ , and a dynamic process on the unit interval,  $\theta_{i,t}$ . We let the expectation of  $\lambda_{i,0}$  be the 2010 population,  $N_{i,0}$ . This is achieved by fixing  $\bar{\lambda}_i = 2 \times N_{i,0}$  and  $E[\theta_{i,0}] = 0.5$  (i.e.,  $E[\eta_{i,0}] = 0$ ). The prior variance of  $\eta_{i,0}$  is fixed to be 0.0001, as we are confident that the Poisson rate of the total population in 2010 is the observed total population.

We let  $v_i \sim N(\bar{v}, 0.01)$  and  $\bar{v} \sim N(0, 0.005)$ . The innovation variance of the  $\eta_{i,1:T}$  process is fixed to be 0.0001 so that  $v_i$  primarily drives changes in the Poisson rate. The implied marginal distribution of  $N_{i,1:T}$  in King County, WA is presented in Figure 6(a). Observe that the distribution is centered at the 2010 King County population and allows significant uncertainty over the six year period. While the prior variance on each  $\eta_{i,t}$  is relatively small, the large magnitude of the scaling factor,  $\bar{\lambda}_i$ , results in a relatively diffuse marginal distribution for  $N_{i,t}$ .

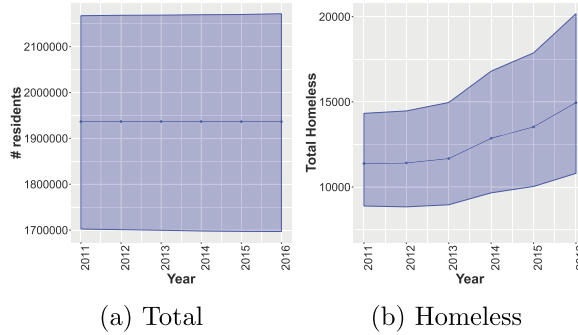


FIG. 6. Marginal prior distributions for  $N_{i,t}$  and  $H_{i,t}|ZRI_{1:T}$  in King County, WA. The solid lines are the prior means and the shaded regions are the 95% prior uncertainty intervals. Left: implied prior distributions for the total population,  $N_{i,1:T}$ . Right: Prior for total homeless population,  $H_{i,1:T}|ZRI_{1:T}$ . The upward trend in the implied prior for the total homeless population is due to observed increases in ZRI.

4.4. Prior for  $\psi_{i,1:T}$ . We utilize the counted number of homeless in 2010 to specify the prior expectation  $E[\psi_{i,0}]$ . The conditionally binomial sampling distribution for  $H_{i,0}|\psi_{i,0}, N_{i,0}$  in (3.7) yields the expectation

$$(4.5) \quad E[H_{i,0}|\psi_{i,0}, N_{i,0}] = \frac{1}{1 + e^{-\psi_{i,0}}} N_{i,0}.$$

Solving for  $\psi_{i,0}$  results in (4.6).

$$(4.6) \quad \psi_{i,0} = \log\left(\frac{E[H_{i,0}|\psi_{i,0}]/N_{i,0}}{1 - E[H_{i,0}|\psi_{i,0}]/N_{i,0}}\right).$$

Because we observe the noisy total population  $N_{i,0}$  from the Census estimate in 2010, we can compute a value for  $\psi_{i,0}$  given the expectation  $E[H_{i,0}|\psi_{i,0}]$ . Though we do not observe  $E[H_{i,0}|\psi_{i,0}]$ , we use an approximation to center the prior distribution of  $\psi_{i,0}$  and compensate for the approximation with moderate prior uncertainty. We approximate the expected total number of homeless in 2010 as an expected inflation of the observed count,  $E[H_{i,0}|\psi_{i,0}] \approx E[\frac{1}{\pi_{i,0}}]C_{i,0}$ .  $C_{i,0}$  is the 2010 homeless count value and  $E[1/\pi_{i,0}]$  is the expectation of the reciprocal count accuracy in 2010. The multiplier  $E[1/\pi_{i,0}]$  is evaluated by Monte Carlo simulation. This leads to a prior expectation as defined in (4.7).

$$(4.7) \quad E[\psi_{i,0}] := \log\left(\frac{E[\frac{1}{\pi_{i,0}}]C_{i,0}/N_{i,0}}{1 - E[\frac{1}{\pi_{i,0}}]C_{i,0}/N_{i,0}}\right).$$

The time zero variance is chosen to be  $\sigma_{\psi_0}^2 = 0.01$ , as this provides a three standard deviation interval of  $\pm 0.3$  around the prior mean. The result is that  $\psi_{i,0} \sim N(E[\psi_{i,0}], \sigma_{\psi_0}^2)$ . One reason that the prior specification of  $\psi_{i,0}$  is important is that the implied prior distribution for the total homeless population,  $H_{i,t}$ , depends on



$\psi_{i,t}$  [refer to (3.7) and (3.8)]. The implied prior distribution for  $H_{i,1:T} | ZRI_{1:T}$  is presented in Figure 6(b). The 95% prior interval spans a very reasonable range for each  $H_{i,t}$ .

The innovation variance of the dynamic process is fixed to be  $\sigma_\psi^2 = 0.001$  [see (3.9)]. We observe in synthetic data experiments that the innovation variance of  $\psi_{i,1:T}$  must be small in order to accurately learn  $\phi_i$  and  $\bar{\phi}$ . If  $\sigma_\psi^2$  is large relative to  $\text{Var}(\phi_i \Delta ZRI_{i,t})$ , changes in the homelessness rate are modeled as noise in  $\psi_{i,t}$  rather than driven by changes in ZRI. The ratio  $\frac{\text{Var}(\phi_i \Delta ZRI_{i,t})}{\sigma_\psi^2}$  can be thought of as a signal-to-noise ratio for  $\psi_{i,t}$ .

We also observe in synthetic data experiments that reliably inferring  $\psi_{i,1:T}$  and  $\phi_i$  requires that a metro’s homelessness rate exceed 0.05% of the total population. Because  $p_{i,t}$  is the logistic transformation of  $\psi_{i,t}$  [see (3.8)], the derivative  $\frac{dp_{i,t}}{d\psi_{i,t}} \rightarrow 0$  as  $|\psi_{i,t}|$  increases. Flat tails of  $p_{i,t}$  as a function of  $\psi_{i,t}$  mean that in metros with very low homeless rates, practically observed changes in homeless counts are consistent with a wide range of changes in ZRI. Under such conditions, it is not possible to reliably estimate  $\phi_i$ . Inference on  $\phi_i$  degrades along the continuum of decreasing  $\psi_{i,t}$ , but we set a limit based on our empirical studies. We do not trust inference for metros where the homelessness rate is less than 0.05%, or when  $\psi_{i,t} < -7.6$ .

**5. Markov chain Monte Carlo.** Our objective is to sample from the posterior distribution

$$(5.1) \quad p(H_{1:25,1:T}, \eta_{1:25,1:T}, \psi_{1:25,1:T}, \phi_{1:25}, \bar{\phi}, \nu_{1:25}, \bar{\nu} | N_{1:25,1:T}, C_{1:25,1:T}).$$

To sample from the posterior, we develop a custom Pólya–Gamma Gibbs sampler for dynamic Bayesian logistic regression [Polson, Scott and Windle (2013), Windle et al. (2013), Windle, Polson and Scott (2014)]. The Pólya–Gamma augmentation strategy allows us to harness a forward filtering and backward sampling (FFBS) algorithm that is commonly used to fit Bayesian dynamic models [Frühwirth-Schnatter (1994), Carter and Kohn (1994)]. We found that a burn-in of 25,000 samples and 50,000 samples collected after burn-in were sufficient for reproducible inferences. The MCMC simulation took approximately four hours to run on a MacBook Pro.

5.1. *Sampling steps.* There are ten different sampling steps required in the MCMC algorithm. The first step is for an auxiliary random variable whose only purpose is to facilitate computation when  $N_{i,t} | \tilde{\lambda}_i, \theta_{i,t} \sim \text{Poisson}(\tilde{\lambda}_i \theta_{i,t})$  [refer to (3.1) and (3.2)]. To construct this marginal distribution, we model the auxiliary  $Z_{i,t} \sim \text{Poisson}(\tilde{\lambda}_i)$  and the observed  $N_{i,t}$  conditionally binomial,  $N_{i,t} | Z_{i,t}, \theta_{i,t} \sim \text{Binomial}(Z_{i,t}, \theta_{i,t})$ . The Binomial–Poisson thinning strategy results in the desired marginal distribution for  $N_{i,t}$  and a computationally tractable method for making

inference on  $\eta_{i,t}$ ,  $v_i$ , and  $\bar{v}$ . The full conditional for the auxiliary  $Z_{i,t}$  is shown in Step 1.

Steps 2 and 6 use Pólya–Gamma data augmentation to allow a forward filtering backward sampling strategy. Steps 3 and 7 sample the auxiliary Pólya–Gamma variables  $\omega_{i,t}$  and  $\zeta_{i,t}$ . The collection of auxiliary variables  $Z_{1:25,1:T}$ ,  $\omega_{1:25,1:T}$ , and  $\zeta_{1:25,1:T}$  are numerically integrated out from the posterior by discarding posterior samples. Each sampling step is outlined below.

1. For each  $i, t$ , sample the auxiliary  $Z_{i,t}$  from a shifted Poisson by first sampling  $j = Z_{i,t} - N_{i,t} | \tilde{\lambda}_i, \theta_{i,t} \sim \text{Poisson}((1 - \theta_{i,t})\tilde{\lambda}_i)$  and then fixing  $Z_{i,t} = j + N_{i,t}$ .
2. For each  $i$ , sample the dynamic process that governs total population growth,  $\eta_{i,1:T} | N_{i,1:T}, \omega_{i,1:T}$ , with an FFBS algorithm.

(a) compute forward filtered distribution  $\eta_{i,t} | N_{i,1:t}, Z_{i,1:t}, v_i, \omega_{i,1:t} \sim N(m_{i,t}, S_{i,t})$

- $S_{i,t} := (\omega_{i,t} + \frac{1}{S_{i,t-1} + \sigma_\eta^2})^{-1}$
- $m_{i,t} := S_{i,t}(N_{i,t} - \frac{1}{2}Z_{i,t} + \frac{m_{i,t-1} + v_i}{S_{i,t-1} + \sigma_v^2})$

(b) sample recursively  $\eta_{i,t} | \eta_{i,t+1}, N_{i,1:t}, \omega_{i,1:t} \sim N(\tilde{m}_{i,t}, \tilde{S}_{i,t})$

- $\tilde{S}_{i,t} := (\frac{1}{S_{i,t}} + \frac{1}{\sigma_\eta^2})^{-1}$
- $\tilde{m}_{i,t} := \tilde{S}_{i,t}(\frac{m_{i,t}}{S_{i,t}} + \frac{\eta_{i,t+1} - v_i}{\sigma_\eta^2})$

3. For each  $i, t$ , sample the auxiliary Pólya–Gamma random variates to augment the total population variable,  $\omega_{i,t} | Z_{i,t}, \eta_{i,t} \sim \text{PG}(Z_{i,t}, \eta_{i,t})$ .
4. For each  $i$ , sample the parameter controlling expected population growth in metro  $i$ ,  $v_i | \bar{v}, \eta_{i,1:T} \sim N(\tilde{m}_{v_i}, \tilde{\sigma}_{v_i}^2)$ .

- $\tilde{\sigma}_{v_i}^2 := (\frac{1}{C_0 + \sigma_\eta^2} + \frac{T-1}{\sigma_\eta^2} + \frac{1}{\sigma_{v_i}^2})^{-1}$
- $\tilde{m}_{v_i} := \tilde{\sigma}_{v_i}^2 (\frac{\eta_{i,1}}{C_0 + \sigma_\eta^2} + \frac{1}{\sigma_\eta^2} \sum_{t=2}^T (\eta_{i,t} - \eta_{i,t-1}) + \frac{\bar{v}}{\sigma_{v_i}^2})$

5. Sample the expected total population growth globally across metros,  $\bar{v} | v_{1:25} \sim N((\frac{N}{\sigma_{v_i}^2} + \frac{1}{\sigma_{\bar{v}}^2})^{-1} \frac{1}{\sigma_{v_i}^2} \sum_{i=1}^{25} v_i, (\frac{N}{\sigma_{v_i}^2} + \frac{1}{\sigma_{\bar{v}}^2})^{-1})$ .
6. For each  $i$ , sample the dynamic process for the log odds of homelessness,  $\psi_{i,1:T} | N_{i,1:T}, H_{i,1:T}, \phi_i, \omega_{i,1:T}$ , with an FFBS algorithm.

(a) compute forward filtered distribution  $\psi_{i,t} | N_{i,1:t}, H_{i,1:t}, \zeta_{i,1:t} \sim N(f_{i,t}, q_{i,t})$

- $q_{i,t} := (\zeta_{i,t} + \frac{1}{q_{i,t-1} + \sigma_\psi^2})^{-1}$
- $f_{i,t} := q_{i,t}(H_{i,t} - \frac{1}{2}N_{i,t} + \frac{f_{i,t-1} + \phi_i \Delta \text{ZRI}_{i,t}}{q_{i,t-1} + \sigma_\psi^2})$

(b) sample recursively  $\psi_{i,t}|\psi_{i,t+1}, N_{i,1:t}, H_{i,1:t}, \zeta_{i,1:t}, \phi_i \sim N(\tilde{f}_{i,t}, \tilde{q}_{i,t})$

- $\tilde{q}_{i,t} := (\frac{1}{q_{i,t}} + \frac{1}{\sigma_\psi^2})^{-1}$
- $\tilde{f}_{i,t} := \tilde{q}_{i,t}(\frac{f_{i,t}}{q_{i,t}} + \frac{\psi_{i,t+1} - \phi_i \Delta ZRI_{i,t}}{\sigma_\psi^2})$

7. For each  $i, t$ , sample the auxiliary Pólya–Gamma random variates to augment the total homeless variable,  $\zeta_{i,t}|N_{i,t}, \psi_{i,t} \sim \text{PG}(N_{i,t}, \psi_{i,t})$ .

8. For each  $i$ , sample the parameter governing the relationship between change in ZRI and change in homelessness in metro  $i$ ,  $\phi_i|\psi_{i,1:T}, \bar{\phi} \sim N(m_{\phi_i}, \Sigma_{\phi_i})$ .

- $\Sigma_{\phi_i} := (\frac{(\Delta ZRI_{i,1})^2}{\sigma_{\psi_0}^2 + \sigma_\psi^2} + \frac{\sum_{t=2}^T (\Delta ZRI_{i,t})^2}{\sigma_\psi^2} + \frac{1}{\sigma_\phi^2})$
- $m_{\phi_i} := \Sigma_{\phi_i}(\frac{\bar{\phi}}{\sigma_{\phi_i}^2} + \frac{\Delta ZRI_{i,1}(\psi_{i,1} - f_{i,0})}{\sigma_{\psi_0}^2 + \sigma_{\psi_i}^2} + \frac{\sum_{t=2}^T \Delta ZRI_{i,t}(\psi_{i,t} - \psi_{i,t-1})}{\sigma_{\psi_i}^2})$

9. Sample the global mean parameter for the change in ZRI and change in homelessness,  $\bar{\phi}|\phi_{1:25} \sim N((\frac{25}{\sigma_\phi^2} + \frac{1}{\sigma_\phi^2})^{-1}(\frac{1}{\sigma_\phi^2} \sum_{i=1}^{25} \phi_i + \frac{m_{\bar{\phi}}}{\sigma_\phi^2}), (\frac{25}{\sigma_\phi^2} + \frac{1}{\sigma_\phi^2})^{-1})$ .

10. For each  $i, t$ , sample the total number of people experiencing homelessness in metro  $i$  and year  $t$ ,  $H_{i,t}$ , from  $p(H_{i,t}|N_{i,t}, C_{i,t}, p_{i,t}, a_{i,t}, b_{i,t}) \propto \frac{\Gamma(H_{i,t}+1)}{\Gamma(C_{i,t}+1)\Gamma(H_{i,t}-C_{i,t}+1)} \frac{\Gamma(C_{i,t}+a_{i,t})\Gamma(H_{i,t}-C_{i,t}+b_{i,t})}{\Gamma(H_{i,t}+a_{i,t}+b_{i,t})} \frac{\Gamma(a_{i,t}+b_{i,t})}{\Gamma(a_{i,t})\Gamma(b_{i,t})} \dots \times \binom{N_{i,t}}{H_{i,t}} p_{i,t}^{H_{i,t}} (1 - p_{i,t})^{(N_{i,t}-H_{i,t})}$ .

Sampling  $H_{i,t}$  in 10 incorporates the beta-binomial marginal distribution for homeless counts, and it does not depend on  $\pi_{i,t}$  but instead on  $a_{i,t}$  and  $b_{i,t}$ . Sampling  $H_{i,t}$  requires sampling from a discrete distribution with support  $[C_{i,t}, N_{i,t}]$ . This large range creates a computational bottleneck as it involves evaluating densities at each value in the support. In practice, though, posterior probability is concentrated on values much closer to the lower end of the support. It is possible to speed up computation by setting a threshold after which the support is truncated. Once posterior probability falls below  $1 \times 10^{-8}$ , we stop evaluating the densities and truncate the support.

5.2. *Posterior predictive distributions.* To examine the predicted increase in total homeless populations associated with increases in ZRI, we utilize the posterior predictive distribution for the total homeless population in each metro,  $H_{i,t}|C_{1:25,1:T}, N_{1:25,1:T}$ . The main quantify of interest is the distribution of the increase in the homeless population when the observed change in ZRI,  $\Delta ZRI_{i,t}$ , increases by an  $x > 0$ . The increase is modeled by  $(H_{i,t}^x - H_{i,t})|C_{1:25,1:T}, N_{1:25,1:T}$ , which is the difference between the predicted homeless total for a change in ZRI of  $\Delta ZRI_{i,t} + x$  and the baseline prediction at  $\Delta ZRI_{i,t}$ .

We draw samples from this posterior with a three step procedure that approximates the integral:

$$(5.2) \quad \begin{aligned} & p(H_{i,t}^x - H_{i,t} | N_{1:25,1:T}, C_{1:25,1:T}) \\ &= \int p(H_{i,t}^x - H_{i,t} | \psi_{i,t}, \phi_i, N_{1:25,1:T}, C_{1:25,1:T}) p(\psi_{i,t}, \phi_i) d\psi_{i,t} d\phi_i. \end{aligned}$$

The procedure relies on the  $m$ th posterior sample of (i) the relationship between ZRI and homelessness,  $\phi_i^{(m)}$ , (ii) the log odds of homelessness,  $\psi_{i,t}^{(m)}$ , and (iii) the Census reported estimate of the total population  $N_{i,t}$ . The procedure is detailed below.

1. Construct the  $m$ th sample of log odds of homelessness where  $\Delta ZRI_{i,t}$  is increased by  $x$ .

$$(5.3) \quad \psi_{i,t}^{(m),x} = \psi_{i,t}^{(m)} + \phi_i^{(m)} x.$$

2. Generate a prediction for the total homeless population at  $\Delta ZRI_{i,t} + x$  by sampling

$$(5.4) \quad H_{i,t}^{(m),x} \sim \text{Binomial}(N_{i,t}, p_{i,t}^x),$$

where  $p_{i,t}^x$  is the same logistic transformation of  $\psi_{i,t}^x$  as in (3.8).

3. Compute the difference.

$$(5.5) \quad H_{i,t}^{(m),x} - H_{i,t}^{(m)}.$$

We go one step further and also examine the predicted change in counted homeless under increased ZRI,  $(C_{i,t}^{*,x} - C_{i,t}^*) | C_{1:25,1:T}, N_{1:25,1:T}$ . Samples from this distribution are drawn by thinning the  $m$ th MCMC samples  $H_{i,t}^{(m),x}$  and  $H_{i,t}^{(m)}$  with a beta-binomial step

$$(5.6) \quad C_{i,t}^{(m),*,x} \sim \text{Beta-Binomial}(H_{i,t}^{(m),x}, a_{i,t}, b_{i,t})$$

and computing the difference  $C_{i,t}^{(m),*,x} - C_{i,t}^{(m),*}$ .

*5.3. Reproducible MCMC inference.* We verify that our MCMC simulation generates reproducible inference about the relationship between increases in homelessness and increases in ZRI by examining the posterior distribution for  $\phi_i$ . Ten different MCMC simulations are run, and inferences from two simulations  $j$  and  $j'$  are compared by computing  $|E[\phi_i^{(j)}] - E[\phi_i^{(j')}]|$ . Figure 7 illustrates the largest deviation across simulations by computing  $\max_j |E[\phi_i^{(j)}] - E[\phi_i^{(1)}]|$  for each metro. Each point in the histogram corresponds to the largest difference in posterior mean in reference to the first simulation for each of the 25 metros. The small values of these maximum differences in Figure 7 give us confidence that our MCMC simulation generates reproducible inferences.

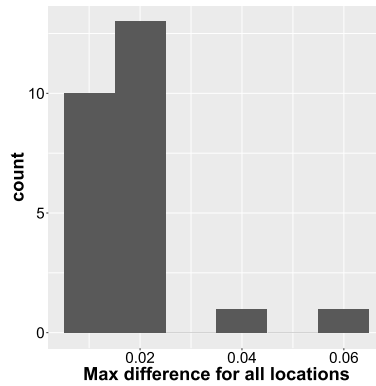


FIG. 7. The maximum difference in posterior means of  $\phi_1, \dots, \phi_{25}$  for 10 different MCMC simulations. Each of the 25 points in the histogram corresponds to  $\max_j |\phi_i^{(j)} - \phi_i^{(1)}|$ , for each of the 25 metros denoted by  $i$ . The superscript index  $j$  denotes the MCMC simulation. The very small differences indicate that our MCMC simulations generate reproducible inference in  $\phi_i$ .

**6. Results.** We seek to answer five questions, (Q1)–(Q5). Each of these questions is answered (in order) in Sections 6.1–6.5. In Section 6.1, we examine changes in homelessness rates from 2011–2016 across all metropolitan areas. In Section 6.2, the inferred relationship between increased ZRI and increases in homelessness is presented. Posterior predictive distributions for additional homeless counts are presented in Section 6.3, and the imputed distributions for the total number of homeless in each metro are presented in Section 6.4. Section 6.5 discusses our forecasts for the total homeless populations in 2017.

6.1. *Percent changes in the homelessness rate.* The inferred increases in homelessness rates from 2011–2016 are illustrated in Figure 8(a). We present results under two scenarios for the trajectory of the count accuracy: (i) the mean of the count accuracy is constant over time [i.e.,  $\bar{\delta}$  in (4.3) is zero]; and (ii) the mean of the unsheltered count accuracy increases by 2% annually until it reaches 100% (i.e.,  $\bar{\delta} = 0.02$ ).

Metros where the rate of homelessness almost certainly increased from 2011–2016 when  $\bar{\delta} = 0$  include New York, Los Angeles, Washington, D.C., Seattle, San Francisco, and Boston. For these cities, the 95% posterior credible interval for the percent change in the homelessness rate is bounded below by 0%, giving confidence that these homelessness rates did in fact increase. For each of these metros, the posterior mean increase in the homelessness rate exceeds 10%. In response to its growing homeless population, the City of Seattle has declared an official state of emergency [Beekman and Broom (2015)]. We adopt this moniker and characterize these metros as in similar states of emergency.

Metros where the homelessness rate almost certainly decreased when  $\bar{\delta} = 0$  include Phoenix, St. Louis, Portland, Detroit, Baltimore, Atlanta, Charlotte, Hous-



FIG. 8. *Left: Posterior distribution in percent change in homelessness rate from 2011 to 2016. The middle point in each segment is the posterior mean and the line segment encompasses the 95% posterior credible interval. For each metro, there is a posterior presented when the count accuracy is modeled with a constant mean over time ( $\bar{\delta} = 0$ ) and a posterior when the count accuracy of unsheltered homeless improves by 2% each year ( $\bar{\delta} = 0.02$ ). Right: Sensitivity of the percent increase in homelessness rate from 2011–2016 to different choices of  $\bar{\delta}$  in Los Angeles. The vertical line marks the count accuracy in 2011.*

ton, Riverside, and Tampa. For these cities, the 95% posterior credible interval for the percent change in the homelessness rate is bounded above by 0%, giving confidence that the homelessness rate did in fact decrease. For all but Phoenix and St. Louis, the posterior mean decrease in the homelessness rate exceeded 10%, and it seems real progress has been made in reducing homelessness in this group.

A third group of cities exists where the percent change in the homelessness rate has neither significantly increased nor decreased in either scenario. The 95% posterior credible interval for the change in the homelessness rate includes zero in Miami, Minneapolis, Dallas, Philadelphia, Sacramento, Pittsburgh, Denver, Chicago, and San Diego. The situation remains largely unchanged in this group, and the current homelessness rate is the status quo.

Observe in Figure 8(a) that New York and Los Angeles exhibit different sensitivities to change in the count accuracy over time. In New York, a city with a predominantly sheltered population, the inferred percent increase is essentially unchanged between the two scenarios. In Los Angeles, a warm-weather city with a large unsheltered population, the difference between the  $\bar{\delta} = 0$  and  $\bar{\delta} = 0.02$  cases is large, as demonstrated by separation of the posterior means. Equation (4.3) demonstrates that, in metros with large unsheltered populations, a  $\bar{\delta}$  increase in the

accuracy of an unsheltered homeless count leads to large changes in the overall count accuracy,  $\pi_{i,t}$ .

In Figure 8(b), we examine how inference on the change in homelessness rates from 2011–2016 can change with different values of  $\bar{\delta}$ . We focus on Los Angeles, above considered in a “state of emergency” and see that the posterior distribution for the 2011–2016 change in the homelessness rate depends on the count accuracy. If the expected count accuracy in LA in 2011 was 0.69 and it remained unchanged from 2011–2016 (i.e.,  $\bar{\delta} = 0$ ), Figure 8(b) shows that the homelessness rate almost certainly increased over that time, with a posterior expected increase of 13.4%. On the other hand, if LA had improved its counting method over this time frame, the increased homeless counts could be explained away by the increased count accuracy rather than an actual increase in the homelessness rate. For example, if the expected count accuracy in LA in 2011 was 0.69 and it incrementally increased to 0.91 over that six year window (i.e.,  $\bar{\delta} = 6\%$ ), it is likely that the homelessness rate decreased, with a posterior expected change in the homelessness rate of  $-5.9\%$ . That is, the increased homeless count in LA is entirely explained by improvements in the count accuracy, and the overall homeless rate likely decreased over the six year window. The conclusion again is that any inferences drawn about changes in homeless populations are highly sensitive to assumptions about the count accuracy. Sensitivity analyses similar to the one presented in Figure 8(b) are presented for each metro in the Supplementary Material [Glynn and Fox (2019)].

*6.2. Rental costs and homelessness.* To examine the predicted increase in homelessness and homeless counts as ZRI increases, we focus on the posterior predictive distributions  $(H_{i,t}^x - H_{i,t})|C_{1:25,1:T}, N_{1:25,1:T}$  and  $(C_{i,t}^{*,x} - C_{i,t}^*)|C_{1:25,1:T}, N_{1:25,1:T}$ . Section 5.2 provides complete details for sampling from these distributions.

We find that, for a fixed percent increase in ZRI of  $x = 10\%$ , the predicted increase in homelessness is largest in New York and Los Angeles (see Figure 10). Predicted increases in homeless counts are robust to whether we set  $\bar{\delta} = 0$  or  $\bar{\delta} = 0.02$ , as one would hope; however, the predicted count increases map to different increases in total homelessness under different prior beliefs about  $\pi_{i,t}$ . In Figure 9, the posterior predictive distributions for the increase in total and counted homeless are illustrated for different increases in ZRI in New York and Los Angeles.

In New York, the large sheltered population and high count accuracy imply that the distributions of increased counts and total homeless populations are nearly identical [Figure 9(a)]. If the ZRI in New York increases by  $x = 10\%$ , given 2016 levels of homelessness, we expect that the homeless population will increase by 5413 people, with 95% posterior probability of the homelessness increase in New York being more than 2896 people and less than 10,523 people. In Los Angeles, the lower overall count accuracy implies more separation between the distributions of increased counted and total homeless [Figure 9(b)]. Under the same  $x = 10\%$

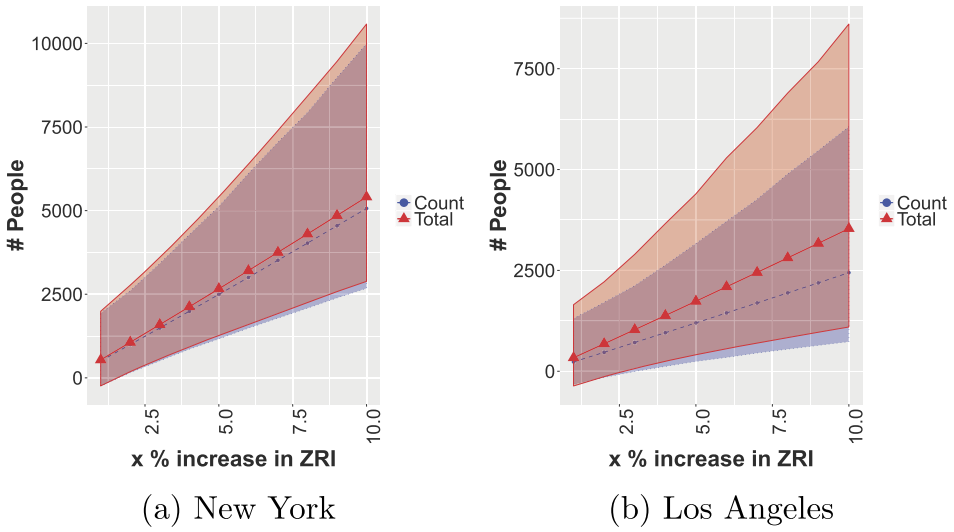


FIG. 9. Posterior predictive distributions of increased homeless counts  $(C_{i,t}^{*,x} - C_{i,t}^* | C_{1:25,1:T}, N_{1:25,1:T})$  and total homeless populations  $(H_{i,t}^x - H_{i,t} | C_{1:25,1:T}, N_{1:25,1:T})$  associated with increases in ZRI for both New York and Los Angeles when  $\delta = 0$ . The dashed line and shaded interval with dotted boundaries correspond to the posterior mean and 95% predictive interval of increases in the homeless count. The triangle-marked line and interval with solid boundaries correspond to the posterior mean and 95% predictive interval of increases in the total homeless population.

increase in ZRI in Los Angeles, we expect that 3536 people will become homeless, with 95% posterior probability of more than 1106 people and less than 8554 people.

Figure 10(a) summarizes the predicted increase in the total homeless population when ZRI increases by  $x = 10\%$  across all metros. The distributions of increases presented in Figure 10 account for the different sizes of metros with binomial sampling as shown in (5.4) and (5.6) (i.e., the values  $N_{i,t}$  are larger for larger metros). We expect the largest increases to occur in the largest metros (New York and Los Angeles), and this is confirmed by our analysis of the data. For the increase in the homeless population associated with increases in ZRI, we report the one-sided 95% posterior credible interval to shed light on the right tail of the distribution.

We note Seattle and Washington, D.C. as metros of interest, even though the 95% credible intervals for the predicted increase in the homeless population cross the zero threshold in Figure 10(b). In Seattle the posterior mean increase in the homeless population is 479 people and the posterior probability that the increase is positive is 0.942. A similar story emerges in Washington, D.C., where the posterior mean increase in homeless population is 386 people and the posterior probability that the increase is positive is 0.941. With at least 0.94 posterior probability of an increase in the homeless population, Seattle and Washington D.C. exhibit



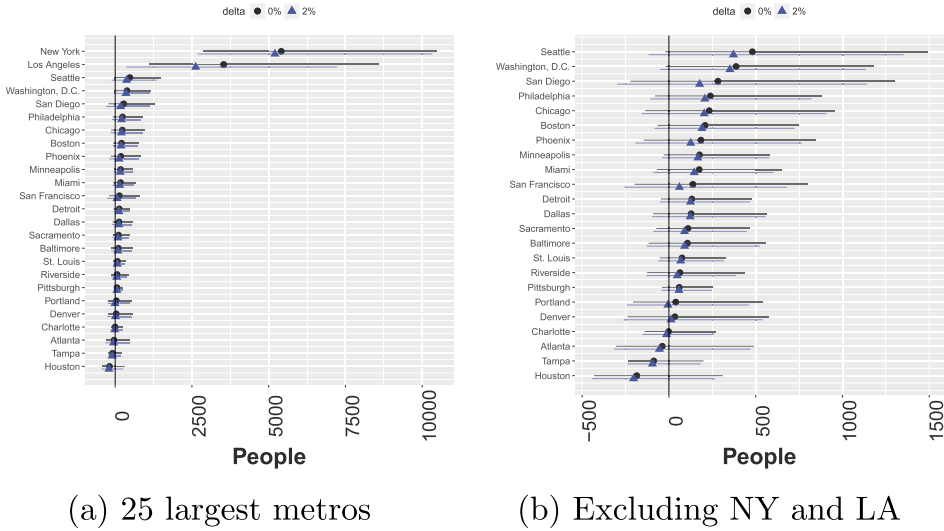


FIG. 10. Predicted increase in homeless population when ZRI increases by  $x = 10\%$  in 2016. The points are the posterior mean of  $(H_{i,T}^x - H_{i,T}) | C_{1:25,1:T}, N_{1:25,1:T}$ . The line segment spans the one-sided (right-tail) 95% posterior credible interval. In the left panel, results are presented for all metros. In the right panel, New York and Los Angeles are excluded for more careful inspection of the remaining 23 metros.

a meaningful statistical relationship between increased ZRI and increased homelessness. The posterior probability of an increase in the homeless population when ZRI increases by 10% is less than 0.9 for all metros but New York, Los Angeles, Washington, D.C., and Seattle. While the expected increases in homelessness in some metros may be large (e.g., San Diego), the variance in the posterior predictive distribution precludes us from confidently concluding that the predicted changes in the population are strictly positive. See Table 2 for posterior probabilities of homeless increases associated with ZRI increases of 10% in each metro. Predicted increases in homelessness as a function of increases in ZRI, as shown in Figure 9 for New York and Los Angeles, are available for each metro in the Supplementary Material [[Glynn and Fox (2019)]].

6.3. *Additional homeless counts.* The number of homeless that HUD reports in each continuum is from an annual point-in-time (PIT) count conducted in January. An important limitation of PIT counts is that no standard errors or other measures of uncertainty are reported, leaving decision makers without important context. One way of assessing sampling variability in PIT data is to construct posterior predictive distributions for the outcome of additional counts, allowing count coordinators to quantify uncertainty in the PIT data.

In this section, we report the posterior predictive distribution for this hypothetical second homeless count. The prediction, denoted by  $C_{i,t}^* | C_{1:25,1:T}, N_{1:25,1:T}$ ,

TABLE 2

Summary of posterior distributions across metros for 2016. The number of counted homeless reported by HUD is presented the first column. The Synthetic counts column corresponds to the posterior predictive distribution for a second hypothetical count in metro  $i$  in 2016,  $C_{i,T}^*|C_{1:25,1:T}, N_{1:25,1:T}$ . The Total homeless column corresponds to the posterior predictive distribution for the total number of homeless in metro  $i$  in 2016,  $H_{i,T}|C_{1:25,1:T}, N_{1:25,1:T}$ . The Forecasted homeless (2017) column is the posterior predictive distribution for the total homeless population in 2017,  $H_{i,T+1}|C_{1:25,1:T}, N_{1:25,1:T}$ . In all cases, the first number reported is the posterior mean, with the 95% posterior predictive interval in parenthesis. The Prob column is the posterior probability that the predicted change in the homeless population when rent rises 10% is greater than zero given that  $\delta = 0$ :  $P((H_{i,T}^{10} - H_{i,T} > 0)|C_{1:25,1:T}, N_{1:25,1:T})$

Metro	HUD	Synthetic count	Total homeless	Forecast (2017)	Prob
New York	73,523	74,272 (66,027, 80,959)	79,348 (75,317, 84,222)	79,404 (72,716, 86,775)	0.999
Los Angeles	46,874	45,315 (39,144, 51,621)	65,585 (60,565, 70,919)	67,285 (60,145, 75,034)	0.993
Chicago	6841	7231 (6399 8025)	8209 (7665 8833)	8223 (7434 9087)	0.849
Dallas	3810	3688 (3288 4056)	4180 (3981 4449)	4301 (3888 4769)	0.829
Philadelphia	6112	6045 (5386 6626)	6645 (6316 7074)	6705 (6109 7373)	0.892
Houston	4031	4709 (4202 5137)	5565 (5245 5708)	5907 (5387 6433)	0.114
Wash. D.C.	8350	8167 (7291 8858)	8719 (8415 9216)	8907 (8176 9733)	0.941
Miami	4235	4274 (3787 4739)	4919 (4615 5265)	4993 (4515 5514)	0.879
Atlanta	4546	5011 (4388 5618)	5689 (5221 6172)	5781 (5137 6471)	0.404
Boston	6240	6289 (5585 6858)	6687 (6345 7110)	6807 (6208 7459)	0.892
San Francisco	6996	6991 (6090 7916)	9547 (8859 10,328)	9610 (8669 10,644)	0.747
Detroit	2612	2871 (2515 3216)	3112 (2851 3397)	3121 (2778 3495)	0.890
Riverside	2165	2574 (2281 2846)	3518 (3376 3582)	3638 (3325 3959)	0.709
Phoenix	5702	5940 (5273 6593)	6997 (6597 7497)	7192 (6513 7954)	0.815
Seattle	10,730	10,604 (9350 11,889)	13,207 (12,355 14,245)	13,688 (12,311 15,273)	0.943
Minneapolis	3056	3379 (2957 3763)	3631 (3328 3953)	3771 (3348 4227)	0.920
San Diego	8669	8933 (7795 10,114)	11,899 (11,084 12,898)	12,164 (10,908 13,595)	0.817
St. Louis	1713	1750 (1544 1940)	1902 (1778 2057)	1939 (1731 2167)	0.833
Tampa	1817	2068 (1819 2296)	2579 (2397 2711)	2634 (2366 2908)	0.171
Baltimore	3488	3578 (3174 3965)	4045 (3784 4348)	4080 (3685 4512)	0.783
Denver	5728	6026 (5318 6727)	6685 (6210 7301)	6813 (6135 7603)	0.578
Pittsburgh	1156	1332 (1163 1466)	1423 (1297 1517)	1458 (1292 1626)	0.832
Portland	3914	4044 (3537 4559)	5175 (4780 5614)	5267 (4669 5927)	0.598
Charlotte	1818	2055 (1804 2272)	2248 (2065 2402)	2313 (2058 2583)	0.489
Sacramento	2500	2611 (2299 2918)	3195 (2975 3433)	3300 (2948 3682)	0.834

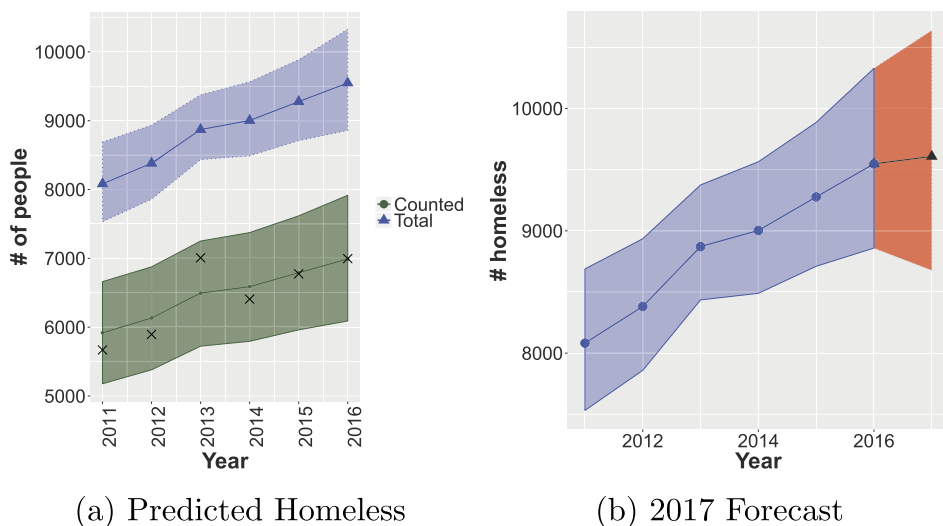


FIG. 11. Predicted homeless totals for San Francisco, CA. Left: Predicted number of counted homeless in additional (hypothetical) counts and predicted number of total homeless in San Francisco. The “x” marks are the actual HUD reported counts, and the solid line with points is the filtered and retrospectively smoothed mean of the posterior predictive distribution  $C_{i,1:T}^*|C_{1:25,1:T}, N_{1:25,1:T}$ . The line with triangles is the mean of the posterior predictive distribution  $H_{i,1:T}|C_{1:25,1:T}, N_{1:25,1:T}$ , and the associated shaded region is the 95% predictive interval. Right: Predicted number of total homeless in San Francisco, CA in 2017,  $H_{i,T+1}|C_{1:25,1:T}, N_{1:25,1:T}, ZR_{i,T+1}$ . The line with triangles is the mean of the out-of-sample prediction for 2017, and the shaded region with dotted boundaries is the out-of-sample predictive interval.

conditions on both the observed counts and the census reported total populations in all metros. Figure 11(a) presents the predicted outcome from additional homeless counts for San Francisco, a metro with one of the larger increases in the homelessness rate from 2011–2016. Observe that the posterior mean of  $C_{i,t}^*|C_{1:25,1:T}, N_{1:25,1:T}$  is a filtered and retrospectively smoothed quantity. The smoothing is apparent in 2013, when the HUD reported count appears to be an outlier relative to prior and subsequent HUD reported counts. Though the 2013 posterior mean is pulled slightly upward toward the reported count, the model does not overfit the data. In the remaining years, the posterior mean closely tracks the reported HUD counts. In 2016, the HUD reported count of homeless in San Francisco was 6996. If a second count were conducted in 2016, we expect the counted number of homeless would have been 6991, with 95% posterior probability of being more than 6090 and less than 7916.

The predictive distribution  $C_{i,t}^*|C_{1:25,1:T}, N_{1:25,1:T}$  provides policymakers and resource constrained counting agencies with a principled and data-driven way of conducting synthetic “additional” homeless counts to quantify uncertainty in the PIT data. The posterior predictive distributions for additional 2016 counts in all

metros are summarized in Table 2. Each metro has its own version of Figure 11(a) in the online Supplementary Material [[Glynn and Fox (2019)]].

6.4. *Imputed total number of homeless.* Imperfect count accuracy leads to count totals that are less than the size of the total homeless population. By modeling the mechanism of count accuracy, we are able to include the uncounted number of homeless in our estimate of the size of the total homeless population. In this section, we predict the total number of homeless in each metro and year. We report the posterior distribution  $H_{i,t}|C_{1:25,1:T}, N_{1:25,1:T}$ . Observe that the posterior distribution does not condition on the count accuracy parameter,  $\pi_{i,t}$ . The count accuracy has been integrated out; however, the variance of  $H_{i,t}|C_{1:25,1:T}, N_{1:25,1:T}$  is inextricably linked to the prior variance of the count accuracy,  $\pi_{i,t}$ . In this analysis, we fixed the prior variance to be 0.0015 so that prior mass would span the  $\pm 0.1$  interval. Though it is appealing to specify a diffuse prior for count accuracy, we found in practice that such a prior does not provide sufficient regularization. In settings with overly diffuse priors for count accuracy, inference for  $\phi_i$  was not reproducible across MCMC simulations. This highlights the importance of reliable prior information about count accuracy as it pertains to estimating the relationship between trends in ZRI and homelessness.

In Figure 11(a), observe that, because the counting process is imperfect, the expected total number of homeless is more than the counted number of homeless. In San Francisco, we expect that, in 2016, there were 9547 people experiencing homelessness, with 95% posterior probability that there were more than 8859 and fewer than 10,328. Table 2 presents the posterior mean and 95% credible interval for the total number of homeless in 2016 for each metro.

6.5. *Forecasts for 2017.* Resources to address the needs of a homeless population are budgeted well in advance of the January point-in-time count. In order to allocate resources in communities with growing (shrinking) homeless populations, a forecast of the next year's total homeless population is needed. In this section, we forecast the total homeless population in each metro in January 2017.

Our forecasts of the homeless population for 2017 take into account both predicted increases in the 2017 total population and the January 2017 ZRI value. We report the one-year-ahead forecast

$$(6.1) \quad H_{i,T+1}|C_{1:25,1:T}, N_{1:25,1:T}, \text{ZRI}_{1:25,T+1}.$$

For true out-of-sample forecasting when  $\text{ZRI}_{i,T+1}$  is not yet available, we could utilize Zillow's forecasted ZRI for metro  $i$ ,  $\hat{\text{ZRI}}_{i,T+1}$ .

Figure 11(b) illustrates the year-ahead forecast in San Francisco. Forecasting the total homeless population one year in the future requires forecasting both the year-ahead total metro population and log odds of homelessness. The uncertainty

in these component year-ahead forecasts accumulates, and the result is an uncertainty interval for the 2017 year-ahead forecast that increases relative to the intervals presented from 2011–2016. This is observed in Figure 11(b) as the predictive interval in 2017 fans out relative to the uncertainty interval from 2011–2016.

We predict that 9610 people experienced homelessness in San Francisco on any given January night in 2017, with 95% posterior probability of more than 8669 and less than 10,644 people. Although the January 2017 ZRI decreased 3.6% relative to its January 2016 value, we still expect a slight increase in San Francisco's total homeless population. The increase is largely driven by the model-based forecasted increase in San Francisco's total population. Each metro has a figure corresponding to Figure 11(b) presented in the Supplementary Material [[Glynn and Fox (2019)]]. The forecasted mean and 95% posterior predictive intervals for each metro are shown in Table 2.

**7. Discussion.** We presented statistical evidence that the relationship between rental costs and homelessness depends on one's beliefs about the time-varying accuracy of homeless counts. We highlight this fact to encourage public policy researchers, policymakers, and continuum leaders to carefully quantify their beliefs and uncertainty about count accuracy or the inferences drawn from studies relying on these counts. While the prior beliefs about count accuracy that we elicit in this paper are informed by existing literature and our discussions with count coordinators and homelessness experts from around the country, we believe that collecting expert opinions from every continuum can lead to a more robust and informed study. We encourage other researchers in this area to explicitly model variation in count accuracy when conducting their own analyses.

We use ZRI as a frequently updated measure of metro-level market rent. While ZRI is frequently updated and responsive to changing market conditions, it is a potentially biased measure of market rent given that it is computed from an incomplete sample of rental homes. To limit the impact of ZRI's absolute level on our analysis, we focus on year-over-year percent differences in ZRI, a strategy that allows us to investigate the relationship between changes in homelessness to changes in ZRI.

We found in synthetic data experiments that making accurate inference on the relationship between ZRI and homelessness with the model outlined in this paper requires homelessness rates that exceed 0.05% of the total population (Section 4.4). To work with counts as large as possible, this implies that sheltered and unsheltered homeless totals should be combined in a single analysis of a metro's total homeless population. Furthermore, reliable estimates of the entire homeless population are significantly aided by a metro having either a large sheltered population or high count accuracy of the unsheltered population; utilizing data on unsheltered homeless populations alone does not yield reliable results. This observation fits with our broader theme of data quality.

Modeling count accuracy is another place where we have directly addressed data quality challenges. Acknowledging that a continuum's homeless counts are imperfect and that the accuracy varies from one year to the next should not be viewed in any way as a failure of the count coordinators or volunteers. We view quantifying the count accuracy as an important step in accounting for the uncertainty inherent in such a difficult undertaking.

In our analysis, we have used the time-varying count accuracy to impute the size of the total homeless population. We believe it is natural in this application to think of the total homeless population as missing data. Assuming that the total population size is the observed count has two flaws. First, it understates the size of the homeless population. Second, it leads to overly confident estimates of regression coefficients  $\phi_1, \dots, \phi_{25}$ . Imputing the missing homeless population size naturally resolves both of these problems.

In this application, proper uncertainty quantification is critical. Counties, city governments, shelters, and health care providers are likely to benefit from an expected range of the homeless population size when they budget resources. By reporting the 95% posterior credible intervals on the total homeless population predictions and the 2017 forecasts, we emphasize the uncertain size of homeless populations now and in the future.

Metro-specific estimates of the relationship between rental costs and homelessness allow for each metro to make more informed policy decisions about affordable housing initiatives. We provide evidence that the homelessness rate significantly increased from 2011–2016 in six metros: New York, Los Angeles, Washington, D.C., Seattle, San Francisco, and Boston. While we are unable to conclude that increased rent *causes* homelessness, we found that large increases in ZRI are *statistically associated* with increases in the homeless population in four of those six metros: New York, Los Angeles, Washington, D.C., and Seattle. We believe that our results provide context for policy discussions that are happening in cities across the United States.

Homelessness is a product of many complicated and intertwined factors. Though an increase in rental costs is surely an important factor, vacancy rates, supportive services, affordable housing, affordable healthcare, and lack of economic opportunity may contribute as well. To investigate the complicated interactions of these factors at the local level with our current modeling framework, more frequent metro-level homeless counts are needed. In addition to counts occurring more often, estimating the relationship of these factors to increased homelessness also requires researchers to have more informed prior beliefs about count accuracy. One potential way of obtaining additional information on count accuracy is with post-count surveys at soup kitchens and other service providers, as in [Hopper et al. \(2008\)](#).

Due to limited information about count accuracy, we make several assumptions to aid our analysis. Most notably, we assume a prior distribution for count accuracy in each metro. The prior mean is determined by the proportion of unsheltered homeless in 2010 and the assumption that 60% of unsheltered homeless are

counted while 95% of sheltered homeless are counted. Naturally, the estimated size of the total homeless population is sensitive to these assumptions, with lower expected accuracy inflating the estimate of the total homeless population. Prior variance of the count accuracy prior also impacts how sharply we can estimate the relationship between homelessness and ZRI. The more diffuse the prior distribution for count accuracy, the more diffuse the posterior distribution for our  $\phi_i$  regression coefficients. While we believe the assumptions that we made to be reasonable based on existing literature and discussions with homeless experts and count coordinators, we again emphasize that stronger prior information is needed about count accuracy in every metro to sharpen this analysis.

There are additional limitations of our current approach. One is that it does not account for relocation in homeless populations. It is possible that people experiencing homelessness move to cities with more services and away from cities with fewer services. In this scenario, increases in the population of one metro are driven by decreases in another. Our present approach does not take into account network effects, and we assume that homeless relocation patterns are not a significant driver of trends across metros. At present, we do not know of data that would allow us to further investigate network effects. A second is that we rely on January count data. It is likely that seasonal patterns in homelessness exist, though our current annual data set does not provide insight into such seasonal fluctuations. Modeling network effects and seasonal fluctuations in homeless populations are important areas of future work.

**Acknowledgments.** The authors thank Surya Tokdar, Svenja Gudell, Cory Hopkins, Melissa Allison, Tom Byrne, and Dennis Culhane for helpful discussions.

#### SUPPLEMENTARY MATERIAL

**Metro-level supporting figures** (DOI: [10.1214/18-AOAS1200SUPP](https://doi.org/10.1214/18-AOAS1200SUPP); .pdf). As supplementary material, we present figures for each of the 25 largest metros in the United States from 2011–2017. For each metro, we present

(a) the posterior predictive distribution for homeless counts,  $C_{i,1:T}^* | C_{1:25,1:T}$ ,  $N_{1:25,1:T}$ , and the imputed total homeless population size,  $H_{i,1:T} | C_{1:25,1:T}$ ,  $H_{1:25,1:T}$ ;

(b) the predictive distribution for the total homeless population in 2017,  $H_{i,2017} | C_{1:25,1:T}$ ,  $N_{1:25,1:T}$ ;

(c) the posterior distribution of increase in the total homeless population with increases in ZRI; and

(d) the sensitivity of the inferred increase in the homelessness rate from 2011–2016 to different annual changes in count accuracy.

## REFERENCES

- ALDOR-NOIMAN, S., BROWN, L. D., FOX, E. B. and STINE, R. A. (2016). Spatio-temporal low count processes with application to violent crime events. *Statist. Sinica* **26** 1587–1610. MR3586230
- APPELBAUM, R. P., DOLNY, M., DREIER, P. and GILDERBLOOM, J. I. (1991). Scapegoating rent control: Masking the causes of homelessness. *J. Am. Plan. Assoc.* **57** 153–164.
- BEEKMAN, D. (2016). King county's homeless count could soar with new method of tallying. *Seattle Times*. [Online; accessed 05/29/2017].
- BEEKMAN, D. and BROOM, J. (2015). Mayor, county exec declare 'state of emergency' over homelessness. *Seattle Times*. [Online; accessed 04/15/2017].
- BOHANON, C. (1991). The economic correlates of homelessness in sixty cities. *Soc. Sci. Q.*
- BUN, Y. (2012). Zillow rent index: Methodology. <https://www.zillow.com/research/zillow-rent-index-methodology-2393/> [Online; accessed 04/2/2017].
- U.S. CENSUS BUREAU (2016). County population totals tables: 2010–2016. <https://www.census.gov/data/tables/2016/demo/popest/counties-total.html> [Online; accessed 04/2/2017].
- BURT, M. M. (1992). *Over the edge: The growth of homelessness in the 1980s*. Russell Sage Foundation.
- BYRNE, T., MUNLEY, E. A., FARGO, J. D., MONTGOMERY, A. E. and CULHANE, D. P. (2013). New perspectives on community-level determinants of homelessness. *J. Urban Aff.* **35** 607–625.
- CARTER, C. K. and KOHN, R. (1994). On Gibbs sampling for state space models. *Biometrika* **81** 541–553. MR1311096
- COLES, S. and SPARKS, R. (2006). Extreme value methods for modelling historical series of large volcanic magnitudes. *Statistics in Volcanology* **1** 47–56.
- CORINTH, K. C. (2015). Ending homelessness: More housing or fewer shelters? AEI Economics Working Papers 863788.
- CORNULIER, T., ROBINSON, R. A., ELSTON, D., LAMBIN, X., SUTHERLAND, W. J. and BENTON, T. G. (2011). Bayesian reconstitution of environmental change from disparate historical records: Hedgerow loss and farmland bird declines. *Methods Ecol. Evol.* **2** 86–94.
- EARLY, D. W. and OLSEN, E. O. (2002). Subsidized housing, emergency shelters, and homelessness: An empirical investigation using data from the 1990 census. *Adv. Econ. Anal. Policy* **2**(1).
- FARGO, J. D., MUNLEY, E. A., BYRNE, T. H., MONTGOMERY, A. E. and CULHANE, D. P. (2013). Community-level characteristics associated with variation in rates of homelessness among families and single adults. *Am. J. Publ. Health* **103**(S2) S340–S347.
- FRÜHWIRTH-SCHNATTER, S. (1994). Data augmentation and dynamic linear models. *J. Time Series Anal.* **15** 183–202. MR1263889
- GLYNN, C. and FOX, E. B. (2019). Supplement to “Dynamics of Homelessness in Urban America.” DOI:10.1214/18-AOAS1200SUPP.
- GRIMES, P. W. and CHRESSANTHIS, G. A. (1997). Assessing the effect of rent control on homelessness. *J. Urban Econ.* **41** 23–37.
- HANRATTY, M. (2017). Do local economic conditions affect homelessness? Impact of area housing market factors, unemployment, and poverty on community homeless rates. *Hous. Policy Debate* 1–16.
- HONIG, M. and FILER, R. K. (1993). Causes of intercity variation in homelessness. *Am. Econ. Rev.* 248–255.
- HOPPER, K., SHINN, M., LASKA, E., MEISNER, M. and WANDERLING, J. (2008). Estimating numbers of unsheltered homeless people through plant-capture and postcount survey methods. *Am. J. Publ. Health* **98** 1438–1442.
- HUDSON, C. G. (1998). Estimating homeless populations through structural equation modeling. *Soc. Choice Welf.* **25** 136.



- KERY, M. and ROYLE, A. J. (2010). Hierarchical modelling and estimation of abundance and population trends in metapopulation designs. *J. Anim. Ecol.* **79** 453–461.
- LASKA, E. M. and MEISNER, M. (1993). A plant-capture method for estimating the size of a population from a single sample. *Biometrics* 209–220.
- LEE, B. A., PRICE-SPRATLEN, T. and KANAN, J. W. (2003). Determinants of homelessness in metropolitan areas. *J. Urban Aff.* **25** 335–356.
- MCCANDLESS, L. C., PATTERSON, M. L., CURRIE, L. B., MONIRUZZAMAN, A. and SOMERS, J. M. (2016). Bayesian estimation of the size of a street-dwelling homeless population. *J. Mod. Appl. Stat. Methods* **15** 15.
- O'FLAHERTY, B. (1995). An economic theory of homelessness and housing. *J. Hous. Econ.* **4** 13–49.
- OFFICE OF COMMUNITY PLANNING AND DEVELOPMENT, DEPT. OF HOUSING AND URBAN DEVELOPMENT (2009). Continuum of care 101. <https://www.hudexchange.info/resources/documents/CoC101.pdf>. [Online; accessed 04/23/2018].
- U.S. DEPARTMENT OF HOUSING AND URBAN DEVELOPMENT (2016). Pit and hic data since 2007. <https://www.hudexchange.info/resource/3031/pit-and-hic-data-since-2007/>. [Online; accessed 04/2/2017].
- POLSON, N. G., SCOTT, J. G. and WINDLE, J. (2013). Bayesian inference for logistic models using Pólya–Gamma latent variables. *J. Amer. Statist. Assoc.* **108** 1339–1349. MR3174712
- QUIGLEY, J. M. (1990). Does rent control cause homelessness? Taking the claim seriously. *J. Policy Anal. Manage.* **9** 89–93.
- QUIGLEY, J. M. and RAPHAEL, S. (2001). The economics of homelessness: The evidence from North America. *European Journal of Housing Policy* **1** 323–336.
- QUIGLEY, J. M., RAPHAEL, S. and SMOLENSKY, E. (2001). Homeless in America, homeless in California. *Rev. Econ. Stat.* **83** 37–51.
- RAPHAEL, S. (2010). Housing market regulation and homelessness. In *How to House the Homeless* 110–135. Russell Sage Foundation, New York.
- SCHWARZ, C. J. and SEBER, G. A. (1999). Estimating animal abundance: Review III. *Statist. Sci.* 427–456.
- SCLAR, E. D. (1990). Homelessness and housing policy: A game of musical chairs. *Am. J. Publ. Health* **80** 1039–1040.
- STOJANOVIC, D., WEITZMAN, B. C., SHINN, M., LABAY, L. E. and WILLIAMS, N. P. (1999). Tracing the path out of homelessness: The housing patterns of families after exiting shelter. *Am. J. Community Psychol.* **27** 199–208.
- TOKDAR, S. T., GROSSMANN, I., KADANE, J. B., CHAREST, A.-S. and SMALL, M. J. (2011). Impact of beliefs about Atlantic tropical cyclone detection on conclusions about trends in tropical cyclone numbers. *Bayesian Anal.* **6** 547–572.
- TROUTMAN, W., JACKSON, J. D. and EKELUND, R. B. (1999). Public policy, perverse incentives, and the homeless problem. *Public Choice* **98** 195–212.
- WINDLE, J., POLSON, N. G. and SCOTT, J. G. (2014). Sampling Pólya–Gamma random variates: Alternate and approximate techniques. Available at [arXiv:1405.0506v1](https://arxiv.org/abs/1405.0506v1).
- WINDLE, J., CARVALHO, C. M., SCOTT, J. G. and SUN, L. (2013). Efficient data augmentation in dynamic models for binary and count data. Available at [arXiv:1308.0774](https://arxiv.org/abs/1308.0774).

PETER T. PAUL COLLEGE OF BUSINESS  
AND ECONOMICS  
UNIVERSITY OF NEW HAMPSHIRE  
DURHAM, NEW HAMPSHIRE 03824  
USA  
E-MAIL: [christopher.glynn@unh.edu](mailto:christopher.glynn@unh.edu)

PAUL G. ALLEN SCHOOL OF COMPUTER SCIENCE  
AND DEPARTMENT OF STATISTICS  
UNIVERSITY OF WASHINGTON  
SEATTLE, WASHINGTON 98195  
USA  
E-MAIL: [ebfox@uw.edu](mailto:ebfox@uw.edu)