1    **Title**

2    Enhancing Race and Ethnicity using Bayesian Imputation in an All Payer Claims Database

3

4    **Authors and Institutions**

5    Sanae El Ibrahimi [1,2]

6    Michelle A. Hendricks [1]

7    Sara E Hallvik [1]

8    Nazanin Dameshghi [1]

9    Christi Hildebran [1]

10   Michael A. Fischer [3,4]

11   Scott G. Weiner [3, 4]

12

13   **Running head**

14   Bayesian Imputation of Race and Ethnicity in APCD

15

[1] Comagine Health, Division of Research and Evaluation, Portland OR, USA

[2] School of Public Health, Department of Epidemiology and Biostatistics, University of Nevada,

Las Vegas NV, USA

[3] Brigham and Women's Hospital, Department of Emergency Medicine, Boston MA, USA

[4] Harvard Medical School, Boston MA, USA

18    **Corresponding author**

19    Sanae El Ibrahimi

20    Comagine Health

21    650 NE Holladay St # 1700,

22    Portland, OR 97232

23    selibrahimi@comagine.org

24    (702) 336-4408

25

26    **Word count: 2969**

27    **Number of tables: 2**

28    **Figures: 0**

29    **Number of References: 36**

30

31         **Acknowledgements**

34

35

36    **Abstract**

37    **Background:** All Payer Claims Databases (APCD) are a rich source of health information,

38    however, race and ethnicity (R&E) data are largely missing. Bayesian Improved Surname

39    Geocoding (BISG) is a common R&E imputation method, yet, validation of BISG in APCDs is

40    lacking. We used the BISG to impute missing R&E in the Oregon APCD.

41    **Methods:** BISG imputed R&E for Asian Pacific Islanders (API), Blacks, Hispanics and Whites

42    were contrasted to the gold standard (vital statistics) and sensitivity and specificity improvements

43    were assessed. Logistic regression examined whether missing R&E was random across patient

44    characteristics.

45    **Results:** Among 85,857 individuals in the study, 32.1% (n=27,594) had missing R&E. Missing

46    R&E was not randomly distributed. There were higher odds of missingness among males,

47    Whites, those age 65 and older, and commercially insured individuals. Differences in the percent

48    missing were also found by co-morbid conditions and mortality causes. Imputing the missing

49    R&E with BISG method improved the sensitivity to identify White, Black, API, and Hispanics.

50    **Conclusions:** APCDs can benefit from enhancing missing R&E with BISG imputation to

51    perform more robust population-health level analyses and identify inequities according to R&E

52    without losing power or dropping non-random records with missing R&E data.

53

## Introduction

In the everchanging healthcare landscape, data are required to develop new approaches to improve healthcare quality, efficiently use resources, and analyze system performance. To meet this growing need, states are increasingly mandating through legislation that commercial and public payers providing insurance plans in their state submit data to All-Payer Claims Data (APCD) databases that include medical claims, pharmacy claims, insurance enrollment data, provider information, and dental claims [1]. To date, many states have existing (mandated or voluntary), under implementation or strong interest in APCDs [2]. States can use APCDs for many purposes, including to improve health system performance, assess the impact of policy changes, understand key cost and utilization drivers, monitor population health trends, develop interventions, and conduct research [3, 4].

APCDs are a rich source of clinical information that have great potential for providing the data needed to comprehensively address long-standing race and ethnic inequities in the quality of healthcare delivery at the population level [5]. This potential is constrained, however, by the lack of reliable race and ethnicity information in APCDs. Despite long-standing recommendations from the Institute of Medicine, the National Quality Forum, and others that health plans systematically collect race and ethnicity data, implementation has been slow due to privacy concerns and resource limitations [6-9]. Even Medicare and Medicaid plans, which are federally mandated to collect race and ethnicity data by the Affordable Care Act [10], continue to struggle to collect this information completely [11].

To address the ongoing challenges regarding limited race and ethnicity data in administrative datasets, Elliott and colleagues developed the Bayesian Improved Surname Geocoding (BISG) imputation method [12]. BISG estimates the probability that an individual is a member of a given

4

80    racial or ethnicity category based on their surname and their address. Surname analysis is

81    conducted using the U.S. Census Surname list, which provides common surnames for racial and

82    ethnic groups based on information collected from the decennial census [13]. Using Bayesian

83    estimation, this information is combined with the racial and ethnic composition in the census

84    block group where the individual lives to estimate likely racial and ethnic membership. BISG

85    can provide estimates of racial and ethnic health disparities at the group or population level [14].

86    It has recently been used to demonstrate significant racial/ethnic disparities in behavioral health

87    care quality measures among Medicare Advantage enrollees [15].

88    As BISG is increasingly used in healthcare delivery and quality improvement, validation of the

89    methodology in multiple settings is fundamental. Previous validation studies comparing BISG

90    estimates to self-reported race and ethnicity have found that the BISG algorithm reliably predicts

91    categories of racial and ethnic membership among commercial health plan members [16, 17].

92    BISG performs well for White, Black, Hispanic, and Asian individuals but identification of

93    American Indian/Alaskan Native (AIAN) and Multi-racial individuals is typically poor [14].  A

94    version of BISG has also been developed for use with Medicare data, with similar limitations

95    [18]. However, to date, validation of BISG with a multi-payer data set such as the APCD is

96    lacking [14].  Given the growing prominence of APCDs and their potential to inform state

97    policies to address health disparities in quality of care and patient outcomes, the validation of

98    reliable methods to estimate race and ethnicity in APCDs is an important public health issue.

99    The aim of this study is to validate and use BISG to estimate distribution of patient

100   characteristics according to race and ethnicity in Oregon's APCD. First, differential missingness

101   of race and ethnicity in the APCD is assessed by patient characteristics. Then, the ability of

102   BISG to accurately impute missing race and ethnicity for Asian/Pacific Islander (API), Black,

103    Hispanic, and White populations was examined using Oregon death certificates as the "gold

104    standard" [19].

105    **Methods**

106        2.1 data sources

107    The sample was extracted from the Oregon Data Collaborative's voluntary APCD, which

108    includes medical and pharmacy claims and enrollment data for Medicaid, Medicare Advantage,

109    and most Commercial plans in the state, covering about 80% of Oregonians.  The study sample

110    was limited to patients with a valid Oregon address and age between 2 and 100 years as of

111    1/1/2014. Detailed patient demographics (e.g. surnames and addresses) for Medicare Fee-for-

112    Service (FFS) beneficiaries were not available at the time of the study and, therefore, were not

113    included in this analysis. The APCD was linked to the vital statistics death records between

114    2013-2018 using a probabilistic methodology matching on name and date of birth using the

115    Fastlink R package [20].

116    Three sources of race/ethnicity were examined: 1) death records race/ethnicity (the gold

117    standard); 2) APCD race/ethnicity from the member eligibility files; and 3) BISG imputed

118    race/ethnicity when missing [16].

119        2.2 Race/ethnicity variable composition

120    The initial sample contained 105,240 individuals who had a vital statistics record and were

121    linked with the APCD. Individuals with missing race/ethnicity in vital statistics data were

122    excluded (n=197 (0.19%)).

123        *Vital statistics race/ethnicity*

124    A combined vitals race/ethnicity variable was created and included: American Indian/Alaska

125    Native (AIAN), Asian Pacific Islander (API), Black, Hispanic (any race), White, and "other"

126    race group. AIAN individuals and individuals with more than one race were placed in the "other

127    race" category due to the limitations of low imputation accuracy of BISG for AIAN and multi-

128    race [14].

129        *APCD race/ethnicity*

130    APCD race/ethnicity variables were transformed to align with the same categories used for vital

131    statistics resulting in these APCD self-reported race/ethnicity variable categories: API, Black,

132    Hispanic, other, unknown, White.

133        *BISG race/ethnicity imputation*

134    To impute BISG race/ethnicity for individuals in the APCD, geographic identifiers were

135    allocated to each members' address of residence in the eligibility files using the Census

136    Geocoding Application Performing Interface Services [21]. Addresses were then normalized (i.e.

137    street, city, state, zip) to match the required format to perform batch geocoding. Next, the 12-

138    digit block group (an area containing about 1,000 persons) Federal Information Processing

139    Standard (FIPS) codes were extracted [12, 16]. Lastly, the five race/ethnicity probabilities (API,

140    Black, Hispanic, Other, White) were calculated using the R package wru: *Who Are You?*

141    *Bayesian Prediction of Racial Category Using Surname and Geolocation* [22] that can be

142    accessed from https://github.com/kosukeimai/wru. The package produced probabilities for each

143    category based on members' surname (last name) compared to the 2010 Census Bureau Surname

144    enumeration list, the Spanish Surname list [23] and the racial/ethnic distribution of the

145    individuals' Census 2000 block group.  The AIAN race category in BISG imputation was not

146    included because imputing AIAN from surname and residence for this classification remains

147    unreliable [24].

148    To estimate the optimal probability thresholds to create binary BISG race/ethnicity variables, the

149    Youden Index optimal cut off [25] was applied which maximizes the difference between the true

150    positives and the false positives on the receiver operating characteristic (ROC) curve. The

151    Youden index produced overlapping cutoffs for Whites (0.44) and Blacks (0.48) where some

152    individuals would be predicted to belong to both White and Black race. Therefore, discriminative

153    thresholds of at least 0.5 were used to assign patients to a single race/ethnic group. Sensitivity

154    analyses between a 0.5 and 0.75 cutoff were conducted based on previous studies [17, 26]. A

155    total of 98.4% of the study population met the starting cutoff of 0.5 for the assignment of race

156    ethnicity, while only 83.1% met criteria at the 0.75 threshold. In the sensitivity analyses, there

157    was a considerable decrease in sensitivity and specificity when the 0.75 is applied compared to

158    the 0.5. Based on this, the 0.5 cutoff was used to assign BISG race/ethnicity. There were 141

159    records for which BISG imputation was undefined, i.e. none of the race categories reached the

160    0.5 cutoff, these were excluded.  BISG probabilities were not imputed for individuals without a

161    geocoded address (missing address completely or partially, PO boxes). There were 19,045

162    records with both missing APCD self-reported and BISG imputed race and ethnicity. These

163    records were excluded from the initial sample resulting in a total sample of 85,857 individuals.

164           *Enhanced APCD race/ethnicity*

165    BISG imputed probabilities were used to assign race/ethnicity for APCD records where it was

166    missing.

167           2.3 Analysis variables

8

168    To assess the differential missingness of race and ethnicity in the APCD records, an indicator

169    was created of whether APCD race ethnicity was missing. Patient characteristics included age at

170    death, gender, payer, year of death, comorbidities and the top ten causes of death. Because age at

171    death from the vital statistics was used, the age distribution is skewed towards older ages;

172    therefore, age was grouped as: <65, 65-74, 75-84, 85+. The payer attributed to each patient was

173    the plan that the individual had for most of the year of death. Comorbidities were flagged using

174    the R "comorbidity" package [27, 28]. This program uses the International Classification of

175    Diseases (ICD) ninth and tenth revision codes based on the Quan et al. definitions [29].

176    Comorbidities present in 10% of the sample or more are reported.

177    The top ten causes of death are based on the national "leading causes of death" report [30].

178    Causes of death were identified using ICD-10 from the underlying causes of death including

179    diseases of heart (I00–I09, I11, I13, I20–I51), malignant neoplasms (C00–C97), accidents

180    (unintentional injuries) (V01–X59, Y85–Y86), chronic lower respiratory diseases (J40–J47),

181    cerebrovascular diseases (I60–I69), Alzheimer's disease (G30), diabetes mellitus (E10–E14),

182    influenza and pneumonia (J09–J18), nephritis, nephrotic syndrome and nephrosis (N00–N07,

183    N17–N19, N25–N27), intentional self-harm (suicide) (U03,X60–X84,Y87.0). Of note, during the

184    study period, there were no deaths directly attributed to Alzheimer's disease in this sample.

185        2.3 Analysis

186    Patient characteristics were tabulated by status of missing race and ethnicity in the APCD and

187    this distribution was tested statistically with the Chi Square tests. Multivariable logistic

188    regression was used to assess factors associated with increased odds of missing race and

189    ethnicity in the APCD. Validation of the various race/ethnicity sources and improvements in

190    sensitivity were compared against the gold standard vital statistics.

191     SAS Studio 9.3 (SAS Institute, Cary, NC) and RStudio version 1.3 were used for all analyses.

192     This study was not deemed Human Subject Research by the Partners Healthcare Human

193     Research Committee (protocol 2018P001185/PHS) and patient consent was not collected.

194

195     **Results**

196     The analytical file included 85,857 APCD records of which 32.1% (n=27,594) had missing

197     APCD race/ethnicity. Individuals in the sample were mostly aged 65 years and older (68.5%),

198     female (52.2%), White race based on vital statistics (91.8%) and were Medicare insured at the

199     time of death (70.2%). Almost 30% of individuals had uncomplicated hypertension.

200     Uncomplicated diabetes, chronic pulmonary disease and cardiac arrhythmia were each present in

201     about 20% of the sample. The leading cause of death was malignant neoplasms (21.8%) followed

202     by heart disease (15.9%) and chronic lower respiratory diseases (6.2) (Table 1).

203     [Place Table 1 here]

204     Missing APCD race and ethnicity data were not randomly distributed across patient

205     characteristics. The logistic regression results showed multiple factors that were significantly

206     associated with the odds of missing APCD self-reported race and ethnicity. For instance, a

207     gradual increase in the frequency of missing race and ethnicity with increasing age was

208     observed. Those in the 85 years of age and older age category had almost 4-fold increased odds

209     of missing race and ethnicity (adjusted odds ratio (aOR) 3.74, 95% confidence level (CI) 3.51-

210     3.98) compared to patients less than 65 years of age. Males were also more likely to be missing

211     race and ethnicity (aOR 1.46, 95% CI 1.41-1.51) compared to females. Compared to vital

212     statistics Whites, all other racial and ethnic groups had lower odds of missing APCD race and

10

213    ethnicity. Missing race and ethnicity was more common in more recent years.  Compared to

214    Medicaid enrollees, individuals with commercial insurance had a considerably higher likelihood

215    of missing APCD race and ethnicity (aOR 43.8, 95% CI 40.17-47.77).  The likelihood of missing

216    race and ethnicity in the APCD differed significantly for patients with different comorbidities.

217    Some conditions had increased likelihood such as uncomplicated hypertension, cardiac

218    arrhythmia, solid tumors and renal failure, while others had lower likelihood such as diabetes,

219    chronic pulmonary disease, and depression. Individuals who died from suicide had a two-fold

220    increased odds of having missing APCD race and ethnicity (aOR 2.10, 95% CI 1.84-2.40), while

221    those who died of diabetes had similar odds (OR 0.91, 95% 0.79-1.04) (Table 1).

222    Table 2 contrasts the distribution and changes in sensitivity and specificity of the APCD race and

223    ethnicity variables before and after enhancement with BISG imputation compared to the "gold

224    standard" vital statistics race and ethnicity. Enhancing missing race and ethnicity in the APCD

225    improves sensitivity for Whites (from 61.1% to 97.2%), Blacks (from 81.7% to 90%), Hispanics

226    (from 51.1% to 66.2%) and API (from 62.6% to 78.2%). In contrast, using BISG imputation

227    alone was not sufficient to capture race and ethnicity for this sample as demonstrated by lower

228    sensitivity and especially specificity for Blacks (26.8%) (Table 2).

229    [Place Table 2 here]

230    **Discussion**

231    This is the first study to our knowledge to examine the ability of the BISG algorithm to

232    accurately impute race and ethnicity when missing in an APCD. Findings suggest that

233    missingness of race and ethnicity variables in the APCD is common. Moreover, missingness

234    does not occur equally; it is more likely among males, Whites, older and commercially insured

235    individuals. There were also significant differences between missing and non-missing records in

236    representation of comorbid conditions and mortality cause. This study found that imputing the

237    missing race/ethnicity with the BISG method greatly improved the sensitivity to detect White,

238    Black, API, and Hispanic groups.

239    Missing race/ethnicity in administrative data sets is a known issue [18, 26, 31]. The APCD data

240    set, by including most of a state's population, can be valuable to answer many important health

241    questions. However, with a high proportion of the data set missing race/ethnicity, the ability of

242    APCD-based analyses to assess health outcomes is limited. Enhancing APCD data with indirect

243    imputation based on surname and address can allow for more robust population health studies.

244    This is particularly important as excluding records with missing race/ethnicity may introduce

245    bias, especially if that missingness is not random as observed in this study. For instance,

246    Grundmeier and colleagues simulated not-randomly missing race/ethnicity data to compare the

247    association of race/ethnicity with pediatric health outcomes when race/ethnicity is imputed vs.

248    excluded.  They found that imputing missing data with BISG reduces bias compared to when

249    only data with non-missing values are used [32]. This study found similar results.

250    Compared to other validation studies, findings from this study show that BISG alone was less

251    optimal to estimate race and ethnic categories correctly, particularly for Blacks (low sensitivity

252    and specificity). For example, Adjaye-Gbewonyo et al. found that BISG captured 71.8% of

253    Blacks while in this study the sensitivity was less optimal at 58% only [33]. This is likely driven

254    by the racial composition of the sample studied. This study is drawn from a mostly White state

255    (Oregon) where Blacks represent a very slim share of the population (2.2%), whereas Adjaye-

256    Gbewonyo and colleagues performed BISG validation on a population form Georgia, a state with

257    a higher than average Black population (32.6%) [34].

258    The lower prediction ability of BISG in a less diverse population is likely due to relying on the

259    Surname Census list. Hispanic and Asian individuals may have more distinctive surnames, which

260    is not the case for surnames of Whites and Blacks. For instance, "Williams", the third most

261    common surname in the U.S. [13], has close to 50/50 probability of either White or Black race.

262    Therefore, in a predominantly White population, the ability of the Surname list to correctly

263    predict Whites vs Blacks may be diminished. BISG could be further improved by incorporating

264    information from the "first name" into the prediction algorithm, particularly for Blacks. Ioan

265    Voicu proposed the Bayesian Improved First Name Surname Geocoding (BIFSG) using a list of

266    first names from mortgage applications [35]. While the author noted amelioration in the

267    predictions compared to BISG, further work is needed using a more representative first name list

268    other than those of mortgage applicants.

269    Another potential reason for the lower BISG prediction of non-White race/ethnic groups is the

270    use of geographic distribution of race/ethnicity which may be influenced by the level of

271    residential segregation in communities. Higher level of segregation is likely to result in more

272    precise prediction of the race/ethnic groups. The state of Oregon is somewhat more integrated

273    than other states [36]. Having diverse communities might make it harder for the BISG algorithm

274    to attribute deterministic probabilities for the racial make-up of a community resulting in less

275    accurate predictions. This study showed that using the BISG method to "enhance" the existing

276    race/ethnicity information is more suitable than relying on BISG predictions alone especially in

277    less diverse and more integrated communities.

278        *Limitations*

279    This study has a number of limitations. The use of death data means that the sample for this

280    study skews older than the Oregon population, in which the race and ethnic distributions may be

281     different. Using a more representative data source with respect to age and that includes a more

282     reliable capture of AIAN race is warranted. Because of data use agreement restrictions, the

283     voluntary APCD used in this study did not include Medicare fee-for-service patients, however,

284     this may have attenuated the skewness towards older patients. Not all individuals had an address

285     that could be properly geocoded. Finally, given the less diverse Oregon population reflected in

286     the high percent White and the low discriminative power of surname probabilities for White and

287     Black surnames, we decided to drop any missing BISG imputed probabilities. Future studies may

288     add the first name to improve the prediction for Blacks.

289

290     **Conclusion**

291     Health equity requires accurately assessing the burden of disease for the different race and ethnic

292     groups to reduce disparities and appropriately allocate resources and interventions. Given the

293     growing use of APCDs and their potential for developing population-based interventions, it is

294     critical to develop methods that allow for the non-biased estimation of race/ethnicity. Having

295     complete race/ethnicity information in administrative data sets would allow to perform more

296     robust population-health level analyses to identify inequities according to race/ethnicity, without

297     losing power or dropping non-random records with missing race data.

298

299

300    **Abbreviations:**

301    AIAN: American Indian Alaska Native; APCD: All Payer Claims Data; API: Asian Pacific

302    Islander; BISG: Bayesian Improved Surname and Geolocation; FIPS: Federal Information

303    Processing Standard; ICD: International Classification of Diseases; OR: Odds ratio

304

## References

1.      Porter J, Love D, Costello A, Peters A, Rudolph B. All-Payer Claims Database Development Manual: Establishing a Foundation for Health Care Transparency and Informed Decision Making. Institute for Health Policy and Practice (IHPP); 2015.

2.      APCD Council. Interactive State Report Map (2020). updated 2021. Available from: https://www.apcdcouncil.org/state/map.

3.      Porter J, Love D. The ABCs of APCDs: How States Are Using Claims Data to Understand and Improve Care. California Health Care Foundation; 2018.

4.      APCD Council, National Association of Health Data Organizations, New Hampshire Institute for Health Policy and Practice. Informing Health System Change - Use of All-Payer Claims Databases. 2018.

5.      US Department of Health Human Services. 2016 National Healthcare Quality and Disparities Report: Agency for Healthcare Research and Quality Rockville, MD; (2017). Available from: https://www.ahrq.gov/research/findings/nhqrdr/nhqdr16/index.html?utm_source=ahrq&utm_medium=en2&utm_term=&utm_content=2&utm_campaign=ahrq_en8_15_2017.

6.      Landon BE, Normand S-LT, editors. National voluntary consensus standards for ambulatory care: Measurement challenges in small group settings2006: National Quality Forum.

7.      Nerenz DR, Carreon R, Veselovskiy G. Race, ethnicity, and language data collection by health plans: findings from 2010 AHIPF-RWJF survey. J Health Care Poor Underserved. 2013;24(4):1769-83.

8.      Smedley BD, Stith AY, Nelson AR. Institute of Medicine (US) Committee on Understanding and Eliminating Racial and Ethnic Disparities in Health Care.(2003). Unequal treatment: Confronting racial and ethnic disparities in health care. 2003.

9.      Escarce JJ, Carreon R, Veselovskiy G, Lawson EH. Collection of race and ethnicity data by health plans has grown substantially, but opportunities remain to expand efforts. Health Aff (Millwood). 2011;30(10):1984-91.

10.     Patient Protection and Affordable Care Act. §3022 and §10307, Pub L No. 111-148, 124 Stat 119. 2010.

11.     Ng JH, Ye F, Ward LM, Haffer SC, Scholle SH. Data On Race, Ethnicity, And Language Largely Incomplete For Managed Care Plan Members. Health Aff (Millwood). 2017;36(3):548-52.

12.     Elliott MN, Fremont A, Morrison PA, Pantoja P, Lurie N. A new method for estimating race/ethnicity and associated disparities where administrative records lack self-reported race/ethnicity. Health Serv Res. 2008;43(5 Pt 1):1722-36.

13.     Comenetz J. Frequently occurring surnames in the 2010 census. 2016.

14.     Fremont A, Weissman JS, Hoch E, Elliott MN. When race/ethnicity data are lacking: using advanced indirect estimation methods to measure disparities. Rand Health Q. 2016;6(1):16.

15.     Breslau J, Elliott MN, Haviland AM, Klein DJ, Dembosky JW, Adams JL, et al. Racial And Ethnic Differences In The Attainment Of Behavioral Health Quality Measures In Medicare Advantage Plans. Health Aff (Millwood). 2018;37(10):1685-92.

16.     Elliott MN, Morrison PA, Fremont A, McCaffrey DF, Pantoja P, Lurie N. Using the Census Bureau's surname list to improve estimates of race/ethnicity and associated disparities. Health Serv Outcomes Res Methodol. 2009;9(2):69-83.

17.     Adjaye-Gbewonyo D, Bednarczyk RA, Davis RL, Omer SB. Using the bayesian improved surname geocoding method (BISG) to create a working classification of race and ethnicity in a diverse managed care population: a validation study. Health Serv Res. 2014;49(1):268-83.

18.     Haas A, Elliott MN, Dembosky JW, Adams JL, Wilson-Frederick SM, Mallett JS, et al. Imputation of race/ethnicity to enable measurement of HEDIS performance by race/ethnicity. Health Serv Res. 2019;54(1):13-23.

351    19.      Arias E, Heron MP, Hakes JK. The validity of race and Hispanic origin reporting on death
352    certificates in the United States: an update. 2016(148):1-23.
353    20.      Enamorado T, Fifield B, Imai K. Using a probabilistic model to assist merging of large-scale
354    administrative records. Available at SSRN 3214172. 2018;113(2):353-71.
355    21.      U.S. Census Bureau. Geocoding Services Web Application Programming Interface (API) (2019).
356    cited 2020. Available from: https://geocoding.geo.census.gov/geocoder/Geocoding_Services_API.pdf.
357    22.      Imai K, Khanna K. Improving ecological inference by predicting individual ethnicity from voter
358    registration records. Polit Anal. 2016;24(2):263-72.
359    23.      Word DL, Coleman CD, Nunziata R, Kominski R. Demographic aspects of surnames from census
360    2000. Unpublished manuscript, Retrieved from http://citeseerx ist psu edu/viewdoc/download. 2008.
361    24.      Elliott MN, Morrison PA, Fremont A, McCaffrey DF, Pantoja P, Lurie N. Using the Census
362    Bureau's surname list to improve estimates of race/ethnicity and associated disparities. Health Services
363    and Outcomes Research Methodology. 2009;9(2):69-83.
364    25.      Fluss R, Faraggi D, Reiser B. Estimation of the Youden Index and its associated cutoff point.
365    Biometrical Journal: Journal of Mathematical Methods in Biosciences. 2005;47(4):458-72.
366    26.      Derose SF, Contreras R, Coleman KJ, Koebnick C, Jacobsen SJ. Race and ethnicity data quality and
367    imputation using U.S. Census data in an integrated health system: the Kaiser Permanente Southern
368    California experience. Med Care Res Rev. 2013;70(3):330-45.
369    27.      Gasparini A. Comorbidity: an R package for computing comorbidity scores. Journal of Open
370    Source Software. 2018;3(23):648.
371    28.      Elixhauser A, Steiner CA, Harris D, Coffey R. Comorbidity measures for use with administrative
372    data. Med Care. 1998(36):8-27.
373    29.      Quan H, Sundararajan V, Halfon P, Fong A, Burnand B, Luthi J-C, et al. Coding algorithms for
374    defining comorbidities in ICD-9-CM and ICD-10 administrative data. Med Care. 2005;43(11):1130-9.
375    30.      Heron MP. Deaths: leading causes for 2017. 2019.
376    31.      Brown DP, Knapp C, Baker K, Kaufmann M. Using Bayesian Imputation to Assess Racial and
377    Ethnic Disparities in Pediatric Performance Measures. Health Serv Res. 2016;51(3):1095-108.
378    32.      Grundmeier RW, Song L, Ramos MJ, Fiks AG, Elliott MN, Fremont A, et al. Imputing missing
379    race/ethnicity in pediatric electronic health records: reducing bias with use of US census location and
380    surname data. Health Serv Res. 2015;50(4):946-60.
381    33.      Adjaye-Gbewonyo D, Bednarczyk RA, Davis RL, Omer SB. Using the Bayesian Improved Surname
382    GeocodingMethod (BISG) to Create a Working Classification ofRace and Ethnicity in a Diverse Managed
383    Care Population: AValidation Study. Health Services Research. 2014;49(1):268-83.
384    34.      Census Bureau. Population estimates, July 1, 2019 (2019). cited 2020 September. Available
385    from: https://www.census.gov/quickfacts/fact/table/US/PST045219.
386    35.      Voicu I. Using first name information to improve race and ethnicity classification. Statistics and
387    Public Policy. 2018;5(1):1-13.
388    36.      Oregon Health Authority. Social Determinants of Health: Residential racial segregation (2019).
389    Available from: https://www.oregon.gov/oha/PH/ABOUT/Documents/indicators/segregation.pdf.

390

391     **Statements and Declarations:**

392         **Funding:**

393     This work was funded by NIH grant 1-R01-DA044167. The funding source was not involved in

394     the study design, collection, analysis and interpretation of data, nor in the writing of the report or

395     the decision to submit the report for publication. The conclusions in this article are those of the

396     authors and do not necessarily represent the official position of the funder. The authors declare

397     that no funds, grants, or other support were received during the preparation of this manuscript

398         **Competing Interests**

399     The authors have no relevant financial or non-financial interests to disclose

400         **Authors' contributions**

401         S. E. conceptualize, analyzed the data and wrote the manuscript, M. H. reviewed the

402     analysis and the manuscript, S. H. reviewed the analysis and the manuscript, N. H. prepared the

403     data for analysis, C. H. reviewed the manuscript, M. F. reviewed the manuscript, S. W.

404     conceptualized and reviewed the manuscript. All authors reviewed the manuscript.

405         **Ethics approval and consent to participate:**

406     This study was not deemed Human Subject Research by the Partners Healthcare Human

407     Research Committee (protocol 2018P001185/PHS) and patient consent was not collected.

408         Consent for publication

409     All authors consent to publish this manuscript.

410

411        **Table 1. Patient characteristics and odds ratios of missing self-reported APCD race and ethnicity**

| Patient characteristics | Total | Non-missing APCD race & ethnicity | Missing APCD race & ethnicity | Odds of missing race & ethnicity |
|---|---|---|---|---|
| | n (%) | n (%) | n (%) | OR (95% CI) |
| **Total sample** | 85857 (100) | 58263 (67.9) | 27594 (32.1) | |
| **Age at death** | | | | |
| <65 | 27035 (31.5) | 21263 (36.5) | 5772 (20.9) | Ref |
| 65-74 | 15975 (18.6) | 10747 (18.5) | 5228 (19.0) | 2.22 (2.08-2.37) |
| 75-84 | 18115 (21.1) | 11439 (19.6) | 6676 (24.2) | 3.04 (2.85-3.25) |
| 85+ | 24732 (28.8) | 14814 (25.4) | 9918 (35.9) | 3.74 (3.51-3.98) |
| **Gender** | | | | |
| Female | 44804 (52.2) | 31397 (53.9) | 13407 (48.6) | Ref |
| Male | 41053 (47.8) | 26866 (46.1) | 14187 (51.4) | 1.46 (1.41-1.51) |
| **Vital race and ethnicity** | | | | |
| White | 78855 (91.8) | 52732 (90.5) | 26123 (94.7) | Ref |
| Black | 1575 (1.8) | 1363 (2.3) | 212 (0.8) | 0.889 (0.793-0.99) |
| Hispanic | 2319 (2.7) | 1766 (3.0) | 553 (2.0) | 0.58 (0.51-0.65) |
| API | 1727 (2.0) | 1266 (2.2) | 461 (1.7) | 0.41 (0.35-0.48) |
| Other | 1381 (1.6) | 1136 (2.0) | 245 (0.9) | 0.62 (0.53-0.73) |
| **Year of death** | | | | |
| 2013 | 9463 (11.3) | 7260 (12.7) | 2203 (8.3) | Ref |
| 2014 | 11762 (14.1) | 8285 (14.5) | 3477 (13.1) | 1.41 (1.31-1.51) |
| 2015 | 13434 (16.1) | 9524 (16.7) | 3910 (14.7) | 1.41 (1.32-1.51) |
| 2016 | 14727 (17.6) | 9579 (16.8) | 5148 (19.4) | 1.81 (1.70-1.94) |
| 2017 | 17102 (20.4) | 10999 (19.3) | 6103 (23.0) | 1.79 (1.68-1.91) |
| 2018 | 17196 (20.6) | 11453 (20.1) | 5743 (21.6) | 1.56 (1.47-1.67) |
| **Payer [†]** | | | | |
| Commercial | 6098 (7.3) | 1902 (3.3) | 4196 (15.8) | 43.76 (40.17-47.77) |
| Medicaid | 18873 (22.6) | 17822 (31.2) | 1051 (4.0) | Ref |
| Medicare | 58713 (70.2) | 37376 (65.5) | 21337 (80.3) | 5.27 (4.91-5.67) |
| **Elixhauser comorbidities, prevalence 10%+ [‡]** | | | | |
| Hypertension, uncomplicated | 25235 (29.4) | 15341 (26.3) | 9894 (35.9) | 1.34 (1.29-1.39) |
| Diabetes, uncomplicated | 18610 (21.7) | 13171 (22.6) | 5439 (19.7) | 0.83 (0.79-0.87) |
| Chronic pulmonary disease | 17154 (20.0) | 12433 (21.3) | 4721 (17.1) | 0.76 (0.73-0.80) |
| Cardiac arrhythmias | 16799 (19.6) | 10176 (17.5) | 6623 (24.0) | 1.37 (1.32-1.43) |
| Congestive heart failure | 14142 (16.5) | 9619 (16.5) | 4523 (16.4) | 0.84 (0.80-0.88) |
| Solid tumor, without metastasis | 11540 (13.4) | 6569 (11.3) | 4971 (18.0) | 1.50 (1.43-1.58) |
| Diabetes, complicated | 10018 (11.7) | 7350 (12.6) | 2668 (9.7) | 0.76 (0.72-0.81) |
| Depression | 9614 (11.2) | 7279 (12.5) | 2335 (8.46) | 0.74 (0.70-0.78) |
| Fluid and electrolyte disorders | 9495 (11.1) | 6570 (11.3) | 2925 (10.6) | 0.89 (0.84-0.94) |
| Renal failure | 8599 (10.0) | 5552 (9.5) | 3047 (11.0) | 1.19 (1.13-1.26) |
| **Top 10 causes of death [‡, §]** | | | | |
| Malignant neoplasms | 18744 (21.8) | 11232 (19.3) | 7512 (27.2) | 1.62 (1.55-1.70) |
| Heart disease | 13645 (15.9) | 8581 (14.7) | 5064 (18.4) | 1.37 (1.31-1.44) |

| | | | | |
|---|---|---|---|---|
| Chronic lower respiratory diseases | | | 1248 (4.5) | |
| | 5276 (6.2) | 4028 (6.9) | | 0.94 (0.87-1.02) |
| Accidents | 5099 (5.9) | 3692 (6.3) | 1407 (5.1) | 1.34 (1.24-1.45) |
| Cerebrovascular diseases | 3628 (4.2) | 2381 (4.1) | 1247 (4.5) | 1.21 (1.12-1.31) |
| Intentional self-harm (suicide) | 2045 (2.4) | 1347 (2.3) | 698 (2.5) | 2.10 (1.84-2.40) |
| Diabetes mellitus | 1461 (1.7) | 1129 (1.9) | 332 (1.2) | 0.91 (0.79-1.04) |
| Nephritis, nephrotic syndrome and nephrosis | | | 268 (1.0) | |
| | 921 (1.1) | 653 (1.1) | | 1.01 (0.86-1.18) |
| Influenza and pneumonia | 693 (0.8) | 486 (0.8) | 207 (0.8) | 1.14 (0.95-1.36) |

412    CI: confidence interval; OR: odds ratio
413    † Payer corresponds to the payer on the year of death
414    ‡ Reference = no event
415    § There were no deaths associated with Alzheimer's in the study period

416

417

418    **Table 2. Improvement in the sensitivity and specificity of race and ethnicity after enhancement with**
419    **the BISG**

| Race and Ethnicity | Vital statistics[†] | | APCD | | | | BISG imputation | | | | Enhanced APCD | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | n | % | n | % | Sensitivity | Specificity | n | % | Sensitivity | Specificity | n | % | Sensitivity | Specificity |
| White | 78,855 | 92 | 52,457 | 61 | 65.2 | 98 | 63,890 | 74 | 78.4 | 98 | 78,308 | 91 | 97.2 | 97.9 |
| Black | 1,575 | 1.8 | 1,403 | 1.6 | 81.7 | 91.7 | 3,400 | 4 | 58 | 26.8 | 2,160 | 2.5 | 90 | 65.7 |
| Hispanic | 2,319 | 2.7 | 2,082 | 2.4 | 51.1 | 56.9 | 2,182 | 2.5 | 56.1 | 56.9 | 2,662 | 3.1 | 66.2 | 57.7 |
| API | 17,27 | 2 | 1,231 | 1.4 | 62.6 | 87.1 | 1,346 | 1.6 | 60.5 | 77.6 | 1,599 | 1.9 | 78.2 | 84.4 |
| Other [‡] | 1,381 | 1.6 | 1,090 | 1.3 | | | 180 | 0.2 | | | 1,128 | 1.3 | | |
| Missing | | | 27,594 | 32 | | | 14,859 | 17 | | | | | | |

420
421    APCD: All Payer claims data; API: Asian/Pacific Islander; BISG: Bayesian Improved Surname Geocoding Method
422    [†] Vital statistics race/ethnicity is considered the Gold Standard
423    [‡] The "other" category is not homogenous across the different race/ethnicity source variables, therefore, sensitivity is
424    not calculated for this category
425