

Automated Enhancement and Curation of Healthcare Data

1 **Melody Greer^{1*}, Cilia Zayas¹, Sudeepa Bhattacharyya^{1,2*}**

2 ¹Department of Biomedical Informatics, University of Arkansas for Medical Sciences, Little Rock,
3 AR, USA

4 ²Arkansas Biosciences Institute, Department of Biological Sciences, Arkansas State University,
5 State University, AR, United States

6 *** Correspondence:**

7 Melody L Greer, Department of Biomedical Informatics, University of Arkansas for Medical
8 Sciences, Little Rock, AR, USA

9 501.313.5030

10 mlgreer@uams.edu

11

12 Sudeepa Bhattacharyya, Arkansas Biosciences Institute, Department of Biological Sciences,
13 Arkansas State University, State University, AR, United States

14 870.972.3712

15 sbhattacharyya@astate.edu

16 **Keywords: Data Quality, SDOH, Social Determinants of Health, Electronic Health Records,**
17 **EHR, Health Informatics, Healthcare, Data curation**

18 **Abstract**

19 Social and behavioral aspects of our lives significantly impact our health, yet minimal social
20 determinants of health (SDOH) data elements are collected in the healthcare system. In this study we
21 developed a repeatable SDOH enrichment and integration process to incorporate dynamically
22 evolving SDOH domain concepts from consumers into clinical data. This process included SDOH
23 mapping, linking compiled consumer data, data quality analysis and preprocessing, storage, and
24 visualization. Our preliminary analysis shows commercial consumer data can be a viable source of
25 SDOH factor at an individual-level for clinical data thus providing a path for clinicians to improve
26 patient treatment and care.

27 **1 Introduction**

28 Socioeconomic and behavioral aspects of our lives significantly impact our health, yet minimal social
29 determinants of health (SDOH) data elements are collected in the healthcare system. Information of
30 this type is needed for quality healthcare research and patient care because it is associated with the
31 full-spectrum of health outcomes from acute to chronic disorders. Studies indicate cancer (Alcaraz et
32 al. 2020), cardiovascular disease (Tamura et al. 2019), dementia (Nicholas et al. 2021), mental health
33 and substance-abuse disorders (Galea and Vlahov 2002, Alegría et al. 2018), viral infection (Greer et
34 al. 2021) and sleep (Grandner and Fernandez 2021) are among a long list of health problems
35 (Kivimäki et al. 2020) which are linked to social risk factors not frequently or consistently collected
36 for patients. The combined effect of missing, inconsistent, or inaccurate data also leads to bias in
37 machine learning, algorithms underlying clinical decision support, predictive analytics or other

38 healthcare processes. (Obermeyer et al. 2019, Cottrell et al. 2020, Seker, Talburt, and Greer 2022) To
39 avoid these problems as well as to gain rich insights from healthcare data we must be cognizant about
40 diverse data collection, veracity and data-quality.

41 Electronic health records (EHR) are assembled from clinical, insurance and basic demographic
42 information during the course of patient care. Although there is growing interest in including SDOH
43 data into EHR, the capture and management mechanisms for this process are uncertain.(Ancker et al.
44 2018, Feldman, Davlyatov, and Hall 2020, Gold et al. 2018) Current possibilities include (1) paper-
45 based or digitally collected social needs screening before a patient receives care, (2) during a visit
46 with a clinician, or (3) publicly available data sets that provide social context. All of these methods
47 have challenges. Screening data that is incorporated into the EHR appears fragmented to users,
48 increases the staff workload plus adds a data entry step where paper is used. (Gold et al. 2018) And
49 although clinicians felt SDOH information was valuable they are already pressed for time, and
50 adding to their list of clerical tasks is often not practical or even possible. (Tong et al. 2018) These
51 barriers are likely the reason for the findings of Frazee et al. that social needs screening for food,
52 housing, utilities, transportation, and experience with interpersonal violence was present in 24.4% of
53 hospitals, and 15.6% of physician practices. (Frazee et al. 2019) SDOH data documented within
54 unstructured EHR fields during a visit needs further study and will require natural language
55 processing tools to be integrated into healthcare practice. (Hatef et al. 2019) Publicly available data,
56 whether in raw form (i.e. Census), or aggregated into measures (i.e. SDI, ADI) has been valuable in
57 studying the relationships between socioeconomic and patient health status (Chamberlain et al. 2020,
58 Johnson-Lawrence, Zajacova, and Sneed 2017, Tung et al. 2018). However, appending community-
59 level data introduces issues with averaging. In their 2020 work on community-level and patient-level
60 social risk data Cottrell et al. observed that community-level data misses some patients that patient-
61 level data would not. (Cottrell et al. 2020) Further complicating the issue, there are no currently
62 accepted standard SDOH data elements. (Cantor and Thorpe 2018) This means that even in cases
63 where social risk data is consistently collected it cannot be easily shared with other healthcare
64 providers during transfers or referrals.

65 Compiled consumer data can address these issues. Finance and marketing businesses have
66 successfully used this data for over twenty years to find customers, understand their needs, and tailor
67 financial products. Using these same strategies, healthcare researchers and providers can improve
68 patient treatment and care. The consumer data includes individual-level SDOH data providing a
69 holistic snapshot of an individual's lifestyle. It includes amongst others, income, education, lifestyle
70 variables, language spoken, household size, smoking status, life events, hobbies, shopping activity
71 etc. that are not available in the insurance claims data or majority of EHR data. A support system, for
72 complex assessments (i.e., risk assessments) and calculations, is needed using information that helps
73 to arrive at conclusions regarding patient's health risk and treatment. The development of an
74 automated pipeline process for multisource healthcare data integration will provide this support. The
75 presence of information integrated from multiple different streams of data will support tools for
76 nurses, social workers, community health workers and patient navigators that supports decision
77 making by providing the ability to consider multiple factors simultaneously for patients and
78 clinicians. To begin the development of an automated pipeline process for multisource healthcare
79 data integration we have conducted a pilot study integrating clinical and consumer information
80 sources to evaluate multiple SDOH and clinical factors simultaneously.

81 **2 Methods**

82 Our goal was to prototype a repeatable clinical data enhancement process to incorporate compiled
83 consumer data into SDOH domain concepts. This process included social risk factor mapping,
84 linking compiled consumer data, data quality assessment, preprocessing, storage, and visualization.
85 Consumer data does not currently line up one-to-one with factors identified as SDOH, so purchased
86 data elements were mapped to concepts identified in social needs screening. Once the mapping was
87 complete, the compiled consumer data elements were linked to the patient population and stored in a
88 Microsoft SQL Server. The resulting data was then accessible using SQL queries, and the quality was
89 evaluated for completeness, consistency, and timeliness or temporal alignment. Visualization tools
90 allow researchers or clinical decision support developers to study the data before analysis or
91 application. Combining a SQL database with an academic license for Tableau provided easy access
92 for visualization, data analysis, and wrangling. As a proof-of-concept, we provided SAS JMP to
93 clinicians and researchers to determine what aspects of these tools deliver the most value. SAS JMP
94 can be easily connected directly to the database and provides many descriptive analysis and
95 visualization methods. Open source tools, like R and Python, can also serve the same purpose.
96 Throughout SDOH enhancement, a security and data privacy layer overlaid the entire process with
97 security features included in each step to ensure data privacy. Figure 1 depicts a snapshot of the
98 entire process.

99 **2.1 Mapping**

100 Various social determinants of health frameworks have been created to assist communities,
101 healthcare professionals, and others in better identifying and managing an expansive range of factors
102 influencing health outcomes. (Office of Disease Prevention and Health Promotion 2022a, World
103 Health Organization 2010, Rural Health Information Hub 2022) Our mapping strategy was
104 developed based on the Healthy People (HP) 2030 SDOH Framework for this study. (Office of
105 Disease Prevention and Health Promotion 2022a, b) HP 2030 not only continued the HP initiative,
106 which set national health targets for 2020 through 2030, but also designed a framework to organize
107 SDOH into five domains: (1) economic stability, (2) education, (3) social and community context, (4)
108 health and healthcare access, and (5) the neighborhood and built environment. HP 2030 outlined
109 essential SDOH within each of these domains. (Office of Disease Prevention and Health Promotion
110 2022b) Employment, food insecurity, housing instability, and poverty, for example, all fall under the
111 domain of economic stability. We mapped as many elements as possible based on the HP 2030
112 framework from the two consumer data sources.

113 **2.2 Clinical Data**

114 We requested electronic health records from the Clinical Data Repository (University of Arkansas for
115 Medical Sciences Translational Research Institute 2022) for all patients with chronic conditions (i.e.,
116 asthma, diabetes, heart disease, congestive heart failure, coronary artery disease, heart attack, stroke),
117 or contagious respiratory illness (i.e., influenza, or COVID-19). All data received was stored on one
118 of the following secure devices: institute supported controlled access server, institute supported
119 password protected desktop computer, encrypted password protected laptop. The data used for
120 linking social determinants information was name, address, DOB only. These demographics were
121 transmitted to the selected data compiler vendors via SFTP. As per the requirement of the Health
122 Insurance Portability and Accountability Act of 1996 (HIPAA) for protection of patient health
123 information these 3rd party vendors signed a Business Associate Addendum (BAA) with our medical
124 institution prior to accessing the patient identification. Following the addition of the SDOH the data
125 was de-identified to increase protection of the participants from any negative consequences in the
126 event of a data breach. De-identification was accomplished by deleting full name and address and

127 replacing them with a random identification number and RUCCA code (cite). The DOB was deleted
128 and replaced with age. Data was stored in a secure database server behind a firewall.

129 **2.3 Non-Clinical Data Integration and Refresh Process**

130 Commercial data is updated monthly and is made up of hundreds of different sources, including
131 consumer surveys, public records, purchase transactions, real estate data, offline and online buying
132 behavior, and warranty information. Wherever possible, compilers compare values from multiple
133 data sources to check accuracy for each element. Vendors that compile data for commercial purposes
134 (i.e., marketing) were identified and interviewed, and then costs were negotiated. SDoH data points
135 were appended by commercial data compilers or other external data sources using only patient name,
136 DOB, and address. These processes occur entirely within a database system using fully HIPAA
137 compliant vendors. All data was encrypted while in transit and immediately destroyed at the compiler
138 location after the completion of the processing. Not all clinical data will be matched to existing
139 consumer data during this process. This is the problem of coverage which refers to the number of
140 patients who could be linked to compiled data. Coverage varies by commercial compiler, but the
141 reasons for coverage variation may also be associated with varying aspects of the patient's lifestyle
142 (e.g., people who use cash exclusively are less likely to have a substantial digital imprint in consumer
143 databases).

144 **2.4 Data Quality Analysis**

145 Data quality is a constellation of factors essential for data collections. The quality of the linked
146 clinical and commercial data (henceforth called merged-data) was evaluated for conformance,
147 completeness, consistency, plausibility and temporal alignment. (Kahn et al. 2016) As was mentioned
148 above, accuracy, and consistency, were evaluated by compilers before data was purchased. After data
149 was purchased and linked with EHR data we further measured consistency of common data elements
150 between compilers from the different sources. In the merged-data there were 5 data elements that
151 were common between the compilers. We matched the data records based on these 5 data elements to
152 measure consistency between the compilers as an added level of data-quality assessment. Each of the
153 data elements differed in categories or levels, between the compilers. Our first step was to collapse
154 the categories in each element into identical categories so that they could be compared for
155 consistency. For example, the data element 'Home Market Value, Estimated' in compiler 1 mapped
156 to 'Home Value Range_VDS_Appended' in compiler 2. The categories of each these feature
157 variables were however not the same. 'Home Market Value, Estimated' in compiler 1 had the
158 following 20 categories; "\$1,000 - \$24,999", "\$25,000 - \$49,999", "\$100,000 - \$124,999",
159 "\$125,000 - \$149,999", "\$150,000 - \$174,999", "\$175,000 - \$199,999", "\$200,000 - \$224,999",
160 "\$225,000 - \$249,999", "\$250,000 - \$274,999", "\$275,000 - \$299,999", "\$300,000 - \$349,999",
161 "\$350,000 - \$399,999", "\$400,000 - \$449,999", "\$450,000 - \$499,999", "\$50,000 - \$74,999",
162 "\$500,000 - \$749,999", "\$75,000 - \$99,999", "\$750,000 - \$999,999", "\$1,000,000 Plus", "NA". The
163 'Home Value Range_VDS_Appended' in compiler 2 had the following 18 categories; "Under \$50k",
164 "\$50 - \$100k", "\$100 - \$150k", "\$150 - \$200k", "\$200 - \$250k", "\$250 - \$300k", "\$300 - \$350k",
165 "\$350 - \$400k", "\$400 - \$450k", "\$450 - \$500k", "\$500 - \$550k", "\$550 - \$600k", "600 -
166 \$650k", "\$650 - \$700k", "NA", "\$700 - \$750K", "\$750K +", "Unknown". In order to assess
167 consistency we first collapsed the categories in each variable in an intuitive and meaningful fashion
168 and made them identical. The new collapsed categories in both variables were: "Less100K", "100-
169 200K", "200-300K", "300-400K", "400-500K", "500Kplus". After discarding the "NA"s or
170 "Unknowns", the consistency between the two variables were calculated.

171 Missing data are a pervasive problem in any source of data also referred to as data ‘completeness’.
172 Lack of complete data can significantly affect a study outcome by introducing unwanted bias. This is
173 why it was important for us to obtain consumer marketing data from a diverse selection of sources to
174 ensure that we have a collection of data that is complete and deep enough to provide meaningful
175 information about majority of the study participants.

176 During preprocessing we also examined conformance and plausibility of the data elements, thus
177 comparing the actual format of the data against the expected, and evaluating the feasibility of
178 multiple existing values of the data elements. Measuring the persistence of the data was not possible
179 in this work since to analyze changes in the data over time multiple batches of consumer data would
180 be needed which was cost-prohibitive.

181 **2.5 Bias**

182 Human bias exists. As we collect, analyze and take actions based on data our biases are perpetuated.
183 This bias pervades healthcare in machine learning, decision support, operations and logistics
184 planning in health systems. Applications to guide clinical practice are not exempt. Real-world
185 clinical data is important for clinical decision-making but it has everyday biases imprinted within it
186 and can preserve or even amplify health disparities. To address this issue requires detection and
187 correction. Sensitive attributes (i.e., race, gender, etc.) are evaluated against classifications to first
188 detect bias that may be present. Once uncovered biased data can be rebalanced using class labels.
189 While this method is not full proof, biases that are not expected may remain hidden, it does allow for
190 the mitigation of known issues and for the discovery of unknown issues.

191 **2.6 Data Analysis and Visualization**

192 After the data was linked and preprocessed we explored the data using Tableau, SAS JMP and R
193 statistical software packages that were linked to the Sql server hosting the merged-data. We
194 determined the summary statistics of a subset of the variables, from both EHR and compiled SDOH
195 sources, that characterized the patients. We are also currently in the process of building classifiers
196 that can predict disease risk based on patients’ clinical, demographic and SDOH factors. There are
197 open-source tools available such as R, Python and others, that can perform interactive data analysis
198 and visualization at no cost and can be customized to the needs of clinicians. For example, Shiny is
199 an R package that can be used to build dashboards or provide the clinicians or researchers a means to
200 interact with the data as per their need. Our ongoing effort focuses on developing an R-Shiny based
201 analysis and visualization tool as well.

202 **3 Results**

203 **3.1 Mapping and Linking Consumer Data with EHR at an individual-level**

204 **Table 1** lists SDOH domains, elements, and coverage percentages from two consumer data databases
205 that were used in our study and named here as compiler 1 and compiler 2. There were 55,422 patients
206 in the initial set of EHR records. After linking with two commercial compilers, 30,895 and 54,880
207 patients with SDOH data remained respectively. As the table shows, compiler 1 had fewer mapped
208 data elements compared to compiler 2. However, compiler 1 had, on an average, ~90% coverage on
209 the mapped data elements while compiler 2 had many more elements mapped but the coverage was
210 much lower, ~54% on an average. Thus, consumer data collection from at least two sources ensured
211 that all categories of SDOH domains, based on Healthy People 2020 SDOH Framework were
212 covered in our merged-data. (Office of Disease Prevention and Health Promotion 2022a) The

213 connection between unlinked patients' needs further study to determine what they have in common
214 aside from a minimal individual digital footprint.

215 **3.2 Data quality assessment and preprocessing.**

216 To measure consistency between the compiler data collected from different sources we matched the
217 data records based on elements that were common between the compilers as part of data-quality
218 assessment. **Table 2** shows the randomly picked data elements and their percent matches. The
219 percent matches of compiler 2 ranged between ~21% - 64% with compiler 1. This underscores the
220 importance of data collected from multiple, trusted and standardized sources.

221 **3.3 Bias**

222 Underserved rural populations are less likely to get needed healthcare due to distance, costs, and poor
223 insurance coverage leading to underdiagnosis of illnesses even though they are more commonly
224 affected by chronic conditions as shown in **Table 3**. To study this problem of bias in healthcare data,
225 we have analyzed and preprocessed a real-world data set of patients with chronic conditions from
226 geographically disparate locations. We have chosen to preprocess because it prepares the data set for
227 any following modeling and does not need to be repeated for each new classifier. We also tested the
228 removal of the sensitive attribute altogether. Each of these tests produces similar AUC results. in
229 **Table 3** or similar. (Seker, Talburt, and Greer 2022)

230 **3.4 Preliminary Analysis and Visualization of the Merged Data**

231 The SDOH data merged with EHR provided insights into the social risk factors of disease both at
232 patient level as well as the population level. We first explored the disparities in demographics and
233 SDOH factors between the 55,422 patients who were Covid positive compared to those who were
234 not. **Table 4** provides a summary of a subset of patient characteristics that were compared between
235 these two groups. Our preliminary analysis showed that apart from demographic factors, several
236 SDOH factors like home-ownership, marital-status, presence of children, number of members per
237 household, economic stability index and education were significant different between the two patient
238 groups while estimated family-income and home market-value were not. **Figure 2** shows a map of
239 the top COVID-affected zip codes in Arkansas overlaid on a heatmap image of average economic
240 stability index of those zip codes. The darker blue counties were less economically stable and also
241 had higher percentages of covid positive patients.

242 **4 Discussion**

243 In this proof-of-concept study our main objective was to evaluate the viability of consumer marketing
244 data, purchased from 3rd party HIPPA-compliant vendors, as a source of SDOH factors that are
245 largely missing from the EHR data. We purchased in-depth patient-level consumer data from two
246 different vendors, mapped a wide array of data elements to the 5 broad SDOH domains as defined by
247 the Healthy People 2030 SDOH framework, linked the consumer data to EHR patient level data,
248 stored the data in a SQL server linked with several statistical and data exploration tools, evaluated
249 data quality and preprocessed the data, and lastly completed a preliminary analysis of a subset of the
250 SDOH elements that characterized the patients. To our knowledge this is the first study to explore the
251 viability of consumer marketing data as a source of patient-level SDOH data.

252 With recent upsurge in research solidifying the significant relationship between SDOH and
253 population health, an increasing number of healthcare stakeholders are exploring the use of public

254 databases for community-level information in order to identify those patients that are most vulnerable
255 to SDOH. For example, census tracts data have been used to identify areas associated with socio-
256 economic risks and poor health outcomes.(Liaw et al. 2018) But a recent study by Cottrell et al
257 (Cottrell et al. 2020) showed that only about 48% of the times community-level data can accurately
258 identify social risks at the patient level. Thus, healthcare decisions on individual patients based on
259 community-level data may fall short on providing adequate care to a significant number of patients.
260 This may give rise to the problem of ‘ecologic fallacy’ where incorrect assumptions can be made
261 about a patient based on aggregate-level information from community-level data. In this study we
262 have attempted to address this problem by partnering with companies/vendors that are honed
263 consumer market researchers.

264 We have developed a repeatable process to incorporate commercially compiled data into EHR data.
265 The added value has been demonstrated based on a published paper (Greer et al. 2021) and other in-
266 progress research efforts focused on disease specific predictive analytics. We have also identified
267 opportunities in data quality research areas that need further study as part of this work. The curated
268 data are being used to support several healthcare analytics applications, including descriptive
269 analytics, data visualization, and predictive modeling. During this work, we have developed the first
270 mapping scheme of commercial data elements with SDOH elements. This is a fascinating aspect of
271 this work because the healthcare community has not reached a consensus on a standard set of social
272 determinants of health concepts demanding that this process be agile and flexible. As we advance this
273 area, our goal is to work on mapping and semantics using an ontology

274 Building an enrichment process for EHR data has allowed us to study important questions about
275 temporal alignment, data dictionary and coverage issues, legal requirements, and security
276 requirements. Initially, the legal, research and business processes required were complex and time-
277 consuming. Patient data must be kept private and secure at all times, and all parties must be bound by
278 a contract to minimize the possibility of a data breach. Fortunately, once contracts and transfer
279 processes are in place, they remain active and available for repeated consumer data collection. This is
280 important because continued collection will be necessary. Compiled data becomes stale over time,
281 and EHR data is collected only at the time of each encounter. Aligning these time windows is
282 necessary for elements that must be current while is less critical for elements that are more likely to
283 remain stable over time. As the process iterates and newly compiled data is integrated into the EHR
284 data, the data dictionary must also be updated. We discovered that the data dictionaries provided by
285 compilers vary in quality and detail. In addition, compilers are continuously adding, removing, and
286 updating elements resulting in multiple versions of the dictionary documentation. Integrating data
287 from these dynamic systems also impacted the mapping of compiled elements onto SDOH concepts,
288 resulting in a mapping component for each iteration. If kept current, the SDOH mapping will require
289 minimal effort to maintain. Throughout all of these components, it was also necessary to tackle
290 practical Information technology issues such as storage, tools, permissions, and access which will
291 need to be customized to each institution.

292 Our study had several limitations. Due to budgetary constraints we restricted our SDOH data sources
293 to two different vendors only. One had more mapped data elements while the other had more
294 coverage, thus highlighting the need for data collection from multiple, trust worthy and reputable
295 vendors with standardized methods of data collection. There were significant missing data in each
296 data elements. Also, among the overlapping data elements, the concordance between the data was not
297 very high which underscores the need for good quality data sources. May be a standardized and
298 validated screening tool based on our approach needs to be developed that can be used by all
299 healthcare entities.

300 In conclusion, we have developed a repeatable SDOH enhancement process to incorporate
301 dynamically evolving SDOH domain concepts from consumers into clinical data. The literature
302 provides early and rapidly growing evidence that integrating individual-level SDOH into EHRs can
303 assist in risk assessment and predicting healthcare utilization and health outcomes, which further
304 motivates efforts to collect and standardize patient-level SDOH information. This study highlights
305 one potential means to incorporate individual-level patient data into EHR, thus opening up
306 possibilities for predictive analytics and enhanced solutions for providers, payers and healthcare
307 organizations to enable them to address the social needs of patients.

308 **5 Conflict of Interest**

309 *The authors declare that the research was conducted in the absence of any commercial or financial*
310 *relationships that could be construed as a potential conflict of interest.*

311 **6 Author Contributions**

312 MG developed the concept. MG SB worked on concept development analysis and manuscript
313 writing. CZ performed SDOH mapping and writing.

314 **7 References**

- 315 Alcaraz, Kassandra I, Tracy L Wiedt, Elvan C Daniels, K Robin Yabroff, Carmen E Guerra, and
316 Richard C Wender. 2020. "Understanding and addressing social determinants to advance
317 cancer health equity in the United States: A blueprint for practice, research, and policy." *CA:
318 A Cancer Journal for Clinicians* 70 (1):31-46.
- 319 Alegría, Margarita, Amanda NeMoyer, Irene Falgàs Bagué, Ye Wang, and Kiara Alvarez. 2018.
320 "Social determinants of mental health: where we are and where we need to go." *Current
321 psychiatry reports* 20 (11):1-13.
- 322 Ancker, Jessica S, Min-Hyung Kim, Yiye Zhang, Yongkang Zhang, and Jyotishman Pathak. 2018.
323 "The potential value of social determinants of health in predicting health outcomes." *Journal
324 of the American Medical Informatics Association* 25 (8):1109-1110.
- 325 Cantor, Michael N, and Lorna Thorpe. 2018. "Integrating data on social determinants of health into
326 electronic health records." *Health Affairs* 37 (4):585-590.
- 327 Chamberlain, Alanna M, Lila J Finney Rutten, Patrick M Wilson, Chun Fan, Cynthia M Boyd, Debra
328 J Jacobson, Walter A Rocca, and Jennifer L St Sauver. 2020. "Neighborhood socioeconomic
329 disadvantage is associated with multimorbidity in a geographically-defined community."
330 *BMC public health* 20 (1):1-10.
- 331 Cottrell, Erika K, Michelle Hendricks, Katie Dambrun, Stuart Cowburn, Matthew Pantell, Rachel
332 Gold, and Laura M Gottlieb. 2020. "Comparison of community-level and patient-level social
333 risk data in a network of community health centers." *JAMA network open* 3 (10):e2016852-
334 e2016852.
- 335 Feldman, Sue S, Ganisher Davlyatov, and Allyson G Hall. 2020. "Toward Understanding the Value
336 of Missing Social Determinants of Health Data in Care Transition Planning." *Applied
337 Clinical Informatics* 11 (04):556-563.
- 338 Frazee, Taressa K, Amanda L Brewster, Valerie A Lewis, Laura B Beidler, Geneva F Murray, and
339 Carrie H Colla. 2019. "Prevalence of screening for food insecurity, housing instability, utility

- 340 needs, transportation needs, and interpersonal violence by US physician practices and
341 hospitals." *JAMA network open* 2 (9):e1911514-e1911514.
- 342 Galea, Sandro, and David Vlahov. 2002. "Social determinants and the health of drug users:
343 socioeconomic status, homelessness, and incarceration." *Public health reports* 117 (Suppl
344 1):S135.
- 345 Gold, Rachel, Arwen Bunce, Stuart Cowburn, Katie Dambrun, Marla Dearing, Mary Middendorf,
346 Ned Mossman, Celine Hollombe, Peter Mahr, and Gerardo Melgar. 2018. "Adoption of social
347 determinants of health EHR tools by community health centers." *The Annals of Family
348 Medicine* 16 (5):399-407.
- 349 Grandner, Michael A, and Fabian-Xosé Fernandez. 2021. "The translational neuroscience of sleep: A
350 contextual framework." *Science* 374 (6567):568-573.
- 351 Greer, Melody L, Steven Sample, Hanna K Jensen, Sacha McBain, Riley Lipschitz, and Kevin W
352 Sexton. 2021. "COVID-19 is connected with lower health literacy in rural areas." *Studies in
353 health technology and informatics* 281:804.
- 354 Hatef, Elham, Masoud Rouhizadeh, Iddrisu Tia, Elyse Lasser, Felicia Hill-Briggs, Jill Marsteller, and
355 Hadi Kharrazi. 2019. "Assessing the Availability of Data on Social and Behavioral
356 Determinants in Structured and Unstructured Electronic Health Records: A Retrospective
357 Analysis of a Multilevel Health Care System." *JMIR medical informatics* 7 (3):e13802.
- 358 Johnson-Lawrence, Vicki, Anna Zajacova, and Rodlescia Sneed. 2017. "Education, race/ethnicity,
359 and multimorbidity among adults aged 30–64 in the National Health Interview Survey."
360 *SSM-population health* 3:366-372.
- 361 Kahn, MG, TJ Callahan, J Barnard, AE Bauck, J Brown, BN Davidson, H5 Estiri, C Goerg, E Holve,
362 SG Johnson, ST Liaw, M Hamilton-Lopez, D Meeker, TC Ong, P Ryan, N Shang, NG
363 Weiskopf, C Weng, MN Zozus, and L Schilling. 2016. "A Harmonized Data Quality
364 Assessment Terminology and Framework for the Secondary Use of Electronic Health Record
365 Data." *EGEMS (Wash DC)* 4 (1):1244.
- 366 Kivimäki, Mika, G David Batty, Jaana Pentti, Martin J Shipley, Pyry N Sipilä, Solja T Nyberg,
367 Sakari B Suominen, Tuula Oksanen, Sari Stenholm, and Marianna Virtanen. 2020.
368 "Association between socioeconomic status and the development of mental and physical
369 health conditions in adulthood: a multi-cohort study." *The Lancet Public Health* 5 (3):e140-
370 e149.
- 371 Liaw, Winston, Alex H Krist, Sebastian T Tong, Roy Sabo, Camille Hochheimer, Jennifer Rankin,
372 David Grolling, Jene Grandmont, and Andrew W Bazemore. 2018. "Living in “cold spot”
373 communities is associated with poor health and health quality." *The Journal of the American
374 Board of Family Medicine* 31 (3):342-350.
- 375 Nicholas, Lauren Hersch, Kenneth M Langa, Julie PW Bynum, and Joanne W Hsu. 2021. "Financial
376 presentation of Alzheimer disease and related dementias." *JAMA internal medicine* 181
377 (2):220-227.
- 378 Obermeyer, Ziad, Brian Powers, Christine Vogeli, and Sendhil Mullainathan. 2019. "Dissecting
379 racial bias in an algorithm used to manage the health of populations." *Science* 366
380 (6464):447-453.
- 381 Office of Disease Prevention and Health Promotion. 2022a. "Healthy People 2030 Framework." U.S.
382 Department of Health and Human Services,, accessed 2/15.

- 383 [https://www.healthypeople.gov/2020/About-Healthy-People/Development-Healthy-People-](https://www.healthypeople.gov/2020/About-Healthy-People/Development-Healthy-People-2030/Framework)
384 [2030/Framework](https://www.healthypeople.gov/2020/About-Healthy-People/Development-Healthy-People-2030/Framework).
- 385 Office of Disease Prevention and Health Promotion. 2022b. "Social Determinants of Health." U.S.
386 Department of Health and Human Services,, accessed 2/15.
387 <https://health.gov/healthypeople/objectives-and-data/social-determinants-health>.
- 388 Rural Health Information Hub. 2022. "Frameworks to Address Social Determinants of Health." Rural
389 Health Information Hub,, accessed 2/15.
390 <https://www.ruralhealthinfo.org/toolkits/sdoh/1/frameworks>.
- 391 Seker, Emel, John R Talburt, and Melody L Greer. 2022. "Preprocessing to Address Bias in
392 Healthcare Data." *Accepted for publication in the Medical Informatics Europe (MIE) 2022*
393 *conference proceedings*.
- 394 Tamura, Kosuke, Steven D Langerman, Joniqua N Ceasar, Marcus R Andrews, Malhaar Agrawal,
395 and Tiffany M Powell-Wiley. 2019. "Neighborhood social environment and cardiovascular
396 disease risk." *Current cardiovascular risk reports* 13 (4):1-13.
- 397 Tong, Sebastian T, Winston R Liaw, Paulette Lail Kashiri, James Pecsok, Julia Rozman, Andrew W
398 Bazemore, and Alex H Krist. 2018. "Clinician experiences with screening for social needs in
399 primary care." *The Journal of the American Board of Family Medicine* 31 (3):351-363.
- 400 Tung, Elizabeth L, Kristen E Wroblewski, Kelly Boyd, Jennifer A Makelarski, Monica E Peek, and
401 Stacy Tessler Lindau. 2018. "Police-recorded crime and disparities in obesity and blood
402 pressure status in Chicago." *Journal of the American Heart Association* 7 (7):e008030.
- 403 University of Arkansas for Medical Sciences Translational Research Institute. 2022. "Clinical Data
404 Repository (AR-CDR)." accessed 1/15. <https://ar-cdr.uams.edu/>.
- 405 World Health Organization. 2010. A Conceptual Framework For Action on the Social Determinants
406 of Health. In *Social Determinants of Health Discussion Paper 2: World Health Organization*,.
- 407
- 408

409 Figure 1: Workflow describing data integration of clinical and SDOH factors translated from consumer information
410 sources.

411 Figure 2: This map of central Arkansas shows that lower economic stability (indicated in orange) trends with higher
412 instances of COVID-19 positivity (indicated in percentage label) but that was not true for all zip codes as there were other
413 SDOH confounders.

414

415

416 Table 1: Data elements mapped from consumer databases to SDOH categories with coverage
417 percentages listed for each individual compiler.

SDOH Domains	Element Category	Compiler 1 Coverage	Compiler 2 Coverage
Economic Stability	Employment		
	▪ Occupation	48.23%	54.31%
	Food Insecurity		
	▪ Health Natural Foods	Not reported	54.31%
	▪ Food and Drink	Not reported	54.31%
	▪ Grocery	Not reported	54.31%
	Housing Instability		
	▪ Homeowner / Renter	98.98%	54.31%
	▪ Length of Residence	98.98%	54.31%
	Poverty		
	▪ Estimated Income	98.98%	54.31%
	▪ Net Worth Indicator	85.98%	54.31%
	▪ Economic Stability Indicator	85.98%	Not reported
	▪ Estimated Discretionary Income %	Not reported	54.31%
	▪ Estimated Household Debt Level	Not reported	54.31%
	▪ Loan to Value Ratio	Not reported	54.31%
	▪ Public Housing	Not reported	54.31%
Education	Early Childhood Education and Development		
	▪ Education Level	Not reported	54.31%
	Enrollment in Higher Education		
	▪ Occupation Student	Not reported	54.31%
	▪ Presence of College Graduate	Not reported	54.31%
	High School Graduation		
	▪ Education Level	Not reported	54.31%
	Language and Literacy		
	▪ Country of Origin	85.80%	Not reported
	▪ Hispanic Language Preference	85.80%	Not reported
▪ Likes to Read	Not reported	54.31%	
Social and Community Context	Civic Participation		
	▪ Activism Social Issues	Not reported	54.31%
	▪ Community Civic Activities	Not reported	54.31%
	▪ Charitable Volunteer	Not reported	54.31%
	▪ Registered Voter Indicator	Not reported	54.31%
	Discrimination	Not reported	Not reported
	Incarceration	Not reported	Not reported
	Social Cohesion		
	▪ Community Groups	Not reported	54.31%
	▪ Community and Family	Not reported	54.31%
	▪ Recreation	Not reported	54.31%
	▪ Sports	Not reported	54.31%
	▪ Travel Family Vacations	Not reported	54.31%
	▪ Caregiver in Home	Not reported	48.89%
Health and Healthcare	Access to Healthcare		
	▪ Insurance	Not reported	54.31%
	▪ Percent Healthcare Uninsured	Not reported	54.31%
	▪ Prescription: Number of Drugs	Not reported	54.31%
	▪ Long Term Care Insurance Index	Not reported	37.46%
	▪ Medicare Supplement Insurance Buyer Index	Not reported	54.31%

	▪ Single Service Plan Vision	Not reported	54.31%
	▪ Single Service Plan Dental	Not reported	54.31%
	▪ Single Service Plan Disability	Not reported	54.31%
	Access to Primary Care		
	▪ Health Rank Number of Physicians	Not reported	48.89%
	▪ Health Rank Doctor Visits	Not reported	48.89%
	Health Literacy		
	▪ Reading Cooking or Culinary	Not reported	54.31%
	▪ Reading Medical or Health	Not reported	54.31%
	▪ Reading Natural Health Remedies	Not reported	54.31%
Neighborhood and Built Environment			
	Access to Food that Support Healthy Eating Patterns	Not reported	Not reported
	Crime and Violence		
	▪ Concealed Weapons	Not reported	54.31%
	Environmental Conditions		
	▪ Census Percent Mobile Homes	Not reported	54.31%
	▪ Census Average Number of Automobiles	Not reported	54.31%
	▪ Digital Neighborhoods	Not reported	54.31%
	Quality of Housing		
	▪ Home Market Value, Estimated	94.16%	54.31%
	▪ Home Building Repair	Not reported	54.31%

418

419

420 Table 2: Percent matches between overlapping elements of the two compilers.

Variable name		No of categories	Percent match
Compiler 1	Compiler 2		
Home Market Value Estimated	Home Value Range	6	55.68%
Marital Status	Marital Status	2	64.26%
Presence of Children	Presence of Children	2	41%
Income Estimate Household	Income Range	9	20.85%

421

422

423 Table 3: As age increases, multimorbidity increases. However, in the case of rural residents a
424 decrease in the average number of chronic conditions reflects bias which could impact patient care
425 where predictive algorithms use this information.

	Geographic Location	
Age Groups	Rural	Urban
18-64	2.3127	2.7463
65+	2.2340	3.0272

426

427 Table 4: Summary statistics of demographics and individual-level SDOH factors in the merged data.

Characteristics	N	COVID Positive		p-val ²
		No, N = 54330 ¹	Yes, N = 1092 ¹	
Gender	55422			0.032
F		29611 (55%)	638 (58%)	
M		24706 (45%)	454 (42%)	
U		13 (<0.1%)	0 (0%)	
Age	55422	62 (18)	46 (18)	<0.001
Race	46869			<0.001
African American		10220 (22%)	298 (32%)	
Asian		551 (1.2%)	25 (2.7%)	
Hispanic		1445 (3.1%)	110 (12%)	
White/Other		33726 (73%)	494 (53%)	
Home Owner Renter	54068			<0.001
Home Owner		34558 (65%)	605 (56%)	
Renter		18420 (35%)	485 (44%)	
Marital Status	54068			<0.001
Married		26140 (49%)	391 (36%)	
Single		26838 (51%)	699 (64%)	
Has Children	54068	14826 (28%)	377 (35%)	<0.001
Member Household	54068			<0.001
1		20473 (39%)	487 (45%)	
2		15603 (29%)	224 (21%)	
3		7772 (15%)	173 (16%)	
4		5201 (9.8%)	112 (10%)	
5		3305 (6.2%)	71 (6.5%)	
6		393 (0.7%)	13 (1.2%)	
7		145 (0.3%)	3 (0.3%)	
8		61 (0.1%)	6 (0.6%)	
8Plus		25 (<0.1%)	1 (<0.1%)	
Estimated Family Income	34792			0.4
Less30K		0 (0%)	0 (0%)	
\$30-50K		14070 (41%)	271 (43%)	
\$50-75K		9228 (27%)	165 (26%)	
\$75-100K		4448 (13%)	90 (14%)	
\$100-125K		2567 (7.5%)	36 (5.7%)	
\$125Kplus		3848 (11%)	69 (11%)	

Estimated Home Market Value	51438			0.9
Less100K		19231 (38%)	405 (39%)	
100-200K		21183 (42%)	434 (42%)	
200-300K		5719 (11%)	124 (12%)	
300-400K		2153 (4.3%)	39 (3.7%)	
400-500K		878 (1.7%)	20 (1.9%)	
500Kplus		1230 (2.4%)	22 (2.1%)	
Economic Stability Index	37135			<0.001
10-15		7126 (20%)	117 (15%)	
16-20		7696 (21%)	118 (15%)	
21-25		10940 (30%)	226 (28%)	
26-30		10573 (29%)	339 (42%)	
Education	34295			<0.001
Attended Vocational/Technical		361 (1.1%)	5 (0.7%)	
Completed College		11392 (34%)	187 (28%)	
Completed Graduate School		4008 (12%)	45 (6.6%)	
Completed High School		17857 (53%)	440 (65%)	

¹Statistics presented: n (%); mean (SD)

²Statistical tests performed: chi-square test of independence; Wilcoxon rank-sum test



