000 001 002

003 004

010

011

012

013

014

015

016

017

018

019

021

023 024

025

LEARNING DISEASE PROGRESSION MODELS THAT CAPTURE HEALTH DISPARITIES

Anonymous authors

Paper under double-blind review

ABSTRACT

Disease progression models are widely used to inform the diagnosis and treatment of many progressive diseases. However, a significant limitation of existing models is that they do not account for health disparities that can bias the observed data. To address this, we develop an interpretable Bayesian disease progression model that captures three key health disparities: certain patient populations may (1) start receiving care only when their disease is more severe, (2) experience faster disease progression even while receiving care, or (3) receive follow-up care less frequently conditional on disease severity. We show theoretically and empirically that failing to account for disparities produces biased estimates of severity (underestimating severity for disadvantaged groups, for example). On a dataset of heart failure patients, we show that our model can identify groups that face each type of health disparity, and that accounting for these disparities meaningfully shifts which patients are considered high-risk.

1 INTRODUCTION

026 In many settings, observed data is used to model the progression of a latent variable over time. 027 Models of human aging use a person's physical and biological characteristics to model progression 028 of their latent "biological age" (Pierson et al., 2019); models of infrastructure deterioration use 029 inspection results to model progression of a system's latent overall health (Madanat et al., 1995); and disease progression models, which we focus on in this paper, use observed symptoms to model 031 progression of a patient's latent severity of a chronic disease (Wang et al., 2014). Disease progression models can help predict a patient's disease trajectory and thus personalize care, detect diseases at earlier stages, and guide drug development and clinical trial design (Mould et al., 2007; Romero 033 et al., 2015). They have been applied to a wide variety of progressive diseases such as Alzheimer's 034 disease (Holford & Peace, 1992) and cancer (Gupta & Bar-Joseph, 2008).

For the benefits of these models to apply to all patients equitably, it is crucial that they accurately describe progression for all populations of patients. However, disease progression models have 037 typically failed to account for the fact that systemic disparities in the healthcare process can bias the observed data that they are trained on. For example, disparities have been shown to arise along axes such as socioeconomic status (Weaver et al., 2010; Miller & Wherry, 2017), race (Yearby, 2018), 040 and proximity to care (Chan et al., 2006; Reilly, 2021). Accounting for such disparities is important 041 because it can meaningfully shift estimates of disease progression. For intuition, imagine learning 042 that a patient in the emergency room traveled three hours to get to the hospital; if their symptoms are 043 ambiguous, this contextual information may increase our estimate of how severe their underlying 044 condition is. Disease progression models have historically been unable to capture this type of social context—as we show later, this can lead to biased estimates of severity. To address this, we propose 046 a method for learning disease progression models that interpretably capture three well-documented 047 health disparities:

- 048
- 1. **Disparities in initial severity.** Certain patient groups may start receiving care only when their disease is more severe (Hu et al., 2024).
- Disparities in disease progression rate. Certain patient groups may experience faster disease progression, even while receiving care (Diamantidis et al., 2021).
 - 3. **Disparities in visit frequency.** Certain patient groups may visit healthcare providers for follow-up care less frequently, even at the same disease severity (Nouri et al., 2023).

054 It is a core technical challenge to design a model that is flexible enough to capture all three dispari-055 ties but still identifiable. Identifiability is necessary for accurate estimates of disparities and disease 056 progression. As such, our key contributions are: (1) we develop an interpretable Bayesian model 057 of disease progression that accounts for multiple types of disparities but remains provably identi-058 fiable from the observed data; (2) we prove and show empirically that failing to account for any of these three disparities leads to biased estimates of severity; and (3) we characterize fine-grained disparities in a heart failure dataset. Our model reveals that non-white patients have more severe 060 heart failure and face multiple types of health disparities: Black and Asian patients tend to start 061 receiving care at more severe stages of heart failure than do White patients, and Black patients see 062 healthcare providers for heart failure 10% less frequently than do White patients at the same disease 063 severity level. Accounting for these disparities meaningfully shifts our estimates of disease severity, 064 increasing the fraction of non-white patients identified as high-risk. While we ground our work in 065 healthcare, our method for learning progression models that account for disparities applies naturally 066 to many other progression model settings where disparities are of interest, including infrastructure 067 deterioration (Madanat et al., 1995) and human aging (Pierson et al., 2019).

068 069

2 RELATED WORK

070 071

Disease progression modeling. Disease progression models have been developed for many 072 chronic diseases, including Parkinson's disease (Post et al., 2005), Alzheimer's disease (Holford 073 & Peace, 1992), diabetes (Perveen et al., 2020), and cancer (Gupta & Bar-Joseph, 2008). A key 074 feature of the progression models we consider, common in the machine learning literature, is that 075 a latent severity Z_t progresses over time and gives rise to the observed symptoms X_t . Models in 076 this family include variants of hidden Markov models (HMMs) (Wang et al., 2014; Liu et al., 2015; 077 Alaa & Hu, 2017; Sukkar et al., 2012; Jackson et al., 2003) and recurrent neural networks (RNNs) (Choi et al., 2016b; Lipton et al., 2017; Lim & van der Schaar, 2018; Choi et al., 2016a; Ma et al., 079 2017; Kwon et al., 2019; Alaa & van der Schaar, 2019). The existing literature has not focused on 080 modeling disparities; we extend it by proposing a new approach to disease progression modeling 081 that can interpretably characterize and account for multiple types of health disparities.

082

Health disparities. Disparities have been documented in many parts of the healthcare process.
Factors such as distance from hospitals (Reilly, 2021), distrust of the healthcare system (LaVeist et al., 2009), or lack of insurance (Venkatesh et al., 2019) can result in underutilization of health services; biases in the judgements of healthcare providers can lead minority groups to receive later screening (Lee et al., 2021), fewer referrals (Landon et al., 2021), or generally worse care (Schäfer et al., 2016); and issues such as limited health literacy or trust can create disparities in follow-through for appointments or the effectiveness of at-home care (Davis, 1968; Brandon et al., 2005).

The existing literature has shown that disparities emerge along the three axes that we capture in this paper: (1) how severe a patient's disease becomes before they start to receive care (Chen et al., 2021; Iqbal et al., 2015; Hu et al., 2024); (2) how quickly their latent severity progresses even while receiving care (Diamantidis et al., 2021; Suarez et al., 2018); and (3) how likely they are to visit a healthcare provider at a given severity level (Nouri et al., 2023). Our goal is to show how accounting for disparities along all three of these axes improves the severity estimates of disease progression models, while also learning more fine-grained descriptions of existing disparities.

090

Capturing disparities with machine learning. We build upon a large body of past work that uses 098 machine learning as a tool to capture and address health disparities, including models that estimate the relative prevalence of underreported medical conditions (Shanmugam et al., 2021), improve risk 100 prediction for patients with missing outcome data (Balachandar et al., 2023), evaluate the impact 101 of race corrections in risk prediction (Zink et al., 2023), assess disparate impacts of AI in health-102 care (Chen et al., 2019), and quantify disparities in the performance of clinical prediction tasks 103 (Zhang et al., 2020). The closest work to our own is Chen et al. (2021), which develops a clus-104 tering algorithm that accounts for the fact that some patients do not come in (and are therefore not 105 observed) until later in their disease progression. While their work addresses one form of data bias that can arise due to health disparities, it differs from our own in two ways: it does not specifically 106 document or study health disparities, and it focuses on clustering patients as opposed to modeling 107 disease severity or progression. Our work proposes a model for capturing three types of health

disparities in the disease progression setting in order to learn precise descriptions of multiple disparities and make severity estimates that exhibit less bias than existing disease progression models.

3 MODEL

114 We build on a standard setup for disease progression 115 modeling, in which each patient has an underlying la-116 tent disease severity Z_t that progresses over time and 117 gives rise to a set of observed features X_t (Klemera & 118 Doubal, 2006; Levine, 2013).

119 We characterize each patient's severity $Z_t \in \mathbb{R}$ at time 120 t by their *initial severity* Z_0 at their first observation 121 (which we denote as t = 0) and their rate of progression 122 R after that point: 123

$$Z_t = Z_0 + R \cdot t$$

 $X_t = f(Z_t) + \epsilon_t$

 $\epsilon_t \sim \mathcal{N}(0, \Psi)$

If a patient visits a healthcare provider at time t, we 125 observe some recorded set of *features* $X_t \in \mathbb{R}^d$ (e.g., 126 lab results, imaging, symptoms). At any given visit, 127 a clinician does not necessarily observe or record all 128 features—we model the features that are observed as a 129 noisy function of their latent severity Z_t : 130

131

124

108

109

110 111 112

113

132

142

143

147 148

149

150

151

152 153

157

159

161

133 where the diagonal covariance matrix $\Psi \in \mathbb{R}^{d \times d}$ parameterizes feature-specific noise (accounting 134 for both measurement error and variation in how the patient's physical state can fluctuate day-to-135 day). In our experiments, we specifically instantiate f as a linear function $f(Z_t) = F \cdot Z_t + b$, 136 where $F \in \mathbb{R}^d$ is a feature-specific scaling factor and $b \in \mathbb{R}^d$ is a feature-specific intercept, but our 137 approach extends to more general parametric forms for f. We constrain the first feature $F_0 > 0$ using 138 domain knowledge; this restriction is necessary for identifiability because it restricts the mapping 139 between features and severity (Shapiro, 1985). We also observe a set of timesteps when a patient visits a healthcare provider; we discretize time and indicate whether a patient visits a healthcare 140 provider at time t with a binary indicator $D_t \in \{0, 1\}$. 141

Capturing disparities. Our model captures the three types of health disparities discussed in §2 by allowing model parameters to vary as a function of a patient's demographic feature vector A (Figure 144 1). For expositional clarity, we describe a setup where A encodes a single categorical label (e.g., 145 a patient's race group), but our approach naturally extends to multiple categorical groupings or to 146 continuous features.

> 1. Disparities in initial severity. Underserved patients may start receiving care only when their disease is more severe. We capture this by learning group-specific distributions of Z_0 , a patient's disease severity at their first visit. For one group $A = a_0$, we pin Z_0 to be drawn from a unit normal distribution; this is a standard and necessary identifiability condition since it fixes the scale of Z_t (Shapiro, 1985). For other groups a,

$$Z_0 \sim \mathcal{N}\left(\mu_{Z_0}^{(a)}, \sigma_{Z_0}^{(a)^2}\right)$$

where $\mu_{Z_0}^{(a)}$ and $\sigma_{Z_0}^{(a)}$ are learned group-specific parameters.

2. Disparities in disease progression rate. Underserved patients may experience faster disease progression even while receiving care. We capture this by learning group-specific distributions of disease progression rate R:

$$R \sim \mathcal{N}\left(\mu_R^{(a)}, {\sigma_R^{(a)}}^2\right)$$

where $\mu_{R}^{(a)}$ and $\sigma_{R}^{(a)}$ are learned group-specific parameters for each group a.



Figure 1: Disease progression generative model. Plate diagram captures N patients over T timesteps. Shaded nodes indicate observed features: demographics $A^{(i)}$, visit indicator $D_t^{(i)}$, and symptoms $X_t^{(i)}$ (only observed when $D_t^{(i)} = 1$). Unshaded nodes indicate latent variables: a patient's initial severity $Z_0^{(i)}$, rate of progression $R^{(i)}$, and severity $Z_t^{(i)}$. Red arrows indicate dependencies capturing health disparities.

3. **Disparities in visit frequency.** Underserved patients may visit healthcare providers for follow-up care less frequently at the same disease severity. We capture this by modeling patient visits as generated by an inhomogeneous Poisson process, parameterized by a time-varying rate parameter λ_t that depends on both Z_t and A:

$$\log(\lambda_t) = \beta_0 + \beta_Z \cdot Z_t + \beta_A^{(a)}$$

where β_Z and β_0 are learned parameters for the entire population and $\beta_A^{(a)}$ is a learned group-specific parameter for each group *a* (we pin $\beta_A^{(a_0)} = 0$ for reference).

171 Overall, our model parameters (on which we 172 place weakly informative priors) are the pa-173 rameters shared across groups $\{F, b, \Psi\},\$ 174 and the group-specific parameters $\{\mu_{Z_0}^{(a)}, \sigma_{Z_0}^{(a)}, \mu_R^{(a)}, \sigma_R^{(a)}, \beta_0, \beta_Z, \beta_A^{(a)}\}$. We learn posterior distributions over these parameters from 175 176 177 our observed data X_t, D_t, A using Hamiltonian 178 Monte Carlo, a standard algorithm for Bayesian 179 inference (Betancourt, 2018), as implemented 180 in Stan (Carpenter et al., 2017). Figure 1 sum-181 marizes the data generating process and Table 1 182 summarizes the notation for our model. 183

Model discussion. Modeling progression as linear over time is a common approach (Holford & Peace, 1992; Pierson et al., 2019), because it provides an interpretable characterization of the trajectory. The interpretability of using a single intercept and progression rate pa-

Notation Meaning				
$\overline{X_t}$	Observed features at time t			
D_t	Binary visit indicator for time t			
<u>A</u>	Demographic features			
$\overline{Z_t}$	Disease severity at time t			
Z_0	Initial severity			
R	Disease progression rate			
F	Severity-feature matrix			
b	Feature intercepts			
Ψ	Feature covariance matrix			
μ_{Z_0}, σ_{Z_0}	Group-specific mean and sd of Z_0			
μ_R, σ_R	Group-specific mean and sd of R			
λ_t	Visit rate at time t			
β_0	Visit rate intercept			
β_Z	Visit rate Z_t coefficient			
β_A	Visit rate A coefficient			

Table 1: **Summary of notation.** Observed data are listed above the double horizontal line.

rameter to characterize a patient's disease trajectory is especially valuable in our setting, allowing us to compare how severe groups are at initial presentation and how quickly they progress. Similarly, using a Poisson process to model event frequency is a common approach, including in work that seeks to capture disparities in event frequency (Liu et al., 2024; Kurashima et al., 2018).

194 195

196

162

163

164

166 167

169

170

4 THEORETICAL ANALYSIS

In this section, we prove two main theoretical results. First, we show that our model is *identifiable*,
a necessary condition for its parameters to be estimated from the observed data and interpreted.
Learning these parameter estimates is what allows us to characterize disparities. Second, we prove
that failing to account for disparities produces *biased estimates of severity*. We summarize proof
strategies in the main text and provide formal proofs in Appendices §A and §B.

202 203 4.1 IDENTIFIABILITY

204 We show that our model is identifiable, meaning different sets of parameters yield different observed 205 data distributions (Bellman & Åström, 1970) and thus that we can recover correct estimates of all 206 model parameters from the observed data. Learning a model of progression that is *flexible* enough to 207 characterize multiple disparities but still identifiable is a fundamental challenge. In fact, if we added 208 one more dependence on A — in particular, adding an arrow from A to X in Figure 1 — the model 209 would no longer be identifiable; without a shared interpretation across groups of how features map to severity, it would be impossible to identify disparities in disease progression. Put another way, 210 our model encodes the richest set of disparities on the observed data while retaining identifiability. 211

Theorem 4.1. All model parameters are identified by the observed data distribution $P(X_t, D_t \mid A)$.

As mentioned in §3, the distribution of initial severity Z_0 is pinned to a unit normal for one demographic group a_0 . This pinned distribution reduces the number of unknown latent parameters for group a_0 , allowing us to show that $\{F, b, \Psi\}$ are identified by $P(X_t \mid A = a_0)$. Having identified these, we show that the parameters $\{\mu_{Z_0}^{(a)}, \sigma_{Z_0}^{(a)}, \mu_R^{(a)}, \sigma_R^{(a)}\}$ are identified by $P(X_t \mid A = a)$ for all groups *a*. Finally, we show that given the previously identified parameters, $\{\beta_0, \beta_Z\}$ are identified by $P(D_t \mid A = a_0)$ and $\{\beta_A^{(a)}\}$ is identified by $P(D_t \mid A = a)$ for all other groups *a*.

220 221

237

4.2 BIAS IN MODELS THAT DO NOT ACCOUNT FOR DISPARITIES

222 Next we show that, when any of the health disparities we discuss are present, a model that does 223 not account for group-specific disparities will produce *biased estimates* of severity—i.e., $\mathbb{E}[Z_t]$ 224 $[X_t, D_t] \neq \mathbb{E}[Z_t \mid X_t, D_t, A = a]$. These theoretical results hold under more general assumptions 225 than our full parametric model: our assumptions, which we formally describe in Appendix B, are 226 that the model dependencies are encoded by the DAG in Figure 1; that severity Z_t increases linearly with progression rate R; and that visit rate λ_t increases with severity Z_t . For each proof, we analyze 227 the effect of one disparity — e.g., for disparities in initial severity, we assume that $P(Z_0 \mid A = a)$ 228 differs across groups — while keeping other distributions constant across groups. These results hold 229 in the presence of multiple disparities as long as existing disparities disfavor or favor the same group, 230 so as to not cancel each other out in their effects. 231

We quantify disparities by using the strict Monotone Likelihood Ratio Property (MLRP) to reason about the probability density functions of initial severity and progression rate for certain groups, relative to the overall population (Karlin & Rubin, 1956):

Definition 4.2. Two distributions characterized by probability density functions f(x) and g(x) have the strict monotone likelihood ratio property in x if $\frac{f(x)}{g(x)}$ is a strictly increasing function of x.

Intuitively, this means that as some variable x (Z_0 or R, in our case) gets larger, it is more likely to be drawn from f than g. The MLRP is a widely-used assumption across many settings (Gaebler & Goel, 2024; Anwar & Fang, 2006; Chemla & Hennessy, 2019); the normal, exponential, binomial, and Poisson families all have this property. For brevity, we say "f(x) strictly MLRPs g(x)" to mean that f(x) and g(x) satisfy the strict MLRP in x. We now prove for each disparity that any model that fails to account for the disparity will produce biased estimates of severity.

Theorem 4.3. A model that does not take into account disparities in initial disease severity Z_0 will underestimate the disease severity of groups with higher initial severity and overestimate that of groups with lower initial severity. Specifically, if $P(Z_0 | A = a)$ strictly MLRPs $P(Z_0)$ for some group a, then $\mathbb{E}[Z_t | X_t] < \mathbb{E}[Z_t | X_t, A = a]$. Similarly, if $P(Z_0)$ strictly MLRPs $P(Z_0 | A = a)$ for some group a, then $\mathbb{E}[Z_t | X_t] > \mathbb{E}[Z_t | X_t, A = a]$.

249 We prove this by showing that $P(Z_0 | X_t, A = a)$ strictly MLRPs $P(Z_0 | X_t)$, which implies that 250 $\mathbb{E}[Z_t | X_t, A = a] > \mathbb{E}[Z_t | X_t]$; §B.1 provides a full proof.

Theorem 4.4. A model that does not take into account disparities in rate of progression R will underestimate the disease severity of groups with higher progression rates and overestimate that of groups with lower progression rates. Specifically, if $P(R \mid A = a)$ strictly MLRPs P(R) for some group a, then $\mathbb{E}[Z_t \mid X_t] < \mathbb{E}[Z_t \mid X_t, A = a]$. Similarly, if P(R) strictly MLRPs $P(R \mid A = a)$ for some group a, then $\mathbb{E}[Z_t \mid X_t] > \mathbb{E}[Z_t \mid X_t, A = a]$.

We use a similar proof technique as for Theorem 4.3 and provide a full proof in §B.2.

Theorem 4.5. A model that does not take into account disparities in visit frequency λ_t (conditional on disease severity) will underestimate the disease severity of groups with lower visit frequency and overestimate that of groups with higher visit frequency. Specifically, if it holds for some group a that $\beta_A^{(a)} < \beta_A^{(\bar{a})}$ for all $\tilde{a} \neq a$, then $\mathbb{E}[Z_t \mid D_t] < \mathbb{E}[Z_t \mid D_t, A = a]$. Similarly, if it holds for some group a that $\beta_A^{(a)} > \beta_A^{(\bar{a})}$ for all $\tilde{a} \neq a$, then $\mathbb{E}[Z_t \mid D_t] > \mathbb{E}[Z_t \mid D_t, A = a]$.

Since group-specific differences in visit rate at a given severity are captured directly by the β_A parameter, we reason about disparities by comparing these parameters by group. We prove the theorem by directly reasoning about the estimates of Z_t when considering the additional term β_A versus not, reasoning in the large-sample limit in which λ_t can be perfectly estimated from the observed data D_t . In §5 we show empirically that our results hold in finite samples as well. Overall, these results convey the importance of accounting for disparities in disease progression models: it is fundamentally not possible to make well-calibrated estimates of severity without accounting for group differences in initial severity, progression rate, and visit frequency.

5 SYNTHETIC EXPERIMENTS

271 272 273

274

275

In this section, we validate our model and theoretical results in synthetic data simulations. We generate synthetic datasets according to the modeling assumptions in §3 (with parameter values for each dataset drawn randomly from each parameter's prior distribution). For each dataset, we generate simulated data for two separate groups, differing in initial severity, progression rate, and visit frequency (characterized by different μ_{Z_0} , μ_B , and β_A , respectively).

276 277 278

279

5.1 IDENTIFIABILITY AND SEVERITY ESTIMATION

We first verify Theorem 4.1 in simulations, showing that when we fit our model on synthetic data, 280 it can accurately recover the true data-generating parameters. We do this by examining the concor-281 dance between the model's estimated parameters and the true, latent parameter values, a common 282 approach in past work (Chang et al., 2021; Pierson et al., 2019). We find high correlation between 283 the true parameters and our model's posterior mean estimates (mean Pearson's r 0.996 across all 284 parameters; median 0.998), and good calibration (mean linear regression slope 1.00; median 1.00 285 when fit without an intercept term). We provide scatterplots of the true and estimated parameters 286 in Appendix C. We also see that our model's mean severity estimates for each group are highly 287 correlated and well-calibrated with ground truth, despite underlying differences in group severity 288 distributions and visit rates (Figure 2).

289 290

291

5.2 BIAS IN MODELS THAT DO NOT ACCOUNT FOR DISPARITIES

292 We now demonstrate in simulation that failing to account 293 for disparities can lead to biased severity estimates, consistent with Theorems 4.3, 4.4, and 4.5. In each trial, we use the same data to fit four models: our full model, which 295 accounts for all disparities, plus three ablated models that 296 each fail to account for one of the disparities (initial sever-297 ity, progression rate, visit frequency). To characterize the 298 resulting bias of failing to account for each type of dis-299 parity, we compute the average error in severity estimates 300 (mean inferred estimate minus mean true severity) of each 301 model, broken down by group. For each ablated model and 302 trial, we define the "underserved group" to be the one that 303 is underserved with respect to the specific disparity that the 304 model fails to capture. When evaluating our full model, we define the "underserved group" to be the one with higher 305 initial severity. 306

307 As seen in Table 2, the models that do not account for 308 disparities produce biased estimates: while our full model 309 achieves average error across all trials -0.02 and 0 for un-310 derserved and other patient groups respectively, the ablated models all have negative error for underserved patients (un-311 derestimated severity) and positive error for other patients 312 (overestimated severity). The ablated models also produce 313 severity estimates that are less correlated with true severity. 314



Figure 2: Well-calibrated severity estimates. Each dot shows the mean true vs. mean recovered severity values for one group in a given simulation trial. Members of groups depicted in red tend to be underserved compared to groups depicted in blue. Our full model produces accurate and well-calibrated severity estimates (estimates lie near dotted y = x line).

- 315
- 316 317

6 MODELING HEALTH DISPARITIES IN HEART FAILURE PROGRESSION

We fit our model on a real-world dataset of heart failure patients in the New York-Presbyterian hospital system. Heart failure is a progressive disease that affects many people, requires both specialty and preventive care (Colucci et al., 2020), and has known health disparities (Lewsey & Breathett, 2021), making it a natural application setting for our model. In §6.1 we summarize the dataset, and in §6.2 we confirm that our model can learn meaningful low-dimensional representations of disease severity by evaluating its reconstruction and predictive performance compared to standard baselines. In §6.3 we present our main results: we interpret our model's learned parameters to provide precise

		Model that fails to account for disparities in.			
	Full model	Initial severity	Progression rate	Visit frequency	
Underserved group bias	-0.02	-0.78	-0.24	-0.88	
Non-underserved group bias	0	+1.03	+0.01	+0.42	
Underserved group correlation	0.98	0.72	0.93	0.94	
Non-underserved group correlation	0.99	0.69	0.94	0.93	

Table 2: Failing to account for disparities produces biased estimates of severity Z_t . We compare severity estimates from our full model to three ablated models that each fail to account for one of the three health disparities. While our full model produces accurate, well-calibrated severity estimates, each ablated model underestimates severity for the underserved group and overestimates it for the other group. The ablated model estimates are also *less correlated* with the true severity values.

336 337

332

333

334

335

338

341

descriptions of health disparities in our setting, and we show that (as our theory predicts) failing to
 account for these disparities meaningfully shifts severity estimates.

342 6.1 DATA

Our data comes from the New York-Presbyterian (NYP)/Weill Cornell Medical Center's electronic health record (EHR) system from 2012 - 2020. We analyze a cohort of N = 2,942 patients who (1) have a specific subtype of heart failure (heart failure with reduced ejection fraction), to ensure our cohort can be described by a single progression model, and (2) are likely to receive most of their cardiology care in the NYP system, to ensure we can reasonably estimate when they receive care.

Observed feature data X_t for each patient includes four types of measurements: left ventricle ejec-349 tion fraction (LVEF), brain natriuretic peptide (BNP), systolic blood pressure (SBP), and heart rate 350 (HR). LVEF and BNP have strong clinical associations with heart failure severity (in terms of both 351 underlying physiological health and observed symptoms) (Murphy et al., 2020). SBP and HR are 352 less informative (more prone to fluctuation and changes not related to heart failure), but they are still 353 expected to show general trends over time as a patient's heart failure progresses. Since we must pin 354 the sign of at least one scaling factor F for identifiability, and decreasing LVEF is strongly asso-355 ciated with increasing severity in the heart failure subtype we study, we pin the sign of the scaling 356 factor between severity and LVEF values ($F_{\text{LVEF}} < 0$).

We discretize time into 1-week bins and observe timesteps when patients receive care. We then analyze disparities across four self-reported race/ethnicity groups: White non-Hispanic patients, Black non-Hispanic patients, Hispanic patients, and Asian non-Hispanic patients (which we will hereby describe as White, Black, Hispanic, and Asian subgroups). A full description of our data processing can be found in Appendix D.

363
364
6.2 MODEL VALIDATION

We first confirm that our model accurately fits the data: we verify that the model's inferred parameters are consistent with medical knowledge (§6.2.1) and compare the model's reconstruction and predictive performance to standard baselines (§6.2.2). Having confirmed this, we then show in §6.3, as our primary result, that our model provides insight into disparities in disease progression.

- 369 370
- 6.2.1 CONSISTENCY WITH MEDICAL KNOWLEDGE

Figure 3 plots our model's inferred parameters, all of which are consistent with existing medical knowledge.¹ Specifically, (1) the model correctly learns that BNP and HR tend to increase with heart failure severity ($F_{\text{BNP}}, F_{\text{HR}} > 0$), while SBP tends to decrease ($F_{\text{SBP}} < 0$) (Murphy et al., 2020); (2) the model learns larger variance parameters for SBP and HR values (ψ), correctly inferring that

 ¹For succinctness, Figure 3 plots only the model parameters of primary interest for interpreting our model
 (omitting, for example, estimated intercepts for each feature); a similar coefficient plot with all learned parameters is shown in Figure S2.



Figure 3: **Inferred model parameters with 95% confidence intervals.** Shared parameters (left) are consistent with medical knowledge of heart failure progression. Group-specific parameters (right) are plotted as differences compared to White patients, so confidence intervals that are non-overlapping with 0 (colored in purple) indicate significant racial/ethnic differences in parameters.

these features are less informative about heart failure progression than are BNP and LVEF (Murphy et al., 2020); and (3) the model estimates that β_Z is positive, meaning it learns that patients with higher disease severity tend to see healthcare providers more frequently, as expected.

6.2.2 **Reconstruction and predictive performance**

We next evaluate the model's ability to reconstruct and predict patient features X_t . Because the 402 model represents each patient visit in terms of a scalar severity Z_t , we do not expect the model 403 to perfectly reconstruct the multi-dimensional X_t ; rather, we hope for predictions that correlate 404 significantly with X_t . Consistent with this, when fit on 3 years of data per patient, our model's 405 predicted feature values correlate with true values both in- and out-of-sample. As we would hope, 406 the model best represents the features that are most informative for heart failure progression—LVEF 407 (r = 0.81 in-sample, r = 0.51 out-of-sample) and BNP (r = 0.62 in-sample, r = 0.31 out-of-)408 sample)—as opposed to the less-informative features SBP (r = 0.42 in-sample, r = 0.24 out-of-409 sample) and HR (r = 0.17 in-sample, r = 0.03 out-of-sample; all p-values besides HR out-of-410 sample < 0.001).

411 To provide a more detailed assessment of performance, we evaluate our model's ability to recon-412 struct features X_t in-sample and predict X_t out-of-sample, in comparison to seven standard base-413 lines. All of the baselines are designed to reconstruct or predict observed feature values (X_t) , as 414 opposed to additionally predicting whether patient visits will occur (D_t) . Our model can predict the 415 latter as well, but in order to provide a direct comparison of reconstruction and predictive perfor-416 mance, we compare only the feature prediction aspect of our model (so we do not fit any models 417 using D_t data) in this subsection. In the main text we report mean absolute percentage error (MAPE) 418 of estimated feature values because it allows us to report a normalized measure of error across multiple feature values; in Appendix E we additionally report RMSE. 419

420

391

392

393

394 395 396

397

398

399 400

401

421 **Reconstruction performance.** We compare our model's reconstruction performance to that of 422 two standard *dimensionality reduction baselines*: principal component analysis (PCA) and factor 423 analysis (FA). We compare our model to two variants of each. First, we compare our model to PCA 424 and FA fit at the *visit level*: one component per patient visit, analogous to our model's Z_t . Second, 425 we compare our model to PCA and FA fit at the *patient level*: two components for each patient, to 426 capture the trajectory of feature values as we do with Z_0 and R. We describe the implementation of 427 these baselines with more detail in Appendix E.

Because both PCA and FA require input vectors of consistent size, all models are fit on feature values
from the first three visits per patient. In Table 3, we report MAPE values averaged across all features
as well as across just the more informative features for heart failure severity: LVEF and BNP. We
achieve equivalent or better reconstruction performance across all features, and we reconstruct the
more informative features more accurately than any of the baselines.

432 433		Our model	FAvisit	PCA _{visit}	FA _{patient}	PCA _{patient}
434	MAPE: informative	20%	28%	23%	25%	21%
435	MAPE: all	16%	19%	17%	18%	16%
436						

Table 3: Our model compared to standard baselines for reconstruction performance. We compare to factor analysis and principal component analysis fit at the patient visit level (FAvisit, PCAvisit) and at the trajectory level (FA_{patient}, PCA_{patient}). Models are fit on the first 3 visits from each patient and evaluated on same data using mean absolute percentage error (MAPE).

	Our model	Linear regression	Quadratic regression	Latest timestep
MAPE: informative	28%	39%	59%	22%
MAPE: all	21%	32%	49%	18%

Table 4: Our model compared to standard baselines for predictive performance. We compare to linear regression, quadratic regression, and latest timestep prediction, each fit at the patient feature level. Models are fit on data from the first 3 years of each patient's disease trajectory and evaluated on visits after 3 years using mean absolute percentage error (MAPE).

449 450 451

447

448

437

438

439

Predictive Performance. We also compare our model's predictive performance to that of three 452 standard *timeseries forecasting baselines*: (1) a linear regression for each patient and feature; (2) a 453 quadratic regression for each patient and feature; and (3) predicting values equal to those at the last 454 timestep in training data. For this comparison, all models are fit on feature values from the first three 455 years of data per patient, and we evaluate predictive performance on all remaining visits. As seen in 456 Table 4, our model outperforms both linear regression and quadratic regression on all features. Our 457 model has slightly higher error than latest timestep, which is a widely-used, strong baseline for pure 458 predictive performance (Hyndman, 2018); latest timestep does not, however, provide any insight 459 into disparities or even patterns of progression over time.

460 Overall, while predicting and reconstructing X_t is not the primary goal of our model, it performs 461 generally well relative to standard baselines, validating its ability to accurately represent the data.

462 463 464

ANALYSIS OF DISPARITIES 6.3

465 We now discuss three main findings from fitting our model on the heart failure data. We learn that 466 (1) Black and Asian patients tend to have higher disease severity than White patients; (2) our model 467 learns precise descriptions of health disparities and finds that disparities of multiple types exist in our 468 setting; and (3) failing to account for the existing disparities meaningfully shifts severity estimates for all racial/ethnic groups. This analysis is descriptive and does not require evaluating held-out 469 performance, so models are fit on all available data. 470

471

Black and Asian patients have higher disease severity. In Figure 4, we compare mean severity 472 estimates for each group to the overall mean severity. Our model infers that Black and Asian patients 473 have significantly higher disease severities than White patients (p < 0.05, computed by cluster 474 bootstrapping at the patient-level). 475

476 **Model parameters capture fine-grained disparities.** As seen in Figure 3 (right), our model infers 477 that Black and Asian patients first visit healthcare providers for heart failure significantly later in 478 their disease progression than do White patients (inferred average initial severity μ_{Z_0} for Black 479 and Asian patient groups is greater than for White patients by 0.22 and 0.27, respectively). To 480 contextualize the magnitude of these disparities, if all patients progressed at the average learned 481 progression rate across the entire population, Black patients' first heart failure visit would occur 482 3.0 years later in the course of their disease progression than White patients', and Asian patients' 483 first visit would occur 3.8 years later. We also observe that β_A for Black patients is significantly lower than that of White patients, indicating that Black patients visit healthcare providers 10% less 484 frequently than White patients with the same disease severity. We describe these calculations in 485 Appendix F.

486 487

496

497

498

499

500

501 502



Figure 4: Accounting for disparities leads to less biased severity estimates. We compare the improvement of our full model (blue) over one that does not account for disparities but is otherwise the same (yellow) in two ways. On the left, we show each group's average difference from the overall mean severity, normalized by the overall standard deviation of severity. On the right, we capture the portion of each group that is identified as "high-risk" (top quartile of disease severity).

503 Accounting for disparities increases estimated severity for non-white patient groups. To as-504 sess whether accounting for disparities meaningfully shifts severity estimates, we compare severity 505 estimates from our model to those of an ablated version of our model that does not account for dis-506 parities (but is otherwise identical). This meaningfully shifts severity estimates (Figure 4 left): while 507 both models learn that non-white patients tend to have higher severity values, the ablated model produces higher severity estimates for White patients and lower severity estimates for all other groups 508 (p < 0.001 for all groups, computed by cluster bootstrapping at the patient-level). This is consistent 509 with our theoretical results. 510

511 To highlight some implications of these shifted severity estimates, we look at each model's ranking 512 of patient severity levels and profile of "high-risk" patient visits: visits where inferred severity lies 513 in the top quartile (25%) of all visits. The ablated model is significantly less likely to rank Black patient visits as high risk (Figure 4 right; p < 0.001, computed by cluster bootstrapping at the 514 patient-level), skewing the demographics of the high-risk patient cohort away from groups that we 515 know to have higher disease severity. 516

517 518

519 520

7 DISCUSSION

In this paper, we formalize three specific axes along which healthcare disparities emerge as biases 521 in observed health data: underserved patients may (1) first receive care only when their disease is 522 more severe, (2) progress faster even while receiving care, or (3) receive care less frequently even 523 at the same disease severity. We develop a disease progression modeling approach to interpretably 524 capture all three types of disparities while provably retaining identifiability. We prove that failing 525 to account for any of these disparities leads to biased estimates of severity and show in a real-world 526 heart failure dataset that accounting for health disparities does indeed meaningfully shift severity 527 estimates by increasing the proportion of non-white patients identified as high-risk. By evaluating 528 our model in a real healthcare setting, we validate its ability to learn fine-grained descriptions of 529 health disparities and to make disease severity estimates that are accurate across diverse populations 530 of patients. We thus urge future work in disease progression modeling to account for disparities in healthcare, and we lay a foundation for doing so. 531

532 There are several natural directions for future work. First, beyond heart failure, our approach could 533 be applied to the many other progressive diseases, including Parkinson's disease (Post et al., 2005), 534 Alzheimer's disease (Holford & Peace, 1992), diabetes (Perveen et al., 2020), and cancer (Gupta & Bar-Joseph, 2008). Second, an interesting technical direction is to extend our model to capture ad-536 ditional data modalities (e.g., medical images) or more flexible progression models (e.g., non-linear 537 trajectories) while retaining its provable identifiability. Finally, our approach generalizes naturally to progression model settings beyond healthcare where disparities are of interest, including infras-538 tructure deterioration (Madanat et al., 1995) and human aging (Pierson et al., 2019); these would be interesting domains for future work.

540 REFERENCES

548

554

555

565

566

567

568

569

570

573

- Ahmed M Alaa and Scott Hu. Learning from Clinical Judgments: Semi-Markov-Modulated Marked
 Hawkes Processes for Risk Prognosis. 2017.
- Ahmed M. Alaa and Mihaela van der Schaar. Attentive State-Space Modeling of Disease Progression. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d' Alché-Buc, E. Fox, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- Shamena Anwar and Hanming Fang. An Alternative Test of Racial Prejudice in Motor Vehicle
 Searches: Theory and Evidence. *American Economic Review*, 96(1):127–151, February 2006.
- Sidhika Balachandar, Nikhil Garg, and Emma Pierson. Domain constraints improve risk prediction when outcome data is missing. 2023.
 - R. Bellman and K.J. Åström. On structural identifiability. *Mathematical Biosciences*, 7(3-4):329– 339, April 1970.
- ⁵⁵⁶ Michael Betancourt. A Conceptual Introduction to Hamiltonian Monte Carlo, July 2018.
- Dwayne T. Brandon, Lydia A. Isaac, and Thomas A. LaVeist. The legacy of Tuskegee and trust in
 medical care: is Tuskegee responsible for race differences in mistrust of medical care? *Journal of the National Medical Association*, 97(7):951–956, July 2005.
- 561
 562
 563
 564
 564
 565
 564
 565
 564
 565
 564
 564
 564
 564
 565
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 564
 - Leighton Chan, L. Gary Hart, and David C. Goodman. Geographic Access to Health Care for Rural Medicare Beneficiaries. *The Journal of Rural Health*, 22(2):140–146, April 2006.
 - Serina Chang, Emma Pierson, Pang Wei Koh, Jaline Gerardin, Beth Redbird, David Grusky, and Jure Leskovec. Mobility network models of COVID-19 explain inequities and inform reopening. *Nature*, 589(7840):82–87, January 2021.
- Gilles Chemla and Christopher A. Hennessy. Controls, belief updating, and bias in medical rcts.
 Journal of Economic Theory, 184:104929, 2019.
- Irene Y. Chen, Peter Szolovits, and Marzyeh Ghassemi. Can AI Help Reduce Disparities in General
 Medical and Mental Health Care? *AMA Journal of Ethics*, 21(2):E167–179, February 2019.
- Irene Y. Chen, Rahul G. Krishnan, and David Sontag. Clustering Interval-Censored Time-Series for
 Disease Phenotyping, December 2021.
- Edward Choi, Mohammad Taha Bahadori, Joshua A. Kulas, Andy Schuetz, Walter F. Stewart, and
 Jimeng Sun. RETAIN: An Interpretable Predictive Model for Healthcare using Reverse Time
 Attention Mechanism. 2016a.
- Edward Choi, Mohammad Taha Bahadori, Andy Schuetz, Walter F. Stewart, and Jimeng Sun. Doctor
 AI: Predicting Clinical Events via Recurrent Neural Networks. *JMLR workshop and conference proceedings*, 56:301–318, August 2016b.
- Wilson S Colucci, SS Gottlieb, and SB Yeon. Overview of the management of heart failure with
 reduced ejection fraction in adults. U: UpToDate, Gottlieb SS ed. UpToDate [Internet]. Waltham,
 MA: UpToDate, 2020.
- Milton S. Davis. Physiologic, Psychological and Demographic Factors in Patient Compliance with Doctors' Orders. *Medical Care*, 6(2):115–122, 1968.
- Clarissa Jonas Diamantidis, Lindsay Zepel, Virginia Wang, Valerie A. Smith, Sarah Hudson Scholle,
 Loida Tamayo, and Matthew L. Maciejewski. Disparities in Chronic Kidney Disease Progression
 by Medicare Advantage Enrollees. *American Journal of Nephrology*, 52(12):949–957, 2021.

594 Johann D. Gaebler and Sharad Goel. A Simple, Statistically Robust Test of Discrimination, July 595 2024. arXiv:2407.06539 [stat]. 596 A. Gupta and Z. Bar-Joseph. Extracting Dynamics from Static Cancer Expression Data. *IEEE/ACM* 597 Transactions on Computational Biology and Bioinformatics, 5(2):172–182, April 2008. 598 Stefanie Hendricks, Iryna Dykun, Bastian Balcer, Matthias Totzeck, Tienush Rassaf, and Amir A. 600 Mahabadi. Higher BNP/NT-pro BNP levels stratify prognosis equally well in patients with and 601 without heart failure: a meta-analysis. ESC Heart Failure, 9(5):3198-3209, October 2022. 602 N H Holford and K E Peace. Methodologic aspects of a population pharmacodynamic model for 603 cognitive effects in Alzheimer patients treated with tacrine. Proceedings of the National Academy 604 of Sciences, 89(23):11466–11470, December 1992. 605 Xiao Hu, John W Melson, Stacey S Pan, Yana V Salei, and Yu Cao. Screening, Diagnosis, and 607 Initial Care of Asian and White Patients With Lung Cancer. The Oncologist, 29(4):332-341, 608 April 2024. 609 RJ Hyndman. Forecasting: principles and practice. OTexts, 2018. 610 611 Javaid Iqbal, Ophira Ginsburg, Paula A. Rochon, Ping Sun, and Steven A. Narod. Differences 612 in Breast Cancer Stage at Diagnosis and Cancer-Specific Survival by Race and Ethnicity in the 613 United States. JAMA, 313(2):165, January 2015. 614 Christopher H. Jackson, Linda D. Sharples, Simon G. Thompson, Stephen W. Duffy, and Elisabeth 615 Couto. Multistate Markov models for disease progression with classification error. Journal of the 616 Royal Statistical Society: Series D (The Statistician), 52(2):193–209, July 2003. 617 618 Samuel Karlin and Herman Rubin. The theory of decision procedures for distributions with mono-619 tone likelihood ratio. The Annals of Mathematical Statistics, 27(2):272–299, 1956. 620 Ben Klemens. When Do Ordered Prior Distributions Induce Ordered Posterior Distributions? SSRN 621 Electronic Journal, 2007. 622 623 Petr Klemera and Stanislav Doubal. A new approach to the concept and computation of biological 624 age. Mechanisms of Ageing and Development, 127(3):240–248, March 2006. 625 Takeshi Kurashima, Tim Althoff, and Jure Leskovec. Modeling interdependent and periodic real-626 world action sequences. In Proceedings of the 2018 world wide web conference, pp. 803–812, 627 2018. 628 629 Bum Chul Kwon, Min-Je Choi, Joanne Taery Kim, Edward Choi, Young Bin Kim, Soonwook Kwon, 630 Jimeng Sun, and Jaegul Choo. RetainVis: Visual Analytics with Interpretable and Interactive 631 Recurrent Neural Networks on Electronic Medical Records. IEEE Transactions on Visualization 632 and Computer Graphics, 25(1):299-309, January 2019. 633 Bruce E. Landon, Jukka-Pekka Onnela, Laurie Meneades, A. James O'Malley, and Nancy L. Keat-634 ing. Assessment of Racial Disparities in Primary Care Physician Specialty Referrals. JAMA 635 Network Open, 4(1):e2029238, January 2021. 636 637 Thomas A. LaVeist, Lydia A. Isaac, and Karen Patricia Williams. Mistrust of Health Care Organi-638 zations Is Associated with Underutilization of Health Services. Health Services Research, 44(6): 639 2093–2105, December 2009. 640 Richard J. Lee, Ravi A. Madan, Jayoung Kim, Edwin M. Posadas, and Evan Y. Yu. Disparities in 641 Cancer Care and the Asian American Population. The Oncologist, 26(6):453-460, June 2021. 642 643 M. E. Levine. Modeling the Rate of Senescence: Can Estimated Biological Age Predict Mortality 644 More Accurately Than Chronological Age? The Journals of Gerontology Series A: Biological 645 Sciences and Medical Sciences, 68(6):667–674, June 2013. 646 Sabra C. Lewsey and Khadijah Breathett. Racial and ethnic disparities in heart failure: current state 647 and future directions. Current Opinion in Cardiology, 36(3):320-328, May 2021.

648 649	Bryan Lim and Mihaela van der Schaar. Disease-Atlas: Navigating Disease Trajectories with Deep Learning, July 2018.
651 652	Zachary C. Lipton, David C. Kale, Charles Elkan, and Randall Wetzel. Learning to Diagnose with LSTM Recurrent Neural Networks, March 2017.
653 654 655	Yu-Ying Liu, Shuang Li, Fuxin Li, Le Song, and James M. Rehg. Efficient Learning of Continuous- Time Hidden Markov Models for Disease Progression. Advances in Neural Information Process- ing Systems, 28:3599–3607, 2015.
656 657 658	Zhi Liu, Uma Bhandaram, and Nikhil Garg. Quantifying spatial under-reporting disparities in resident crowdsourcing. <i>Nature Computational Science</i> , 4(1):57–65, 2024.
659 660	Fenglong Ma, Radha Chitta, Jing Zhou, Quanzeng You, Tong Sun, and Jing Gao. Dipole: Diagnosis Prediction in Healthcare via Attention-based Bidirectional Recurrent Neural Networks. 2017.
661 662 663 664	Samer Madanat, Rabi Mishalani, and Wan Hashim Wan Ibrahim. Estimation of Infrastructure Tran- sition Probabilities from Condition Rating Data. <i>Journal of Infrastructure Systems</i> , 1(2):120–125, June 1995.
665 666	Sarah Miller and Laura R. Wherry. Health and Access to Care during the First 2 Years of the ACA Medicaid Expansions. <i>New England Journal of Medicine</i> , 376(10):947–956, March 2017.
667 668 669	D R Mould, N G Denman, and S Duffull. Using Disease Progression Models as a Tool to Detect Drug Effect. <i>Clinical Pharmacology & Therapeutics</i> , 82(1):81–86, July 2007.
670 671	Sean P. Murphy, Nasrien E. Ibrahim, and James L. Januzzi. Heart Failure With Reduced Ejection Fraction: A Review. <i>JAMA</i> , 324(5):488, August 2020.
672 673 674 675	Sarah Nouri, Courtney R. Lyles, Elizabeth B. Sherwin, Magdalene Kuznia, Anna D. Rubinsky, Kathryn E. Kemper, Oanh K. Nguyen, Urmimala Sarkar, Dean Schillinger, and Elaine C. Khoong. Visit and Between-Visit Interaction Frequency Before and After COVID-19 Telehealth Implemen- tation. JAMA Network Open, 6(9):e2333944, September 2023.
676 677 678 679	Sajida Perveen, Muhammad Shahbaz, Muhammad Sajjad Ansari, Karim Keshavjee, and Aziz Guer- gachi. A Hybrid Approach for Modeling Type 2 Diabetes Mellitus Progression. <i>Frontiers in</i> <i>Genetics</i> , 10:1076, January 2020.
680 681 682	Emma Pierson, Pang Wei Koh, Tatsunori Hashimoto, Daphne Koller, Jure Leskovec, Nicholas Eriks- son, and Percy Liang. Inferring Multidimensional Rates of Aging from Cross-Sectional Data, March 2019.
683 684 685 686	Teun M. Post, Jan I. Freijer, Joost DeJongh, and Meindert Danhof. Disease System Analysis: Basic Disease Progression Models in Degenerative Disease. <i>Pharmaceutical Research</i> , 22(7):1038–1049, July 2005.
687 688	Megan Reilly. Health Disparities and Access to Healthcare in Rural vs. Urban Areas. <i>Theory in Action</i> , 14(2):6–27, April 2021.
689 690 691 692 693	K Romero, K Ito, Ja Rogers, D Polhamus, R Qiu, D Stephenson, R Mohs, R Lalonde, V Sinha, Y Wang, D Brown, M Isaac, S Vamvakas, R Hemmings, L Pani, Lj Bain, B Corrigan, and Alzheimer's Disease Neuroimaging Initiative* for the Coalition Against Major Diseases**. The future is now: Model-based clinical trial design for Alzheimer's disease. <i>Clinical Pharmacology</i> & <i>Therapeutics</i> , 97(3):210–214, March 2015.
695 696 697 698 699	Rasmus Rørth, Pardeep S. Jhund, Mehmet B. Yilmaz, Søren Lund Kristensen, Paul Welsh, Ak- shay S. Desai, Lars Køber, Margaret F. Prescott, Jean L. Rouleau, Scott D. Solomon, Karl Swed- berg, Michael R. Zile, Milton Packer, and John J.V. McMurray. Comparison of bnp and nt-probnp in patients with heart failure and reduced ejection fraction. <i>Circulation: Heart Failure</i> , 13(2): e006541, 2020.
700 701	Gráinne Schäfer, Kenneth M. Prkachin, Kimberley A. Kaseweter, and Amanda C. De C Williams. Health care providers' judgments in chronic pain: the influence of gender and trustworthiness. <i>Pain</i> , 157(8):1618–1625, August 2016.

702 703 704	Divya Shanmugam, Kaihua Hou, and Emma Pierson. Quantifying disparities in intimate partner violence: a machine learning method to correct for underreporting. 2021.
705 706	Alexander Shapiro. Identifiability of factor analysis: Some results and open problems. <i>Linear Algebra and its Applications</i> , 70:1–7, 1985.
707 708 709 710	Jonathan Suarez, Jordana B. Cohen, Vishnu Potluri, Wei Yang, David E. Kaplan, Marina Serper, Siddharth P. Shah, and Peter Philip Reese. Racial Disparities in Nephrology Consultation and Disease Progression among Veterans with CKD: An Observational Cohort Study. <i>Journal of the American Society of Nephrology</i> , 29(10):2563–2573, October 2018.
711 712 713 714	R. Sukkar, E. Katz, Yanwei Zhang, D. Raunig, and B. T. Wyman. Disease progression modeling using Hidden Markov Models. In 2012 Annual International Conference of the IEEE Engineering in Medicine and Biology Society, pp. 2845–2848, San Diego, CA, August 2012. IEEE.
715 716 717 718	Arjun K. Venkatesh, Shih-Chuan Chou, Shu-Xia Li, Jennie Choi, Joseph S. Ross, Gail D'Onofrio, Harlan M. Krumholz, and Kumar Dharmarajan. Association Between Insurance Status and Ac- cess to Hospital Care in Emergency Department Disposition. JAMA Internal Medicine, 179(5): 686, May 2019.
719 720 721 722	Xiang Wang, David Sontag, and Fei Wang. Unsupervised learning of disease progression models. In <i>Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining</i> , KDD 2014, pp. 85–94, New York New York USA, August 2014. ACM. ISBN 978-1-4503-2956-9.
723 724 725 726	Kathryn E. Weaver, Julia H. Rowland, Keith M. Bellizzi, and Noreen M. Aziz. Forgoing medical care because of cost: Assessing disparities in healthcare access among cancer survivors living in the United States. <i>Cancer</i> , 116(14):3493–3504, July 2010.
727 728 729	Ruqaiijah Yearby. Racial Disparities in Health Status and Access to Healthcare: The Continuation of Inequality in the United States Due to Structural Racism. <i>The American Journal of Economics and Sociology</i> , 77(3-4):1113–1152, May 2018.
730 731 732	Haoran Zhang, Amy X. Lu, Mohamed Abdalla, Matthew McDermott, and Marzyeh Ghassemi. Hurt- ful Words: Quantifying Biases in Clinical Contextual Word Embeddings. 2020.
733 734 735	Anna Zink, Ziad Obermeyer, and Emma Pierson. Race Corrections in Clinical Algorithms Can Help Correct for Racial Disparities in Data Quality, April 2023.
736 737	
738 739	
740 741	
742 743	
744	
746	
747 748	
749 750	
751	
752 753	
754 755	

756 A PROOF OF IDENTIFIABILITY

758 A.1 PROOF OF THEOREM 4.1

 Theorem 4.1. All model parameters are identified by the observed data distribution $P(X_t, D_t | A)$.

Proof. We want to show that each unique set of parameter assignments leads to a different distribution over the observed data. To do this, we divide our argument into four lemmas:

Lemma A.1. Parameters F, b, Ψ are identified by $P(X_t | A = a_0)$.

Proof. We want to show that if two parameter sets $\{F, b, \Psi\}$ and $\{\tilde{F}, \tilde{b}, \tilde{\Psi}\}$ yield the same observed data distribution $P(X_0 \mid A = a_0)$, the parameter sets must be identical.

We first note that at t = 0, we have $Z_t = Z_0 \sim \mathcal{N}(0, 1)$ for group a_0 . Then the mapping between severity and features

$$X_0 = F \cdot Z_0 + b + \epsilon_t$$
$$\epsilon_t \sim \mathcal{N}(0, \Psi)$$

captures a factor analysis model with factor loading matrix F and diagonal covariance matrix Ψ . At t = 0, the feature distribution for group a_0 has the standard factor analysis distribution (Shapiro, 1985):

$$X_0 \sim \mathcal{N}(b, FF^T + \Psi).$$

Assuming the two sets of parameters map to distributions of X_0 with the same mean, it must hold that $b = \tilde{b}$. Thus, parameter b is identified by data distribution $P(X_0 \mid A = a_0)$.

Further, the covariance matrix of X_0 induced by each set of parameters must be the same: $F(F)^T + \Psi = \tilde{F}(\tilde{F})^T + \tilde{\Psi}$. Element-wise equality of the covariance matrix gives us the following, where subscripts *i* refer to the *i*-th element of each parameter vector:

$$F_i F_j = F_i F_j \quad \forall i, j, i \neq j \tag{1}$$

$$(F_i)^2 + \Psi_i = (\tilde{F}_i)^2 + \tilde{\Psi}_i \tag{2}$$

Using the equality constraint (1) for multiple pairs of indices, we have that for all assignments of distinct indices i, j, k:

$$(F_i F_j = \tilde{F}_i \tilde{F}_j) \wedge (F_j F_k = \tilde{F}_j \tilde{F}_k) \implies \frac{\tilde{F}_i}{F_i} = \frac{\tilde{F}_k}{F_k}$$
(3)

$$F_i F_k = \tilde{F}_i \tilde{F}_k \implies \frac{F_i}{\tilde{F}_i} = \frac{F_k}{F_k}$$
(4)

Together, equations 3 and 4 give us:

$$\frac{\tilde{F}_i}{F_i} = \frac{F_i}{\tilde{F}_i} \implies (\tilde{F}_i)^2 = (F_i)^2 \implies F_i = \alpha \tilde{F}_i$$

where $\alpha \in \{-1, +1\}$. Since we have fixed $F_0 > 0$ for all factor loading matrices F, the sign of α is fixed:

$$F_0 = \alpha \tilde{F}_0 \implies \alpha = 1 \implies F_i = \tilde{F}_i \ \forall i \in [0, d), \tag{5}$$

meaning we have identified F.

 Lastly, using equations (2) and (5) we get $F_i = \tilde{F}_i \implies \Psi_i = \tilde{\Psi}_i$. We have now shown that if two parameter sets induce the same distribution of X at time t = 0, they must have the same exact value assignments. Therefore F, b, Ψ are identified by $P(X_t \mid A = a_0)$.

Lemma A.2. Global parameters F, b, Ψ and parameters $\mu_{Z_0}^{(a)}, \sigma_{Z_0}^{(a)}, \mu_R^{(a)}, \sigma_R^{(a)}$ for each group a are identified by $P(X_t \mid A)$.

Proof. By Lemma A.1, we know that F, b, Ψ are identified by $P(X_0 \mid A = a_0)$. We want to show that for any group a, if two parameter sets $\{\mu_{Z_0}^{(a)}, \sigma_{Z_0}^{(a)}, \mu_R^{(a)}, \sigma_R^{(a)}\}$ and $\{\tilde{\mu}_{Z_0}^{(a)}, \tilde{\sigma}_{Z_0}^{(a)}, \tilde{\mu}_R^{(a)}, \tilde{\sigma}_R^{(a)}\}$ yield the same observed data distribution $P(X_t \mid A = a)$, the parameter sets must be identical. In this proof we consider an arbitrary group a and omit the (a) superscript for brevity.

We model the following:

$$Z_{0} \sim \mathcal{N} \left(\mu_{Z_{0}}, \sigma_{Z_{0}}^{2}\right)$$

$$R \sim \mathcal{N} \left(\mu_{R}, \sigma_{R}^{2}\right)$$

$$Z_{t} = Z_{0} + R \cdot t \implies Z_{t} \sim \mathcal{N} \left(\mu_{R} \cdot t + \mu_{Z_{0}}, \sigma_{R}^{2} \cdot t^{2} + \sigma_{Z_{0}}^{2}\right)$$

$$X_{t} = F \cdot Z_{t} + b + \epsilon_{t}, \text{ where } \epsilon_{t} \sim \mathcal{N}(0, \Psi)$$
(6)

We see that equation (6) captures a factor analysis model with factor loading matrix F and diagonal covariance matrix Ψ , meaning

$$X_t \sim \mathcal{N}(b + F(\mu_R \cdot t + \mu_{Z_0}), F(\sigma_R^2 \cdot t^2 + \sigma_{Z_0}^2)F^T + \Psi).$$

Recalling that $F_0 > 0$, we first consider t = 0, where $X_0 \sim \mathcal{N}(b + F\mu_{Z_0}, F(\sigma_{Z_0}^2)F^T + \Psi)$. In order for the two parameter sets to map to distributions of X_0 with the same mean, it must be the case that

$$b + F\mu_{Z_0} = b + F\tilde{\mu}_{Z_0} \implies \mu_{Z_0} = \tilde{\mu}_{Z_0}$$

Further, for the two parameter sets to map to distributions with the same covariance matrix, it must hold that

$$F(\sigma_{Z_0}{}^2)F^T + \Psi = F(\tilde{\sigma}_{Z_0}{}^2)F^T + \Psi \implies \sigma_{Z_0} = \tilde{\sigma}_{Z_0}$$

since we know $\sigma_{Z_0}, \tilde{\sigma}_{Z_0} > 0$. So we have identified μ_{Z_0} and σ_{Z_0} . We next consider any time $t \neq 0$. For the two parameter sets to map to distributions of X_t with the same mean, given that we have already shown μ_{Z_0} must equal $\tilde{\mu}_{Z_0}$, it must hold that

$$b + F(\mu_R \cdot t + \mu_{Z_0}) = b + F(\tilde{\mu}_R \cdot t + \tilde{\mu}_{Z_0}) \implies \mu_R = \tilde{\mu}_R.$$

For the two parameter sets to map to distributions with the same covariance matrix, given that we have already shown σ_{Z_0} must equal $\tilde{\sigma}_{Z_0}$, it must hold that

$$F(\sigma_R^2 \cdot t^2 + \sigma_{Z_0}^2)F^T + \Psi = F(\tilde{\sigma}_R^2 \cdot t^2 + \tilde{\sigma}_{Z_0}^2)F^T + \Psi \implies \sigma_R = \tilde{\sigma}_R$$

since $\sigma_R, \tilde{\sigma}_R > 0$. Thus we have shown that for any group *a*, group-specific values of $\mu_{Z_0}, \sigma_{Z_0}, \mu_R, \sigma_R$ are identified by $P(X_t \mid A = a)$.

Lemma A.3. Global parameters β_0, β_Z and the parameter $\beta_A^{(a)}$ for each group *a* are identified by $P(D_t \mid A)$.

Proof. We want to show that if two parameter sets $\{\beta_0, \beta_Z, \beta_A^{(a)}\}\$ and $\{\tilde{\beta}_0, \tilde{\beta}_Z, \tilde{\beta}_A^{(a)}\}\$ yield the same observed data distribution $P(D_t \mid A = a)$, the parameter sets must be identical. Unless otherwise specified, we consider an arbitrary group a and omit the (a) superscript for brevity. We also assume $\mu_R \neq 0$, since in general the severity of a progressive disease should change over time and it does not make sense to learn progression in the case that it does not.

Each event when a patient visits the hospital $(D_t = 1)$ is generated by an inhomogeneous Poisson process parameterized by λ_t , where $\log(\lambda_t) = \beta_0 + \beta_Z \cdot Z_t + \beta_A$.

In order for two data distributions to have identical $P(D_t \mid A = a)$ they must have identical expected rates $\mathbb{E}_{Z_0,R}[\lambda_t]$: $\mathbb{E}_{Z_0,R}[\lambda_t]$ is the expected rate of events (across the population) at time t—if two distributions have a different expected rate of events at any time t, then $P(D_t \mid A = a_0)$ must differ at that point in time as well. Thus if two sets of parameters $\{\beta_0, \beta_Z, \beta_A\}$ and $\{\tilde{\beta}_0, \tilde{\beta}_Z, \tilde{\beta}_A\}$ yield the same observed data distribution $P(D_t \mid A = a)$, they must also generate the same observed values $\mathbb{E}_{Z_0,R}[\lambda_t]$ at all timesteps t. We finish the proof by showing that this holds only if $\{\beta_0, \beta_Z, \beta_A\} = \{\tilde{\beta}_0, \tilde{\beta}_Z, \tilde{\beta}_A\}$.

$$\mathbb{E}_{Z_0,R}[\lambda_t] = \int \int \lambda_t \cdot P(Z_0) \cdot P(R) \, dZ_0 dR$$

By Lemma A.2, we know that $\mu_{Z_0}, \sigma_{Z_0}, \mu_R, \sigma_R$ are identified by $P(X_t \mid A)$. Then

$$P(Z_0) = \frac{1}{\sqrt{2\pi(\sigma_{Z_0})^2}} \exp\left(-\frac{(Z_0 - \mu_{Z_0})^2}{2(\sigma_{Z_0})^2}\right)$$
$$P(R) = \frac{1}{\sqrt{2\pi(\sigma_R)^2}} \exp\left(-\frac{(R - \mu_R)^2}{2(\sigma_R)^2}\right)$$

$$\mathbb{E}_{Z_0,R}[\lambda_t] = \exp(f(\beta_0,\beta_Z,\beta_A,t)) \tag{7}$$

where $f(\beta_0,\beta_Z,\beta_A,t) = \left(\frac{(\beta_Z\sigma_R)^2}{2}\right)t^2 + (\beta_Z\mu_R)t + \left(\beta_0 + \frac{(\beta_Z\sigma_{Z_0})^2}{2} + \beta_Z\mu_{Z_0} + \beta_A\right)$

The expression in 7 must be equal for $\{\beta_0, \beta_Z, \beta_A\}$ and $\{\tilde{\beta}_0, \tilde{\beta}_Z, \tilde{\beta}_A\}$ at all timesteps t. Since exp is an injective function, this means that $f(\beta_0, \beta_Z, \beta_A, t) = f(\tilde{\beta}_0, \tilde{\beta}_Z, \tilde{\beta}_A, t)$ for all t. By equality of polynomials, each of the individual polynomial coefficients must be equal must be equal for this to hold.

We first consider the case for group a_0 , since we pin $\beta_A^{(a_0)}$ at 0 as a reference for all other groups. Given that we have already identified $\mu_{Z_0}^{(a_0)}, \sigma_{Z_0}^{(a_0)}, \mu_R^{(a_0)}, \sigma_R^{(a_0)}$,

$$\left(\beta_0 + \frac{(\beta_Z \sigma_{Z_0})^2}{2} + \beta_Z \mu_{Z_0}\right) = \left(\tilde{\beta}_0 + \frac{(\tilde{\beta}_Z \sigma_{Z_0})^2}{2} + \tilde{\beta}_Z \mu_{Z_0}\right) \implies \beta_0 = \tilde{\beta}_0$$

Now we return to our analysis of any arbitrary group a. Given that we have already identified $\mu_{Z_0}, \sigma_{Z_0}, \mu_R \neq 0, \sigma_R$,

 $\beta_Z \mu_R = \tilde{\beta}_Z \mu_R \implies \beta_Z = \tilde{\beta}_Z$

$$\left(\beta_0 + \frac{(\beta_Z \sigma_{Z_0})^2}{2} + \beta_Z \mu_{Z_0} + \beta_A\right) = \left(\tilde{\beta}_0 + \frac{(\tilde{\beta}_Z \sigma_{Z_0})^2}{2} + \tilde{\beta}_Z \mu_{Z_0} + \tilde{\beta}_A\right) \implies \beta_A = \tilde{\beta}_A$$

Thus we have shown that β_0, β_Z , and $\beta_A^{(a)}$ for any group a are identified by $P(D_t \mid Z_t, A)$.

By showing that each parameter of the model is uniquely recovered from the observed data, we have proved that our model is identifiable.

B PROOFS OF BIAS

921 922

923

924 925

926

927

928

929

930

931

932 933

934

940

941

953

954

955

956 957 958 In this section, in order to capture the effect of failing to account for one disparity at a time, we consider the setting where everything between two groups is the same except for disparity of focus. It is clear to see from our analysis that these results hold even more generally—as long as all existing disparities disfavor or favor the same group (e.g. a disadvantaged group with respect to one disparity is not advantaged with respect to another, in which case the effects could cancel each other out), our proofs of bias will hold. Throughout our proofs, we assume that all PDFs and conditional PDFs have positive support over their entire domain, and that all PDFs are differentiable, a very reasonable assumption over our setting.

B.1 THEOREM 4.3

Theorem 4.3. A model that does not take into account disparities in initial disease severity Z_0 will underestimate the disease severity of groups with higher initial severity and overestimate that of groups with lower initial severity. Specifically, if $P(Z_0 | A = a)$ strictly MLRPs $P(Z_0)$ for some group a, then $\mathbb{E}[Z_t | X_t] < \mathbb{E}[Z_t | X_t, A = a]$. Similarly, if $P(Z_0)$ strictly MLRPs $P(Z_0 | A = a)$ for some group a, then $\mathbb{E}[Z_t | X_t] > \mathbb{E}[Z_t | X_t, A = a]$.

Proof. We want to show that $\mathbb{E}[Z_t \mid X_t, A = a] > \mathbb{E}[Z_t \mid X_t]$. We first show that $P(Z_0 \mid X_t = x, A = a)$ strictly MLRPs $P(Z_0 \mid X_t)$ with respect to Z_0 :

$$\frac{\partial}{\partial Z_0} \left(\frac{P(Z_0 \mid X_t, A = a)}{P(Z_0 \mid X_t)} \right) = \frac{\partial}{\partial Z_0} \left(\frac{\frac{P(X_t \mid Z_0, A = a)P(Z_0 \mid A = a)}{P(X_t \mid A = a)}}{\frac{P(X_t \mid Z_0)P(Z_0)}{P(X_t)}} \right)$$
(Bayes Rule)

$$= \frac{\partial}{\partial Z_0} \left(\frac{\frac{P(Z_0 | A = a)}{P(X_t | A = a)}}{\frac{P(Z_0)}{P(X_t)}} \right) \qquad (X_t \perp A \mid Z_0, R)$$
$$= \frac{P(X_t)}{P(X_t \mid A = a)} \cdot \frac{\partial}{\partial Z_0} \left(\frac{P(Z_0 \mid A = a)}{P(Z_0)} \right)$$
$$> 0 \qquad (Disparity assumption)$$

Since MLRP implies first-order stochastic dominance (FOSD) (Klemens, 2007), this proves that $P(Z_0 \mid X_t, A = a)$ strictly FOSDs $P(Z_0 \mid X_t)$ and thus that $\mathbb{E}[Z_0 \mid X_t, A = a] > \mathbb{E}[Z_0 \mid X_t]$. By linearity of expectation,

$$\mathbb{E}[Z_0 \mid X_t, A = a] + \mathbb{E}[f(R, t) \mid X_t, A = a] > \mathbb{E}[Z_0 \mid X_t] + \mathbb{E}[f(R, t) \mid X_t], \quad \forall t \ge 0$$

$$\implies \mathbb{E}[Z_t \mid X_t, A = a] > \mathbb{E}[Z_t \mid X_t]$$

It is clear to see that this argument extends naturally to show that if a group tends to come in at earlier disease stages than the rest of the population, that their severity will be overestimated: If there exists a group \tilde{a} such that $P(Z_0)$ strictly MLRPs $P(Z_0 \mid A = \tilde{a})$ with respect to Z_0 and $\mathbb{E}[R \mid X_t] \ge \mathbb{E}[R \mid X_t, A = \tilde{a}]$, then we will see that $\mathbb{E}[Z_t \mid X_t, A = \tilde{a}] < \mathbb{E}[Z_t \mid X_t]$. Hence any model that does not take into account demographic disparities in initial disease severity levels at a patient's first visit will lead to biased estimates of severity.

966 B.2 PROOF OF THEOREM 4.4

Theorem 4.4. A model that does not take into account disparities in rate of progression R will underestimate the disease severity of groups with higher progression rates and overestimate that of groups with lower progression rates. Specifically, if $P(R \mid A = a)$ strictly MLRPs P(R) for some group a, then $\mathbb{E}[Z_t \mid X_t] < \mathbb{E}[Z_t \mid X_t, A = a]$. Similarly, if P(R) strictly MLRPs $P(R \mid A = a)$ for some group a, then $\mathbb{E}[Z_t \mid X_t] > \mathbb{E}[Z_t \mid X_t, A = a]$. 972 R is a patient's linear rate of progression, so we model a patient's severity over time as $Z_t = f(R,t) + Z_0$, where f is linearly increasing in R.

Proof. We want to show that $\mathbb{E}[Z_t \mid X_t, A = a] > \mathbb{E}[Z_t \mid X_t]$. We first show that $P(R \mid X_t, A = a)$ strictly MLRPs $P(R \mid X_t)$ with respect to R:

$$\frac{\partial}{\partial R} \left(\frac{P(R \mid X_t, A = a)}{P(R \mid X_t)} \right) = \frac{\partial}{\partial R} \left(\frac{\frac{P(X_t \mid R, A = a)P(R \mid A = a)}{P(X_t \mid A = a)}}{\frac{P(X_t \mid R)P(Z_t = z_t)}{P(X_t)}} \right)$$
(Bayes Rule)
$$= \frac{\partial}{\partial R} \left(\frac{\frac{P(R \mid A = a)}{P(X_t \mid A = a)}}{\frac{P(X_t \mid A = a)}{P(X_t)}} \right)$$
(X \pm A \| Z_0, R)
$$= \frac{P(X_t)}{P(X_t \mid A = a)} \cdot \frac{\partial}{\partial R} \left(\frac{P(R \mid A = a)}{P(R)} \right)$$
> 0 (Disparity assumption)

Since MLRP implies FOSD (Klemens, 2007), this also implies that $P(R \mid X_t, A = a)$ strictly FOSDs $P(R \mid X_t)$. It follows directly that $\mathbb{E}[R \mid X_t, A = a] > \mathbb{E}[R \mid X_t]$. By linearity of expectation,

988 989

975

976

1017

 $\mathbb{E}[f(R,t) + Z_0 \mid X_t, A = a] > \mathbb{E}[f(R,t) + Z_0 \mid X_t], \quad \forall t > 0$ $\implies \mathbb{E}[Z_t \mid X_t, A = a] > \mathbb{E}[Z_t \mid X_t]$

It is clear to see that this argument extends naturally to show that if a group tends to progress *more* slowly than the rest of the population, that their severity will be overestimated: if there exists a group \tilde{a} such that P(R) strictly MLRPs $P(R \mid A = \tilde{a})$ with respect to R and $\mathbb{E}[Z_0 \mid X_t] \ge \mathbb{E}[Z_0 \mid$ $X_t, A = \tilde{a}]$, then we will see that $\mathbb{E}[Z_t \mid X_t, A = \tilde{a}] < \mathbb{E}[Z_t \mid X_t]$. Thus any model that does not take into account demographic disparities in patient progression rates will lead to biased estimates of severity.

1000 1001 B.3 PROOF OF THEOREM 4.5

Theorem 4.5. A model that does not take into account disparities in visit frequency λ_t (conditional on disease severity) will underestimate the disease severity of groups with lower visit frequency and overestimate that of groups with higher visit frequency. Specifically, if it holds for some group a that $\beta_A^{(a)} < \beta_A^{(\tilde{a})}$ for all $\tilde{a} \neq a$, then $\mathbb{E}[Z_t \mid D_t] < \mathbb{E}[Z_t \mid D_t, A = a]$. Similarly, if it holds for some group a that $\beta_A^{(a)} > \beta_A^{(\tilde{a})}$ for all $\tilde{a} \neq a$, then $\mathbb{E}[Z_t \mid D_t] > \mathbb{E}[Z_t \mid D_t, A = a]$.

We model a patient's visit pattern using an inhomogeneous poisson process characterized by visit rate λ_t , such that $\log(\lambda_t) = g(Z_t) + \beta_A^{(A)}$ for some function of severity $g(Z_t)$ and group-specific adjustments $\beta_A^{(A)}$. In our proof, we assume the large-sample limit in which λ_t can be perfectly estimated from the observed data, and thus treat it as observed; we show empirically that our results hold in finite samples as well. We assume $g(Z_t)$ is a strictly monotonically increasing function of severity.

1015 *Proof.* We want to show that $\mathbb{E}[Z_t \mid D_t, A = a] > \mathbb{E}[Z_t \mid D_t]$. We do this by calculating each term separately.

We first consider $\mathbb{E}[Z_t \mid D_t, A = a]$. Observing D_t over time gives us an observed value of visit rate λ_t . The strictly monotone assumption of g ensures g is invertible, and the fact that all visit rates λ_t are characterized by $\log(\lambda_t) = g(Z_t) + \beta_A^{(A)}$ ensures that this holds over the entire range of λ_t values. This gives us:

1023
1024
1025
$$\mathbb{E}[Z_t \mid D_t, A = a] = \mathbb{E}\left[g^{-1}\left(\log(\lambda_t) - \beta_A^{(A)}\right) \mid D_t, A = a\right]$$
1026
$$= \frac{-1}{2}\left(\log(\lambda_t) - \beta_A^{(A)}\right)$$

$$= g^{-1} \left(\log(\lambda_t) - \beta_A^{(\alpha)} \right)$$

We next consider the case where a model infers severity without taking into account disparities in visit rate conditional on severity. Estimating severity Z_t based solely on visit observations gives:

$$\mathbb{E}[Z_t \mid D_t] = P(A = a) \cdot \mathbb{E}[Z_t \mid D_t, A = a] + P(A \neq a) \cdot \mathbb{E}[Z_t \mid D_t, A \neq a]$$
$$= P(A = a) \cdot \mathbb{E}\left[g^{-1}\left(\log(\lambda_t) - \beta_A^{(A)}\right) \mid D_t, A = a\right]$$

$$+ P(A \neq a) \cdot \mathbb{E} \left[g^{-1} \left(\log(\lambda_t) - \beta_A^{(A)} \right) \mid D_t, A \neq a \right]$$

$$< P(A = a) \cdot \mathbb{E} \left[g^{-1} \left(\log(\lambda_t) - \beta_A^{(A)} \right) \mid D_t, A = a \right]$$

$$+ P(A \neq a) \cdot \mathbb{E} \left[g^{-1} \left(\log(\lambda_t) - \beta_A^{(a)} \right) \mid D_t, A = a \right]$$
(*)
$$= P(A = a) \cdot \left(g^{-1} \left(\log(\lambda_t) - \beta_A^{(a)} \right) \right) + P(A \neq a) \cdot \left(g^{-1} \left(\log(\lambda_t) - \beta_A^{(a)} \right) \right)$$

$$= g^{-1} \left(\log(\lambda_t) - \beta_A^{(a)} \right)$$

1041

1039

1043

1045

1046

1052 1053 1054

1061

1062

As justification for (*):

 $\beta_A^{(a)} < \beta_A^{(A)}, \quad \forall A \neq a$

 $= \mathbb{E}[Z_t \mid D_t, A = a]$

(Disparity assumption)

$$\Rightarrow \log(\lambda_t) - \beta_A^{(a)} > \log(\lambda_t) - \beta_A^{(A)}, \quad \forall A \neq a, \forall \lambda_t$$

$$\Rightarrow g^{-1} \left(\log(\lambda_t) - \beta_A^{(a)} \right) > g^{-1} \left(\log(\lambda_t) - \beta_A^{(A)} \right), \quad \forall A \neq a, \forall \lambda_t$$

$$(g \text{ strictly monotonically increasing} \Rightarrow g^{-1} \text{ strictly monotonically increasing})$$

$$\implies \mathbb{E}\left[g^{-1}\left(\log(\lambda_t) - \beta_A^{(a)}\right) \mid D_t, A = a\right] > \mathbb{E}\left[g^{-1}\left(\log(\lambda_t) - \beta_A^{(A)}\right) \mid D_t, A \neq a\right]$$

It is clear to see that this argument extends naturally to show that if a group tends to visit the hospital more frequently conditional on severity, that their severity will be overestimated: if there exists a group \tilde{a} such that $\beta_A^{(\tilde{a})} > \beta_A^{(A)}$ for all $A \neq \tilde{a}$, then we will see that $\mathbb{E}[Z_t \mid D_t, A = \tilde{a}] < \mathbb{E}[Z_t \mid D_t]$. Thus any model that does not take into account demographic disparities in patient visit rates given their severity will lead to biased estimates of severity.

C SIMULATIONS

Figure S1 shows the results of 30 simulation runs, where we randomly instantiate the parameters of our model and then generate data to fit on. We generate simulated data for 1000 patients on each run, each of whom is assigned to one group (50% chance of being from either group). We visualize the recovery of each parameter by plotting true parameter values versus recovered posterior mean values, with one dot per run.

To generate data with prevalent disparities, we set our priors to $\mu_{Z_0} \sim \mathcal{N}(0, 2.5)$ and $\sigma_{Z_0} \sim \mathcal{TN}(1, 0.5)$ (normal distribution restricted to positive values) for the non-pinned group; $\mu_R \sim \mathcal{N}(0, 3)$ and $\sigma_R \sim \mathcal{TN}(1, 0.01)$ (normal distribution restricted to positive values) for both groups; $F \sim \mathcal{TN}(1, 1)$ (normal distribution restricted to values above 0.5 to enforce positive constraint) for $F_0; F \sim \mathcal{N}(0, 2)$ for all other features; $b \sim \mathcal{N}(0, 1); \psi \sim \mathcal{TN}(8, 1)$ (normal distribution restricted to positive values); $\beta_0 \sim \mathcal{N}(1.5, 0.1); \beta_Z \sim \mathcal{N}(0.5, 0.1);$ and $\beta_A \sim \mathcal{N}(0, 2)$ for the non-pinned group.

1076

1077 D NYP HEART FAILURE DATA PROCESSING

- 1078
- 1079 This study was conducted in accordance with Health Insurance Portability and Accountability Act (HIPAA) guidelines and with Institutional Review Board (IRB) approval.

1080 **Cohort filtering.** We analyze patients with *heart failure with reduced ejection fraction* (HFrEF) whom we identify, following clinical guidance, by filtering the available NYP data for patients 1082 who have at least one LVEF measurement below 50% and who have been recorded as receiving a diuretic prescription. To ensure we have relatively complete records for each patient, we then 1084 filter for patients who are likely to receive most of their cardiology care within the NYP system, by filtering for patients whose home zipcode is in the New York metropolitan area and who have at least two LVEF or BNP records at least 6 months apart within our data. Lastly, NYP switched electronic 1086 health record (EHR) systems, introducing inconsistencies in record-keeping across sites and years; 1087 to ensure our records are consistently recorded, we analyze data from Weill Cornell Medical Center, 1088 one of NYP's two largest sites, between January 1, 2012 (the start of reliable record-keeping) to 1089 October 2, 2020 (NYP Cornell's transition to a new EHR). This ensures records are consistently 1090 recorded in our data.

1091

Feature processing. We convert pBNP to BNP with the conversion pBNP = 6.25 * BNP (Rørth et al., 2020) and then log-transform BNP values to get one combined $\log_2(BNP)$ feature (Hendricks et al., 2022). We then normalize (z-score) all feature values so that each feature has mean 0 and variance 1. Because patient blood pressure and heart rate are much more likely to be measured at hospital visits unrelated to heart failure (while visiting another specialist in the NYP system), we limit patient observations to visits where a patient had one measurement of at least one of LVEF and BNP.

1099 We encode demographic categories by making A a one-hot encoding of race/ethnicity groups. 1100 Lastly, we describe the time scale of our model. As mentioned in §6, we discretize time in 1-1101 week bins; if a patient has multiple measurements of one feature within a timestep, we average all 1102 measurements within that timestep. Discretizing time in this way allows us to capture more long-1103 term changes rather than acute changes in patient status. We normalize time so that the total time range in our model is 0 to 1. The longest patient trajectory in our data is 446 weeks (timesteps), so 1104 we normalize timestep values so that they range from 0 to 1; we therefore have fractional, discrete 1105 time values, each representing one week as 1/446 units of time. 1106

1107

1109

1108 E MODEL EVALUATION

Fitting model on real data. We fit our model on real data using weakly informative priors: $\mu_{Z_0} \sim \mathcal{N}(0,1)$ and $\sigma_{Z_0} \sim \mathcal{TN}(1,1)$ (normal distribution restricted to positive values) for the non-pinned groups; $\mu_R \sim \mathcal{N}(0,1)$ and $\sigma_R \sim \mathcal{TN}(1.5,1)$ (normal distribution restricted to positive values) for both groups; $b \sim \mathcal{N}(0,1)$; $\psi \sim \mathcal{TN}(1,0.5)$ (normal distribution restricted to positive values); $\beta_0 \sim \mathcal{N}(2.5,1)$; $\beta_Z \sim \mathcal{N}(0,1)$; and $\beta_A \sim \mathcal{N}(0,1)$ for the non-pinned group.

For F, we set model priors using Factor Analysis: at t = 0, we have $Z_t = Z_0 \sim \mathcal{N}(0, 1)$ for group a_0 , meaning the mapping between severity and features

1117
1118

$$X_0 = F \cdot Z_0 + b + \epsilon_t$$

1118
 $\epsilon \to \mathcal{N}(0, \mathbf{H})$

$$\epsilon_t \sim \mathcal{N}(0, \Psi)$$

1119 captures a factor analysis model with factor loading matrix F and diagonal covariance matrix Ψ . We 1120 run factor analysis using feature measurements from the *first timestep* of all White patients (our a_0 1121 group) and use the estimates of F from Factor Analysis as the mean of our priors on F. We define 1122 the variance of our priors on F to be 1, and we pin the sign of F_{LVEF} to be negative for identifiability. 1123 Since we have no inherent value scale for what F values should be, Factor Analysis allows us to fit 1124 the model on more substantiated priors for feature scaling factors.

1125 We then fit the model and get the parameter estimates from 1000 samples. For any time t, we can calculate an estimate of Z_t and X_t for each sample, based on the sample's parameter estimates; we then take the average over all samples to get a patient's estimate of Z_t and X_t . In order to get actual feature value estimates, we can linearly transform X_t to undo the normalization for each feature and recover an estimate of each feature value at t. We can then use our model's estimates of Z_t and predicted feature values to analyze and evaluate our model's behavior.

1131

Comparison to baselines. We filter out patients who do not have at least three visits (since several of the baselines we fit require this many visits per patient, as we describe below), leaving a total of 1834 patients: 1118 White, 347 Black, 216 Hispanic, and 153 Asian patients.

1134 To evaluate our model's ability to reconstruct feature values, we compare our model to PCA and FA. 1135 PCA and FA require consistent dimensionality of the input data, so we fit all models on the first three 1136 visits for each patient. We train two variants of both PCA and FA: the first attempts to reconstruct 1137 patient visits from a single latent dimension (analogous to Z in our model), taking as input the X_t 1138 vector at one visit (4 features total) and representing it with a single latent component. The second variant attempts to reconstruct *patient trajectories* from two latent dimensions (analogous to Z_0 and 1139 R in our model), taking as input a concatenated vector of features X_t from the first three visits (12) 1140 features total) and representing it with two latent components. We impute missing values as the 1141 overall mean of the data for both PCA and FA, since these methods cannot naturally handle missing 1142 data. 1143

To evaluate our model's ability to predict future feature values, we compare our model to last time-1144 step, logistic regression, and quadratic regression. Unlike PCA and FA, these methods do not require 1145 consistent dimensionality in the input data, so we fit the models to the first three years of observed 1146 data. Last-timestep predicts all future feature values to be equal to the most recent feature value 1147 observed in the training data for that patient; if there is no observed feature value, the baseline 1148 predicts the population mean. Linear regression regresses values on time for each patient and each 1149 feature to predict future feature values. For patients with fewer than 2 observations for a given 1150 feature value, we use the population mean for the preceding or subsequent timestep. Quadratic 1151 regression follows a similar approach. Because linear regression and quadratic regression can overfit 1152 the data and make unrealistic predictions, we clip their predicted feature values to a range determined 1153 by that observed within the training data. 1154

Ablated Model. We compare our full model to an ablated version of the model that does not account for any of our three disparities. We do this by removing all group-specific parameters from the model, while leaving everything else the same: we learn one value of μ_R and σ_R and exclude β_A from the model. Since the distribution of Z_0 must be fixed for at least one group for identifiability (to fix the scale of Z_t), the distribution is pinned for all groups. Factor Analysis for model priors on *F* is also fit on all patients rather than only on white patients.

1161

1163

1162 F DISPARITIES ESTIMATES

We first describe our calculations for §6.3 to estimate how much later Black and Asian patients start receiving care for heart failure compared to White patients. Our model learns the following:

1166	$\mu_{Z_0}^{(\text{Black})} = \mu_{Z_0}^{(\text{White})} + 0.22$
1168	$\mu_{Z_0}^{(\text{Asian})} = \mu_{Z_0}^{(\text{White})} + 0.27$

1169

The learned average rate of progression across all patients is 0.62. This means that Black patients come in 0.22/0.62 = 0.35 units of time later in their disease progression than White patients, and Asian patients come in 0.27/0.62 = 0.44 units of time later than White patients. Given that one unit of time is the longest patient trajectory, 8.5 years, this leads us to 3.0 and 3.8 years for Black and Asian patients, respectively.

Next we describe our calculations to estimate how much less frequently Black patients visit the hospital than White patients at the same disease severity. Our model learns that

$$\beta_A^{(\mathrm{Black})} = \beta_A^{(\mathrm{White})} - 0.11$$

1179 At the same disease severity Z_t , Black patients will have a visit rate of

1181
$$\lambda_t = \exp(\beta_0 + \beta_Z \cdot Z_t + (\beta_A^{(\text{White})} - 0.11))$$

 $= \exp(\beta_0 + \beta_Z \cdot Z_t + \beta_A^{\text{(White)}}) \cdot \exp(-0.11)$

1183
1184
$$= 0.897 \cdot \exp(\beta_0 + \beta_Z \cdot Z_t + \beta_A^{(White)})$$

1185 1186

¹¹⁸⁷ So at the same disease severity, we estimate that Black patients have a visit rate that is 90% that of a White patient's visit rate.



G SUPPLEMENTARY FIGURES AND TABLES

Figure S1: Parameter recovery from fitting our model to synthetic data. The priors from which we draw the synthetic data are: $\mu_{Z_0} \sim \mathcal{N}(0,2)$ and $\sigma_{Z_0} \sim \mathcal{TN}(1,1)$ (normal distribution restricted to positive values) for the non-pinned group; $\mu_R \sim \mathcal{N}(0,2)$ and $\sigma_R \sim \mathcal{TN}(1,1)$ (normal distribution restricted to positive values) for both groups; $F \sim \mathcal{TN}(1,1)$ (normal distribution restricted to values above 0.5 to enforce positive constraint) for F_0 ; $F \sim \mathcal{N}(0,2)$ for all other features; $b \sim \mathcal{N}(0,1)$; $\psi \sim \mathcal{TN}(5,1)$ (normal distribution restricted to positive values); $\beta_0 \sim \mathcal{N}(1.5,0.2)$; $\beta_Z \sim \mathcal{N}(0.5,0.1)$; and $\beta_A \sim \mathcal{N}(0,2)$ for the non-pinned group.

1233

- 1234 1235
- 1236
- 1237
- 1238
- 1239
- 1240
- 1241

RMSE: informative	Our model 0.67	FA _{visit}	PCA _{visit}	FA _{patient}	PCA _{patient}
	Our model	FAvisit	PCAvisit	FApatient	PCA _{patien}

Table S1: Our model compared to standard baselines for reconstruction performance. We compare to factor analysis and principal component analysis fit at the patient visit level (FA_{visit}, PCAvisit) and at the trajectory level (FApatient, PCApatient). Models are fit on the first 3 visits from each patient and evaluated on same data using root mean squared error (RMSE).

	Our model	Linear regression	Quadratic regression	Latest timestep
RMSE: informative	0.99	1.6	2.3	0.89
RMSE: all	0.98	1.8	2.5	0.98

Table S2: Our model compared to standard baselines for predictive performance. We compare to linear regression, quadratic regression, and latest timestep prediction, each fit at the patient feature level. Models are fit on data from the first 3 years of each patient's disease trajectory and evaluated on visits after 3 years using root mean squared error (RMSE).

10/10

