

1 A Poisson binomial based statistical testing framework for
2 comprehensive comorbidity discovery across massive Electronic
3 Health Record datasets

4
5 Gordon Lemmon^{1,2}, Sergiusz Wesolowski^{1,2}, Alex Henrie^{1,2}, Martin Tristani-Firouzi^{3,4},
6 Mark Yandell^{1,2}

7
8 ¹Department of Human Genetics, University of Utah, Salt Lake City, UT, USA;

9 ²Utah Center for Genetic Discovery & Department of Human Genetics, University of
10 Utah, Salt Lake City, UT, USA

11 ³Division of Pediatric Cardiology, University of Utah School of Medicine, Salt Lake City,
12 UT, USA;

13 ⁴Nora Eccles Harrison CVRTI, University of Utah School of Medicine, Salt Lake City,
14 UT, USA;

15
16 *Corresponding Authors: Martin.Tristani@utah.edu, myandell@genetics.utah.edu

17

18 Abstract

19 Discovery of comorbidities (the concomitant occurrence of distinct medical conditions in
20 the same patient) is a prerequisite for creating forecasting tools for downstream outcomes
21 research. Current comorbidity discovery applications are designed for small datasets and
22 use stratification to control for confounding variables such as age, sex, or ancestry.
23 Stratification lowers false positive rates, but reduces power, as the size of the study cohort
24 is decreased. Here, we describe a Poisson Binomial based approach to comorbidity
25 discovery (PBC) designed for big-data applications that circumvents the need for
26 stratification. PBC adjusts for confounding demographic variables on a per-patient basis,
27 and models temporal relationships. We benchmark PBC using two datasets, the publicly
28 available MIMIC-IV; and the entire Electronic Health Record (EHR) corpus of the
29 University of Utah Hospital System, encompassing over 1.6 million patients, to compute
30 comorbidity statistics on 4,623,841 pairs of potentially comorbid medical terms. The
31 results of this computation are provided as a searchable web resource. Compared to
32 current methods, the PBC approach reduces false positive associations, while retaining
33 statistical power to discover true comorbidities.

34 Introduction

35 Comorbidity refers to the concomitant occurrence of distinct medical conditions in the
36 same patient¹. Comorbidities can occur together, or sequentially across the patient's
37 medical history. Exploring these temporal connections offers additional insight into
38 disease progression and disease associations, and promises improved predictive tools
39 for evidence-based medicine²⁻⁵. Traditionally, comorbidities have been discovered

40 manually, through human chart review, literature search, and clinical knowledge⁶. For
41 example, the authors of the Charlson comorbidity index⁷ selected comorbid diagnoses
42 based on manual chart review for a 559-patient cohort. Likewise, the well-known
43 Elixhauser Comorbidity index, was compiled through review of published studies
44 identifying comorbid conditions⁸.

45 Large collections of Electronic Health Records (EHRs) present promising new
46 opportunities for comorbidity discovery. However, manual review of millions of EHR
47 records in search of comorbidities is infeasible, and *ab initio* means for discovery and
48 temporal ordering of comorbidities using large collections of EHRs is an area ripe for
49 innovation. Current computational approaches to *ab initio* comorbidity discovery use
50 statistics such as risk ratio, odds ratio, comorbidity-score⁹, propensity score or ϕ -
51 correlation to measure effect size. P-values are obtained using Fisher's exact test (or
52 hypergeometric), χ^2 test, or binomial test¹⁰⁻¹⁴. All of these approaches rely on an
53 assumption that each member of the population has a disease probability equal to the
54 population incidence rate. Confounding variables such as age and sex are controlled for
55 by sub-setting the data, a process termed *stratification*, or alternatively through use of
56 matched case-control cohorts^{14,15}. Both of these approaches control for confounders, but
57 at the expense of statistical power, because they necessarily reduce sample size. One
58 approach to overcoming this intrinsic limitation is to aggregate massive collections of
59 EHRs, but as we show, even millions of records are too few to explore comorbidities
60 associated with rare diseases when controlling for multiple confounding variables.

61 PBC models the effects of confounding variables, allowing every sample to be
62 personalized, resulting in improved statistical power compared to stratification. Briefly,

63 PBC uses logistic regression to model how each patient's demographics impact his or
64 her probability of having a medical term. These personalized probabilities are then used
65 to calculate pairwise expectations and p-values under the Poisson binomial distribution,
66 rather than the hypergeometric and binomial distributions that are used for Fisher's exact
67 test, and the χ^2 test, respectively. Moreover, with minor modification, this approach can
68 also be used to temporally order comorbidities and to determine the significance of
69 directionalities. As we demonstrate, use of the Poisson binomial is a significant
70 advantage, because it obviates the need for stratification. The result is increased power
71 for discovery, which we leverage to explore the relationships among diagnoses, medical
72 procedures and medications.

73 Alongside the need for improved statistical methods, tools are also needed to browse,
74 search, and visualize the network of comorbid medical terms discovered in big-data
75 applications. In a manner similar to *Siggaard et al.*¹⁴, we provide a browser-based query
76 engine for navigating these comorbidities and their temporal relationships within the
77 University of Utah Hospitals system [[https://pbc.genetics.utah.edu/lemmon2021/pbc-
78 utah](https://pbc.genetics.utah.edu/lemmon2021/pbc-utah)].

79 In what follows, we describe PBC, explore its behavior, and benchmark its performance
80 using the contents of the University of Utah Health system, and the publicly available
81 MIMIC-IV dataset¹⁶. For brevity's sake, we will refer to co-occurring medical diagnoses,
82 procedures and medications using the single blanket term, comorbidity. We demonstrate
83 how PBC can be used to transform massive EHR datasets into a temporal dependency
84 graph for large-scale *ab initio* discovery of comorbid relationships and investigations of
85 disease progressions.

86 Results

87 *Modeling the effects of confounders using logistic regression*

88 We collected records for 1.6 million patients, encompassing 50 million visits and 150
89 million diagnosis (DX), procedure (PX) and medication (RX) codes from the University of
90 Utah Electronic Data Warehouse (EDW). For the proof of principle analyses presented
91 here, diagnoses were converted from ICD9¹⁷ and ICD10¹⁸ diagnosis codes to Clinical
92 Classification System¹⁹ (CCS) multi-level diagnosis codes. CPT²⁰ provider billing codes
93 were converted to CCS multi-level procedure codes. We include both leaf nodes and
94 internal nodes so that the researcher can discover more specific comorbidities (e.g “CCS
95 2.1.1: cancer of colon”) as well as more general comorbidities (e.g. “CCS 2.1: colorectal
96 cancer” or “CCS 2: neoplasms”). Medications were coded using RxNorm concept unique
97 identifiers (CUIs)²¹. This procedure reduced the number of distinct medical terms to 1007
98 DX, 259 PX and 1775 RX codes. We also collected demographic information, including
99 sex, race, ethnicity, insurance class, age and length of medical records.

100 A logistic regression model (LRM) was determined for each DX, PX and RX term. The
101 LRM includes demographic information for each patient (age, gender, ancestry, ethnicity,
102 insurance type) and EHR exposure (the length and density of a person’s medical record).
103 Since medical terms are included or excluded from year to year, and coding practices
104 vary over time, we also include the date of the patient’s last visit as a control for this effect.
105 The complete list of features is described in **Table S1** and **Figure S1**. The response
106 variable is whether each patient has the term in their medical record. Note that we do not
107 model recurrence - in this analysis we consider only the first instance of a term in a

108 patient's medical history. We use these LRMs to estimate for each medical term, each
109 patient's *a priori* personalized probability of having that term in their medical record.

110 We used a regularized regression model under the assumption of collinearity in the
111 confounding features. However, there is no requirement for such an assumption; for
112 example, Neural networks could be used in place of LRM, so as to better capture
113 nonlinear relationships between variables. L1 and L2 penalized logistic regression include
114 a value C that prevents overfitting by penalizing large coefficients. Smaller C -values
115 specify stronger regularization (e.g., stronger prevention of over-fitting). To determine the
116 optimal C -value for each LRM, it was necessary to choose a score function for LRM
117 evaluation. We experimented with a number of standard and custom score functions as
118 described in supplemental "Math.pdf". We optimize C for each LRM using stratified 3-fold
119 cross validation. Grid search is used to evaluate C -values in the set $\{10^{-14}, 10^{-13}, 10^{-12},$
120 $\dots, 10^{12}, 10^{13}, 10^{14}\}$. **Figure S2** top panel shows boxplots of the scores reported by each
121 of these score functions. We use entropy as a measure of the ability of a score function
122 to differentiate model quality. From those score functions achieving high entropy, we
123 evaluate the distribution of C -values (**Figure S2** bottom panel). We choose J_{cutoff} for all
124 downstream analysis because it includes fewer outliers than the other methods
125 examined. J_{cutoff} is based on Youden's J statistic²², only rather than a 50% probability
126 threshold, the classification threshold is determined empirically so that the total number
127 of predicted positives is equal to the actual count of positives.

128 We train LRMs using stratified 3-fold cross validation, evaluated using J_{cutoff} . Under L1
129 penalized logistic regression, model features can be unselected by setting their
130 coefficients to zero. **Figure 1** summarizes how often each demographic feature is

131 included in the trained LRMs. Age at last visit, number of visits, number of terms, and
132 length of medical record were grouped together as “EHR exposure”. Patients may be
133 seen at a non-university clinic, may move in or out of the state over the years, or may
134 have differing proclivities toward visiting the doctor. EHR exposure is an attempt to control
135 for these effects. As **Figure 1** makes clear, EHR exposure is always important in
136 predicting whether a patient has a particular medical term.

137 ***PBC retains power as features are added and as sample sizes are reduced***

138 The binomial distribution models the discrete probability of the number of successes in N
139 independent experiments each with probability P. A naïve approach to comorbidity
140 analysis assumes the probability of seeing term 1 and term 2 (P_{t_1,t_2}) in a medical record
141 is the product of the population incidence rates for terms 1 and 2. Knowing the number of
142 patients with both terms 1 and 2 in their medical record, one can calculate a comorbidity
143 p-value using the binomial test of statistical significance. However, because this method
144 does not adjust for demographic factors, the p-values they generate can be driven by
145 effects such as age and sex. Thus, a common approach in comorbidity literature is to
146 stratify the population by age and sex, and then calculate the binomial p-value for each
147 stratum⁷. In contrast, the PBC approach uses the Poisson binomial distribution to
148 calculate p-values for each term pair. The Poisson binomial distribution is a generalization
149 of the binomial distribution in which every trial/sample (i.e., patient) has a different
150 probability of “success” (i.e., having both terms in their medical record). The probability of
151 an individual patient having both terms in their medical record is calculated as the product
152 of per-patient per-term probabilities generated from corresponding LRMs described
153 above.

154 **Figure 2** explores the relationships between comorbidities discovered by PBC versus
155 stratification as a function of increasing numbers of demographic features. As can be
156 seen, the stratification approach rapidly loses power as more criteria are added to the
157 stratum filter. In contrast, by modeling demographic features, PBC maintains power.

158 **Table 1** compares the power of stratification and PBC to detect three well known
159 comorbidities using our EHR corpus. As in **Figure 2**, we compare the p-values generated
160 by sequentially adding confounding variables. Notice how stratification dramatically
161 lowers the strength of p-values as a function of the size of the stratum. This same behavior
162 is illustrated globally in **Figure 2**. Even with millions of EHRs, controlling for more than a
163 few confounding demographic variables leads to strata that are too small to achieve
164 statistical significance. By modeling the effects of multiple confounding variables, PBC
165 retains statistical power to identify comorbidities.

166 **Table S3** presents five well known comorbidities of breast cancer^{23–28}. Stratification by
167 age and gender deflates the strength of all p-values, and as a result they fall below the
168 Bonferroni corrected significance threshold (**Table S3** column “binomial female 50-59”).
169 In contrast, by explicitly modeling age, gender, race, ethnicity, insurance type, and EHR
170 exposure, PBC retains power to capture these true positive associations.

171 The complete set of comorbid term pairs discovered by PBC can be visualized using
172 network analysis. In supplemental **Figure S3**, we use the minimum description length
173 algorithm²⁹ to perform clustering of pairwise comorbidities by p-value strength. Terms with
174 similar patterns of comorbidities are closer together in the network. To illustrate this, we
175 annotated selected comorbidities discovered within each cluster. A literature search

176 confirmed that these labeled comorbidities represent existing clinical knowledge (see
177 corresponding citations in **Table S7**). **Figure S3** provided the motivation to produce an
178 interactive tool for querying, exploring and extracting information from the comorbidity
179 network. In order to provide better means to navigate this complex network we developed
180 a browser-based tool for exploration of comorbidities, discussed below.

181 ***PBC identifies comorbidities unique to underrepresented minority groups***

182 PBC can also capture true comorbidities hidden within mixed populations. For instance,
183 consider Sickle Cell Anemia, a disease that affects 1 in 365 African American newborns
184 and 1 in 100,000 newborn Caucasians in the United States³⁰. Malaise and fatigue, while
185 common to many disorders, are among the most common symptoms of Sickle Cell
186 Anemia (SCA)^{31–33}, and usually manifest in an age dependent manner. **Table 2** compares
187 five comorbidity p-values. Row 1 presents p-values calculated using all data. While the
188 true comorbidity is discovered, we cannot say with certainty if the relationship is a true
189 comorbidity or simply driven by a third confounding variable (e.g. ancestry). The following
190 rows present X2 p-values after stratifying by ancestry, ethnicity, gender, and age and
191 PBC p-values after including these same features in the regression model. Stratification
192 fails to find a significant comorbidity between SCA and malaise/fatigue once the data is
193 partitioned by ancestry. The rarity of Sickle Cell disease in Caucasians and the small
194 sample size for African Americans within Utah's EHR corpus make detection of this
195 comorbidity difficult. PBC, in contrast, discovers the comorbidity. In fact, the additions of
196 ancestry, ethnicity, and age each increase the strength of the association between SCA
197 and malaise/fatigue.

198 ***Application of PBC to a publicly available dataset***

199 Because the University of Utah data used in this study cannot be shared publicly, as it
200 includes protected health information (PHI), we also demonstrated the general
201 applicability of PBC by applying it on the publicly available MIMIC-IV dataset¹⁶. This
202 dataset includes 248,714 patients with associated ICD10 diagnosis or procedure codes.
203 A total of 5,363,338 ICD9 and ICD10 diagnosis and procedure codes were converted to
204 multi-level CCS codes. These include 725 distinct CCS diagnosis codes and 395 distinct
205 CCS procedure codes. We repeated our experiments on this dataset, training a logistic
206 regression model for each diagnosis and procedure code, and calculating comorbidities
207 using the Poisson binomial distribution. Although the MIMIC-IV data is missing much of
208 the detail available in the University of Utah dataset, the PBC results generally mirror
209 those of the Utah dataset (**Figure S4**). Regression features are shown in **Figure S4**, top
210 panel. In addition, because a random offset is added to each patient's admission dates,
211 it is not possible to control for the changes in use of various billing codes over time.
212 Despite these limitations, the deployment of PBC on the MIMIC-IV public dataset further
213 illustrates how PBC retains statistical power to identify comorbid relationships that are
214 lost by stratification. Tables 1 and 2 are replicated on MIMIC-IV data as supplemental
215 tables S4 and S5. Co-occurrence and directional comorbidities discovered within MIMIC-
216 IV data can be queried at the following link:

217 <https://pbc.genetics.utah.edu/lemmon2021/pbc-mimic>.

218 ***Temporalized P-values allow for understanding of disease progression***

219 The PBC approach can be extended to provide temporalized (or directional) p-values
220 across pre-specified time windows (see Methods for details). The inclusion of a direction
221 window is necessary for several reasons. On short time scales, the order of appearance
222 of diagnostic codes is an unreliable indicator of which condition actually preceded the
223 other in the patient. The development of underlying disease, the relevant signs and
224 symptom, the provider arriving at a given diagnosis, and the eventual recording of the
225 said diagnosis might follow staggered paths that have little or no relevance when viewing
226 the data in a time-slice of less than 30 days or so, thus for many analyses it is probably
227 best to treat these events as contemporaneous. However, in some cases a short window
228 size is optimal for capturing a comorbidity.

229 Consider the following example. PBC reports the following p-values for the
230 diagnosis/procedure pair amputation of lower extremity → postoperative infection: within
231 30 days = $1e-3663$, greater than 30 days = $1e-24$, greater than 90 days = $1e-6$, greater
232 than 365 days = 0.86. It is clear that shorter window sizes better capture the increased
233 risk of infection after amputation, as one would logically expect. Thus, the window size
234 must be informed by clinical knowledge and research objectives. In this manuscript, we
235 examined associations based on a 90-day window. Using our website (see below) a user
236 can also query additional window sizes (30 days, 365 days and 730 days).

237 **Table S6** presents specific directional associations discovered in an *ab initio* fashion by
238 PBC and supported by clinical knowledge. For instance, Milrinone is prescribed to
239 patients awaiting heart transplant³⁴. The tendency of type 2 diabetics to develop chronic

240 kidney disease is well known^{35,36}. HIV induced immunocompromisation often leads to
241 pneumocystis³⁷ and obesity is a known risk factor for hypertension³⁸.

242 ***A web-based resource for comorbidity research***

243 Among 4,623,841 pairs of medical terms in our collection of 1.6 million EHRs, we
244 identified 3,311,830 comorbidities co-occurring within a 90-day window, and 1,969,941
245 temporally directed ones, acting over a time period greater than 90 days. All associations
246 meet a Bonferroni significance threshold of 1.08e-08. The result is a highly-connected
247 network of comorbid diagnoses and associated procedures and medications based on
248 the University of Utah EHR database.

249 In response to the size and complexity of these 'big-data', we have created browser-
250 based means to navigate, query and explore them
251 (<https://pbc.genetics.utah.edu/lemmon2021/pbc-utah>). A screenshot highlighting the
252 functionality of the browser is shown in **Figure 3**. The site allows the user to search for
253 relationships between any pairwise diagnosis, procedure or medication. The result is a
254 searchable table of all other DX, PX, and RX codes along with statistics about the
255 connection to the query term. These statistics include the counts, expectation and p-value
256 of association after adjusting for confounders shown in **Figure 1**. We use two-sided p-
257 values, so that "less than expected" associations are also discoverable. For comparison,
258 we provide both χ^2 p-values and G-test p-values.

259 In addition to patient lifetime co-occurrence p-values, we provide within window, and out
260 of window directional p-values with a selectable window size (30, 90, 365, or 730 day).
261 We also provide statistics on effect size, including relative risk, odds ratio, and "flow rate"

262 which is the percent of patients coded with term 1 who later are coded with term 2. The
263 “Flowchart” view puts the query term in focus and shows the terms that tend to precede
264 and follow the query term to the left and right respectively. The user can filter by p-value
265 strength and by flow rate.

266 A slice of that network is shown in **Figure 3**. This figure contains a “flow chart” view for
267 essential hypertension. By switching out the central node, an investigator can step
268 through the temporalized network of diagnoses, procedures and medications. The
269 investigator can further filter by effect size to find co-occurrences that are both significant
270 and prevalent.

271 ***Comparison with other published comorbidity tools***

272 **Table 3** compares functionalities provided on our comorbidity website with those found in
273 other published comorbidity discovery tools. To the best of our knowledge, our website is
274 the only available public resource of its kind that considers the individual risk profile for
275 each patient having each medical term. Additionally, our approach (and website) captures
276 inverse comorbidities (terms occurring together less often than expected) and models
277 temporal relations between pairs of terms.

278 The R package “comoRbidity” was published April 2018¹⁰. We installed the package, and
279 reformatted our data to fit the required specifications. Given our input of 150,598,377 CCS
280 and RxNorm codes, the “comoRbidity” package consumes all available RAM (we used a
281 Linux server with 504 GB of RAM) and fails to complete. We tried to acquire CytoCom¹¹
282 and comoR¹², however the corresponding author has indicated that these projects are no
283 longer maintained nor available for download. The Java package “Comorbidity4j” was

284 published January 2019¹³. We installed Comorbidity4j and attempted to compute on our
285 full dataset. The application warned against calculating comorbidities for more than
286 300,000 pairs of terms (774 distinct terms). After several hours the calculation times out
287 - having allocated 50.2 GB of RAM. In contrast our method uses a maximum of 174 MB
288 of RAM and has no upper limit on the number of patients, visits, or unique terms.

289 The Disease Trajectory Browser (DTB) allows navigation of temporal relationships among
290 medical records of 7.2 million Danish patients¹⁴. Direct comparison between our results
291 is not possible since our groups have access to different EHR datasets, however we can
292 consider differences in methodology. DTB measures effect size using relative risk,
293 significance is measured using the binomial test, and confounders are controlled for using
294 case/control matching. For comparison, on PBC-Web, we provide relative risk estimates
295 and χ^2 p-values (which are a close approximation to the binomial test p-values). As
296 described above, our approach models patient features rather than controlling for them
297 through stratification or case/control matching.

298 Discussion

299 Although the term comorbidity is often used to denote significant associations between
300 medical outcomes, (e.g., hypertension and heart attack), the concept is easily extended
301 to include associated variables, such as medical procedures and medications. For
302 brevity's sake, in what follows, we refer to statistically significant associations among
303 these collective variables as comorbidities.

304 Comorbidity discovery is a feature discovery/selection process, and it is important to
305 distinguish it from outcomes prediction. Before one can understand and predict a medical
306 outcome, one must first decide which prior diagnoses, medical procedures and
307 medications are germane to the outcome of interest. Discovery of comorbidities is thus a
308 prerequisite for downstream outcomes research and for creating forecasting tools³⁹⁻⁴²,
309 but is logically separate from them.

310 Big data offer many opportunities and challenges for comorbidity discovery. One limitation
311 imposed by data size is that morbidity discovery is necessarily pairwise; hence the term
312 comorbidity, as opposed to multimorbidity discovery. Tools for multimorbidity-based
313 discovery⁴³ are necessarily limited in scale due to computational constraints, considering
314 for instance, 34 disease clusters⁴⁴. In contrast, we have calculated pairwise comorbidities
315 among 37,997 ICD10 diagnosis codes.

316 Commonly used statistical approaches to *ab initio* comorbidity discovery are hindered by
317 the assumption that every member of the population (or stratum) has a disease probability
318 equal to the population (or stratum) incidence rate. As we have explained, stratification is
319 commonly used to subset EHR collections to meet this requirement; but stratification
320 necessarily reduces sample size and statistical power (c.f. **Figure 2**). In practice,
321 stratified data quickly become limiting, even for very large datasets, as inclusion criteria
322 grow more complex. This problem is exacerbated for *ab initio* approaches aimed at
323 simultaneous discovery of comorbid relationships among thousands of diagnoses,
324 procedures and medications. This process necessitates many millions of statistical tests;
325 the requirement for multiple testing corrections mean that statistical power is of
326 paramount importance.

327 Our motivation in developing PBC was to overcome the need for stratification, while still
328 achieving high accuracy and statistical power, so as to allow discovery of comorbidities
329 of rare diseases, using small datasets, and *ab initio* discovery using very large EHR
330 collections. Our results document the efficacy of PBC for achieving these ends. Still, it is
331 important to bear in mind, that the comorbidities it discovers do not necessarily indicate
332 mechanistic relationships. For instance, two diagnoses may both be driven by smoking,
333 but since smoking was not included logistic regression model, we cannot say anything
334 about smoking as a potential driver (cause) of the relationship. Thus PBC, like all existing
335 methods in this domain, cannot with certainty assign comorbidities to one of the four
336 etiological models described by Valderas et al¹; it can only say that the relationship is not
337 due to factors included in the logistic regression model.

338 Our hope is that PBC will provide an effective solution for a foundational step for outcomes
339 research. The curse of dimensionality is a well-known phenomenon in which training a
340 predictor with too many features can lead to higher error rates^{45,46}. Considering there are
341 about 69,000 ICD10 diagnosis codes, 70,000 ICD10 procedure codes, and 350,000
342 RxNorm CUI codes, dimensionality reduction is necessary for effective machine learning
343 on EHRs. By discovering which variables influence which outcome, PBC can reduce
344 dimensionality and facilitate the creation of downstream tools for outcome predictions.
345 Thus, PBC's role in feature selection becomes clear. For a given clinical outcome, PBC
346 can produce a manageable set of pairwise associations which become the inputs for
347 training predictive models of disease.

348 It is important to note the limitations inherent in the use of data from a single EHR for
349 comorbidity discovery. Data from a single EHR represents a non-random sampling of the

350 general population and PBC does not model this sampling bias. Billing practices can vary
351 within hospital systems - for instance, between inpatient and outpatient services and
352 between provider billing and hospital billing. Hospital billing is performed by medical billing
353 specialists, whereas provider billing is performed by clinicians. In this paper, we have
354 restricted our analysis to provider billing terms. Clinical notes provide a still richer, more
355 nuanced source of data and may more accurately describe a patient's medical condition.
356 Clearly application of PBC to the outputs of Natural Language Processing (NLP) tools will
357 be a fruitful path for future research.

358 Conclusion

359 Capobianco and Lio⁴⁷ present a vision for comorbidity discovery and analysis that is multi-
360 disciplinary and enabled by dynamic networks, with time as a key component in
361 explaining disease relationships. We share this vision, and our PBC method directly
362 addresses the challenges for creating a scalable network-based approach that can (1)
363 dynamically adjust for confounding demographic variables and (2) model temporal
364 relationships in large, complex EHR datasets. However, comorbidities do not exist as
365 isolated pairs, rather they combine in a conditionally dependent manner to create a
366 complex web of influence on any given outcome. While PBC is powered to discover that
367 web by identifying the major drivers of a particular outcome, determining the joint
368 contributions of conditionally dependent variables on that outcome requires a separate
369 computational machinery. Bayesian networks⁴⁸⁻⁴⁹, for example can be used to compute
370 the joint contributions of multiple conditionally dependent variables (so-called multimorbid
371 calculations), providing fully explainable patient outcome predictions.

372 Obtaining a global overview of comorbidity and disease progressions across a major
373 research hospital network is as difficult as it is desirable. We offer the PBC web-browser
374 as a first-generation navigation tool for this new domain of EHR database visualization.
375 Our hope is that the PBC website will provide a community resource for outcomes
376 research, laying the foundation for improving current comorbidity-based outcomes tools,
377 creation of new ones, and, more generally, fueling healthcare discovery for improved
378 care.

379 Methods

380 ***University of Utah medical records***

381 The University of Utah maintains an Electronic Data Warehouse (EDW) – a central
382 storage and search facility for all data collected from all university hospitals and clinics,
383 and all departments and specialties. SQL queries were composed to the following
384 information: (1) medical record number, sex, race, ethnicity, and age for each patient; (2)
385 list of patient visits, along with visit date, and medical terms associated with each visit,
386 including diagnostic codes, procedure codes, and medications ordered. Data were
387 deidentified.

388 We collect ICD9¹⁷ and ICD10¹⁸ diagnosis codes CPT procedural codes²⁰ and RXNorm²¹
389 medication codes (“concept unique identifiers”) from University of Utah electronic medical
390 records. ICD and CPT diagnosis and procedure codes were mapped to the hierarchical
391 Clinical Classification System (CCS)¹⁹. CCS codes allow for more powerful statistics at
392 the expense of concept resolution. After mapping to CCS, we retain 1007 distinct

393 diagnosis codes and 259 distinct procedure codes. These codes include both internal
394 nodes and leaf nodes in the hierarchical CCS tree. In all, we collected records for 1.6
395 million patients, 50 million visits and 150 million diagnosis (DX), procedure (PX) and
396 medication (RX) codes.

397 Counts of EDW patient demographics are displayed in **Table S1**. **Figure S1** displays how
398 our data is distributed by gender and age decade. **Figure S1**, panel B shows how the
399 length of patient medical records (in years) is distributed. Note that these lengths are
400 limited by the history of electronic data collection at the University of Utah that began in
401 the early 2000s but started to ramp up around 2009 and has since increased rapidly.
402 Thus, we see the 95th percentile for medical history length is around 12-15 years for most
403 age bins. **Figure S1** panels C and D show how the number of visits and the number of
404 terms in a medical record trend with age. In almost all decades, women have more
405 medical visits and medical terms than men, though this effect is most pronounced
406 between 20 and 50.

407 ***Logistic regression for person-term probabilities***

408 The initial step in comorbidity analysis is to ascertain the probability of a given term being
409 found in a given person's medical record. A naïve approach could assign everyone the
410 same probability based on the term's frequency in the database or within each age-
411 gender strata as seen in other methods.

412 Our approach involves developing a logistic regression model for each term. The
413 independent features, X , are the list of persons in the EHR along with their gender, race,
414 ethnicity, financial class and risk exposure. Risk exposure includes the age of a person

415 at the time of their last visit as well as the length and density of their medical record.
416 Length is defined as the number of days between first visit and last visit, while density is
417 approximated by the number of visits within a medical record. The dependent outcome,
418 y , is a binary vector indicating whether each person has the term in their medical record.
419 The value C is the inverse of regularization strength in L2 penalized logistic regression.
420 Smaller values of C indicate stronger regularization. The coefficients are determined by
421 minimizing the following loss function, where β represents the coefficients and c is a
422 constant (see Scikit-learn documentation⁵⁰ for a more complete discussion of logistic
423 regression):

$$424 \quad \min_{\beta, c} \frac{1}{2} \beta^T \beta + C \sum_{i=1}^n \log \left(\exp \left(-y_i (X_i^T \beta + c) \right) + 1 \right) \quad (1)$$

425 For each term, stratified 3-fold cross-validation is used to determine the optimal value for
426 C within the set $\{10^{-14}, 10^{-13}, 10^{-12}, \dots, 10^{12}, 10^{13}, 10^{14}\}$. Cross-validation relies on a scoring
427 function to assess the accuracy of logistic regression, given differing values for C . We
428 evaluated standard and custom score functions based on their ability to differentiate
429 logistic regression results with differing C values (see Figure S2).

430 The above approach resulted in logistic regression models for each term, capable of
431 predicting the probability that a given person has a given term. For rare terms, we find
432 that the probabilities output by logistic regression may not sum up to the actual number
433 of patients with the term. To account for this we adjust each probability by a bias
434 correction factor such that the sum of the adjusted patient probabilities is equal to the

435 actual number of patients with the term. The exact form of this correction factor is given
436 in the supplemental “Math.pdf”.

437 A limitation of logistic regression is the lack of a confidence metric on each predicted
438 probability output by the regression model. For instance, predicted probabilities might be
439 more accurate for patients of western European than African ancestry, because the
440 corpus data is skewed for this demographic variable (see Table S1). To overcome this
441 limitation, we divide our data into 6 randomized partitions, balanced so that the number
442 of affected individuals in each partition is approximately equal. This partitioning is
443 accomplished using “StratifiedKFold” from the python package sklearn⁴⁶. Next each
444 partition is used to fit a logistic regression model, using the previously determined
445 regularization strength. Each model is used to predict the probability of each person
446 having the term. These 6 probabilities are used to calculate a sample variance for each
447 person-term probability, $s_{P_{t \in m}}^2$, where m represents a patient’s medical record, t
448 represents a medical term, and $P_{t \in m}$ is the probability of term t in m . To be clear these 6
449 LRMs are only used to calculate sample variance, while the LRM trained on the full
450 dataset is used per-patient per-term probability predictions.

451 ***Poisson Binomial for term pair p-values***

452 Our null hypothesis is that pairs of medical terms are independently distributed in
453 University of Utah medical records, $H_0: t1 \perp t2 | \mathcal{M}$. A significant p-value would indicate
454 that a pair of terms co-occur more often than would occur by chance. Given two
455 independent terms, $t1 \perp t2$ the probability the two would occur by chance in a given
456 person’s medical record m , is the product of their individual probabilities:

457
$$P_{\{t_1, t_2\} \subseteq m} = P_{t_1 \in m} * P_{t_2 \in m} \quad (2)$$

458 A naïve approach assumes person-term probabilities are equal to population incidence
459 rates:

460
$$P_{t_1 \in m} = \frac{|\mathcal{M}_{t_1}|}{|\mathcal{M}|}; P_{t_2 \in m} = \frac{|\mathcal{M}_{t_2}|}{|\mathcal{M}|} \quad (3)$$

461 Using the naïve approach, person-term-pair probabilities follow the binomial distribution
462 with

463
$$P_{\{t_1, t_2\} \subseteq m} = \frac{|\mathcal{M}_{t_1}|}{|\mathcal{M}|} * \frac{|\mathcal{M}_{t_2}|}{|\mathcal{M}|} \quad (4)$$

464 However, using logistic regression, we have different probabilities for each person/term
465 pair. The Poisson binomial distribution is the discrete distribution of a sum of Bernoulli
466 trials where the probability of each trial differs. Thus, using logistic regression, our data
467 follows a Poisson binomial distribution.

468 Because the cumulative distribution function (CDF) of a Poisson binomial is
469 computationally tractable only for a small number of values, numerous approximations
470 have been developed⁵². We use the normal approximation because it is fast and accurate
471 for large datasets. To determine a p-value using the normal approximation, the mean and
472 variance for the Poisson binomial are calculated and used as parameters for a normal
473 distribution. The mean of a Poisson binomial represents the expected number of
474 University of Utah patients who will have in their medical record both terms in the pair and
475 is calculated as the sum of probabilities for each person m :

476
$$\mu_{t_1, t_2} = \sum_m^{\mathcal{M}} P_{t_1, t_2 | m} \quad (5)$$

477

478 The variance of a Poisson binomial is likewise similar in form to a binomial distribution:

479
$$\sigma_{t_1, t_2}^2 = \sum_m^{\mathcal{M}} P_{\{t_1, t_2\} \subseteq m} (1 - P_{\{t_1, t_2\} \subseteq m}) \quad (6)$$

480 The Poisson binomial variance is augmented with the logistic regression variances
481 described in the previous section ($s_{P_{t_1 \subseteq m}}^2, s_{P_{t_2 \subseteq m}}^2$) using the product rule and the law of total
482 variance. One can think of these variances as *measurement error* for $P_{t_1 \subseteq m}$ and $P_{t_2 \subseteq m}$
483 and they are larger for rare terms.

484 The probability that a person has a pair of terms, $P_{\{t_1, t_2\} \subseteq m}$, can be so rare it exceeds the
485 limits of floating-point arithmetic. Thus, we implement our methods in log-space. Our
486 logistic regression models report per-patient per-term log probabilities. We calculate
487 $\ln(\mu_{t_1, t_2})$ and $\ln(\sigma_{t_1, t_2}^2)$ rather than summing in normal space. We use a numerical
488 approximation to calculate the $\ln(pvalue)$ of a normal distribution⁵³ (implemented as
489 “gsl_sf_log_erfc” in the Gnu scientific library⁵⁴). To report significance, we use throughout
490 this paper an alpha threshold of 0.05. Since we calculate p-values for 4,622,320 pairs of
491 medical terms, our Bonferroni corrected alpha is set to 1.08e-08.

492 **Direction P-values**

493 Given two terms that occur together in medical records more often than would occur by
494 chance, which term tends to occur first in the medical record, or do they tend to occur in
495 the same time frame? We calculate p-values for the temporal nature of each association.
496 For each patient with medical record m , the date of the first occurrence of each term in m
497 is recorded. Pairs of terms occurring within a window of size W are labeled as “in-window”.
498 Pairs of terms occurring outside of W contribute to the $t_1 \rightarrow t_2$ count or the $t_2 \rightarrow t_1$
499 count. For the analyses presented here, we chose a 90-day window, as this duration
500 decreases noise associated with the date of information capture within the medical
501 record, but the approach is valid over any interval. Note that some term relationships –
502 such as a surgery procedure followed by an infection diagnosis – will only show a
503 significant direction with a window smaller than 90 days. PBC-Web includes 4 window
504 sizes (30, 90, 365, 730).

505 For a person with medical record m , containing terms 1 and 2 ($\{t_1, t_2\} \subseteq m$), the
506 probability that term 1 occurs before term 2 is a function of the ratio of the probabilities of
507 the 2 terms:

$$508 \quad P_{t_1 \rightarrow t_2}^{m_{t_1, t_2}} = P_{t_1 \rightarrow t_2 | \{t_1, t_2\} \subseteq m}^m = \frac{P_{t_1}^m}{P_{t_1}^m + P_{t_2}^m} \quad (7)$$

509 Given $\{t_1, t_2\} \subseteq m$, the probability t_1 and t_2 occur within a time window of size W , is a
510 function of the span or length of a person’s medical history, $Span_m$:

511
$$P_{t_1 \sim t_2}^{m_{t_1, t_2}} = P_{t_1 \sim t_2 | \{t_1, t_2\} \subseteq m}^m = \begin{cases} 1, & \text{Span}_m \leq W \\ \frac{W}{\text{Span}_m} \left(2 - \frac{W}{\text{Span}_m} \right), & \text{otherwise} \end{cases} \quad (8)$$

512 The above formula represents the percent of timepoints t_1 and t_2 that fall within W days
513 of each other. The derivation of the formula is given in the supplemental “Math.pdf”.

514 The product of (7) and (8) gives the probability term 1 would precede term 2 within a
515 window of size W :

516
$$P_{t_1 \rightarrow t_2}^{m_{t_1, t_2}} = P_{t_1 \rightarrow t_2 | \{t_1, t_2\} \subseteq m}^{in W} = P_{t_1 \rightarrow t_2}^{m_{t_1, t_2}} \cdot P_{t_1 \sim t_2}^{m_{t_1, t_2}} \quad (9)$$

517 Similar logic can be applied to derive the “out-of-window” probability of term 1 occurring
518 at least W days before term 2. A more complete explanation is found in supplemental
519 “Math.pdf”. Direction p-values are calculated using a normal approximation of the Poisson
520 Binomial CDF as in the previous section.

521 **Code, web-development, and calculations**

522 We implement our statistical analysis using Python, Cython and C. Cython is a static
523 compiler for Python and the extended Cython programming language⁵⁵. Scikit-learn⁵¹ was
524 used for Logistic regression studies. All the figures accompanying this article were
525 generated using Matplotlib⁵⁶. Our website is built using the Flask web framework⁵⁷. The
526 backend is pure python and the front end is JavaScript and D3⁵⁸. Logistic regression
527 modeling and pairwise calculation of Poisson Binomial p-values were performed at the
528 University of Utah Center for High Performance Computing (CHPC) PHI protected
529 environment. Training 3041 logistic regression models while tuning the regularization

530 strength with cross validation took 1959 CPU hours. Maximum memory usage was 174
531 MB. Calculating comorbidity statistics for 4,623,841 term pairs took 7455 CPU hours while
532 maximum memory usage was 87 MB.

533 Acknowledgements

534 The following collaborators have provided valuable discussion, feedback, and insight
535 which has guided development of PBC: Bruce Bray, Vikrant Deshmukh, Karen Eilbeck,
536 Edgar Javier Hernandez, Rashmee Shah. We thank members of the University of Utah
537 EDW for facilitating access to medical records. The computational resources used were
538 partially funded by the NIH Shared Instrumentation Grant 1S10OD021644-01A1.

539 This research was supported by the AHA Children's Strategically Focused Research
540 Network grant (17SFRN33630041) and the Nora Eccles Treadwell Foundation. Gordon
541 Lemmon was supported by NRSA training grant T32H757632. Sergiusz Wesolowski was
542 supported by NRSA training grant T32DK110966-04 and the AHA Children's Strategically
543 Focused Research Network Fellowship award (17SFRN33630041).

544 Author Contributions

545 Gordon Lemmon is the senior research associate leading PBC development and
546 validation. Sergiusz Wesolowski is an applied mathematician who has helped formalize
547 our approach to statistical testing. Alex Henrie was a software engineer on the project.
548 Martin Tristani-Firouzi and Mark Yandell conceived of the project and secured research
549 funding and played a key role in scientific discussions regarding development of PBC. All
550 authors edited the manuscript.

551 Competing interests

552 GL, MY own shares in Backdrop Health, a University of Utah effort to commercialize
553 Bayesian inference on health records. However, there are no financial ties regarding
554 this research.

555 Data Availability

556 In this paper we calculate comorbidity statistics for all pairs of medical billing codes –
557 including diagnoses, procedures, and medications. All of these p-values are available to
558 query and download from the following link: <https://pbc.genetics.utah.edu/lemmon2021>.

559 Code Availability

560 We provide a CodeOcean capsule with code and data; the link is submitted by the
561 editor to the reviewers during the peer-review process.

562 References

- 563 1. Valderas, J. M., Starfield, B., Sibbald, B., Salisbury, C. & Roland, M. Defining
564 comorbidity: implications for understanding health and health services. *Ann. Fam.*
565 *Med.* **7**, 357–363 (2009).
- 566 2. Lone, N. I. *et al.* Predicting risk of unplanned hospital readmission in survivors of
567 critical illness: a population-level cohort study. *Thorax* **74**, 1046–1054 (2019).
- 568 3. Wang, H. *et al.* Predicting Hospital Readmission via Cost-Sensitive Deep Learning.
569 *IEEE/ACM Trans. Comput. Biol. Bioinform.* **15**, 1968–1978 (2018).

- 570 4. Facchinetti, G. *et al.* Continuity of care interventions for preventing hospital
571 readmission of older people with chronic diseases: A meta-analysis. *Int. J. Nurs.*
572 *Stud.* **101**, 103396 (2020).
- 573 5. Atashi, A., Sarbaz, M., Marashi, S., Hajlaliasgari, F. & Eslami, S. Intensive Care
574 Decision Making: Using Prognostic Models for Resource Allocation. *Stud. Health*
575 *Technol. Inform.* **251**, 145–148 (2018).
- 576 6. Yurkovich, M., Avina-Zubieta, J. A., Thomas, J., Gorenchtein, M. & Lacaille, D. A
577 systematic review identifies valid comorbidity indices derived from administrative
578 health data. *J. Clin. Epidemiol.* **68**, 3–14 (2015).
- 579 7. Charlson, M. E., Pompei, P., Ales, K. L. & MacKenzie, C. R. A new method of
580 classifying prognostic comorbidity in longitudinal studies: development and
581 validation. *J. Chronic Dis.* **40**, 373–383 (1987).
- 582 8. Elixhauser, A., Steiner, C., Harris, D. R. & Coffey, R. M. Comorbidity Measures for
583 Use with Administrative Data: *Med. Care* **36**, 8–27 (1998).
- 584 9. Roque, F. S. *et al.* Using Electronic Patient Records to Discover Disease
585 Correlations and Stratify Patient Cohorts. *PLOS Comput. Biol.* **7**, e1002141 (2011).
- 586 10. Gutiérrez-Sacristán, A. *et al.* comoRbidity: an R package for the systematic
587 analysis of disease comorbidities. *Bioinformatics* **34**, 3228–3230 (2018).
- 588 11. Moni, M. A., Xu, H. & Liò, P. CytoCom: a Cytoscape app to visualize, query and
589 analyse disease comorbidity networks. *Bioinforma. Oxf. Engl.* **31**, 969–971 (2015).
- 590 12. Moni, M. A. & Liò, P. comoR: a software for disease comorbidity risk assessment.
591 *J. Clin. Bioinforma.* **4**, 8 (2014).
- 592 13. Ronzano, F., Gutiérrez-Sacristán, A. & Furlong, L. I. Comorbidity4j: a tool for

- 593 interactive analysis of disease comorbidities over large patient datasets.
594 *Bioinforma. Oxf. Engl.* **35**, 3530–3532 (2019).
- 595 14. Siggaard, T. *et al.* Disease trajectory browser for exploring temporal, population-
596 wide disease progression patterns in 7.2 million Danish patients. *Nat. Commun.*
597 **11**, 4952 (2020).
- 598 15. Winter, A. C., Rist, P. M., Buring, J. E. & Kurth, T. Prospective comorbidity-
599 matched study of Parkinson’s disease and risk of mortality among women. *BMJ*
600 *Open* **6**, (2016).
- 601 16. Johnson, Alistair *et al.* MIMIC-IV (version 1.0). *Physionet*. doi:10.13026/S6N6-
602 XD98 (2021).
- 603 17. ICD-9-CM - International Classification of Diseases, Ninth Revision, Clinical
604 Modification. <https://www.cdc.gov/nchs/icd/icd9cm.htm> (2019).
- 605 18. ICD-10-CM - International Classification of Diseases, Tenth Revision, Clinical
606 Modification. <https://www.cdc.gov/nchs/icd/icd10cm.htm> (2020).
- 607 19. Clinical Classifications Software Refined (CCSR). [https://www.hcup-](https://www.hcup-us.ahrq.gov/toolssoftware/ccsr/ccs_refined.jsp)
608 [us.ahrq.gov/toolssoftware/ccsr/ccs_refined.jsp](https://www.hcup-us.ahrq.gov/toolssoftware/ccsr/ccs_refined.jsp).
- 609 20. CPT Codes - Current Procedural Terminology - AAPC.
610 <https://www.aapc.com/resources/medical-coding/cpt.aspx>.
- 611 21. Liu, S., Wei Ma, Moore, R., Ganesan, V. & Nelson, S. RxNorm: prescription for
612 electronic drug information exchange. *IT Prof.* **7**, 17–23 (2005).
- 613 22. Youden, W. J. Index for rating diagnostic tests. *Cancer* **3**, 32–35 (1950).
- 614 23. Han, H. *et al.* Hypertension and breast cancer risk: a systematic review and meta-
615 analysis. *Sci. Rep.* **7**, (2017).

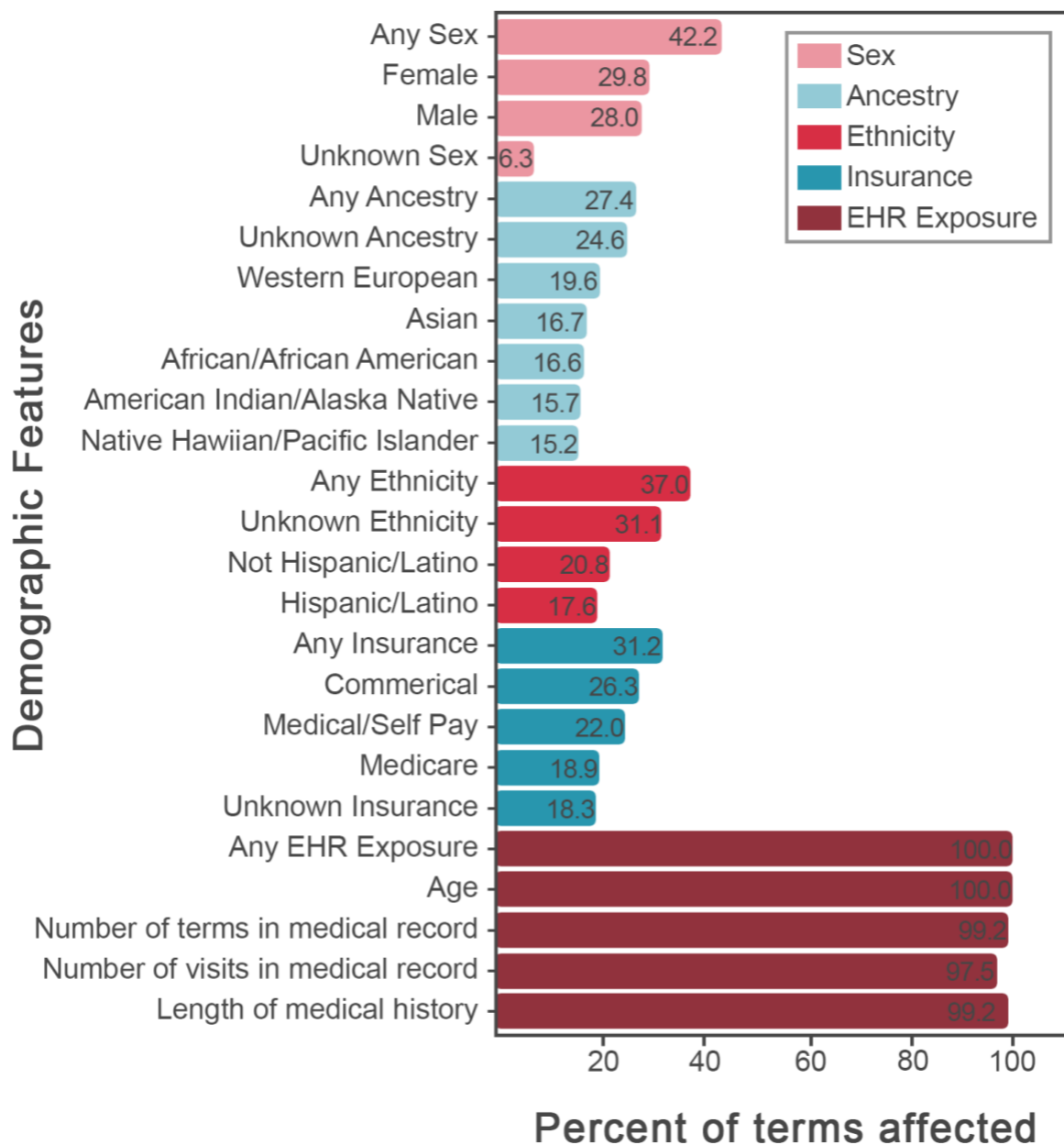
- 616 24. Li, X. *et al.* Comorbidities among patients with cancer who do and do not develop
617 febrile neutropenia during the first chemotherapy cycle. *J. Oncol. Pharm. Pract.*
618 *Off. Publ. Int. Soc. Oncol. Pharm. Pract.* **22**, 679–689 (2016).
- 619 25. Chia, V. M. *et al.* Chronic comorbid conditions associated with risk of febrile
620 neutropenia in breast cancer patients treated with chemotherapy. *Breast Cancer*
621 *Res. Treat.* **138**, 621–631 (2013).
- 622 26. Toma-Dasu, I., Wojcik, A. & Kjellsson Lindblom, E. Risk of second cancer following
623 radiotherapy. *Phys. Medica PM Int. J. Devoted Appl. Phys. Med. Biol. Off. J. Ital.*
624 *Assoc. Biomed. Phys. AIFB* **42**, 211–212 (2017).
- 625 27. Donin, N. *et al.* Risk of second primary malignancies among cancer survivors in
626 the United States, 1992 through 2008. *Cancer* **122**, 3075–3086 (2016).
- 627 28. Grantzau, T. & Overgaard, J. Risk of second non-breast cancer among patients
628 treated with and without postoperative radiotherapy for primary breast cancer: A
629 systematic review and meta-analysis of population-based studies including
630 522,739 patients. *Radiother. Oncol. J. Eur. Soc. Ther. Radiol. Oncol.* **121**, 402–413
631 (2016).
- 632 29. Rissanen, J. Modeling by shortest data description. *Automatica* **14**, 465–471
633 (1978).
- 634 30. Hassell, K. L. Population estimates of sickle cell disease in the U.S. *Am. J. Prev.*
635 *Med.* **38**, S512-521 (2010).
- 636 31. Ahmadi, M., Poormansouri, S., Beiranvand, S. & Sedighie, L. Predictors and
637 Correlates of Fatigue in Sickle Cell Disease Patients. *Int. J. Hematol.-Oncol. Stem*
638 *Cell Res.* **12**, 69–76 (2018).

- 639 32. Herson, J., Sharma, S., Crocker, C. L. & Jones, D. Physical complaints of patients
640 with sickle cell trait. *J. Reprod. Med.* **14**, 129–132 (1975).
- 641 33. Aich, A., Jones, M. K. & Gupta, K. Pain and sickle cell disease. *Curr. Opin.*
642 *Hematol.* **26**, 131–138 (2019).
- 643 34. Tariq, S. & Aronow, W. S. Use of Inotropic Agents in Treatment of Systolic Heart
644 Failure. *Int. J. Mol. Sci.* **16**, 29060–29068 (2015).
- 645 35. Anders, H.-J., Huber, T. B., Isermann, B. & Schiffer, M. CKD in diabetes: diabetic
646 kidney disease versus nondiabetic kidney disease. *Nat. Rev. Nephrol.* **14**, 361–377
647 (2018).
- 648 36. Koye, D. N., Magliano, D. J., Nelson, R. G. & Pavkov, M. E. The Global
649 Epidemiology of Diabetes and Kidney Disease. *Adv. Chronic Kidney Dis.* **25**, 121–
650 132 (2018).
- 651 37. El Fane, M. *et al.* [Pneumocystosis during HIV infection]. *Rev. Pneumol. Clin.* **72**,
652 248–254 (2016).
- 653 38. Seravalle, G. & Grassi, G. Obesity and hypertension. *Pharmacol. Res.* **122**, 1–7
654 (2017).
- 655 39. Choi, E., Xiao, C., Stewart, W. F. & Sun, J. MiME: Multilevel Medical Embedding of
656 Electronic Health Records for Predictive Healthcare. *ArXiv181009593 Cs Stat*
657 (2018).
- 658 40. Rajkomar, A. *et al.* Scalable and accurate deep learning with electronic health
659 records. *Npj Digit. Med.* **1**, 1–10 (2018).
- 660 41. Franz, L., Shrestha, Y. R. & Paudel, B. A Deep Learning Pipeline for Patient
661 Diagnosis Prediction Using Electronic Health Records. *ArXiv200616926 Cs* (2020).

- 662 42. Miotto, R., Li, L., Kidd, B. A. & Dudley, J. T. Deep Patient: An Unsupervised
663 Representation to Predict the Future of Patients from the Electronic Health
664 Records. *Sci. Rep.* **6**, 26094 (2016).
- 665 43. Hassaine, A., Salimi-Khorshidi, G., Canoy, D. & Rahimi, K. Untangling the
666 complexity of multimorbidity with machine learning. *Mech. Ageing Dev.* **190**,
667 111325 (2020).
- 668 44. Hassaine, A. *et al.* Learning multimorbidity patterns from electronic health records
669 using Non-negative Matrix Factorisation. *J. Biomed. Inform.* **112**, 103606 (2020).
- 670 45. Chandrasekaran, B. & Jain, A. K. Quantization Complexity and Independent
671 Measurements. *IEEE Trans. Comput.* **C-23**, 102–106 (1974).
- 672 46. Trunk, G. V. A Problem of Dimensionality: A Simple Example. *IEEE Trans. Pattern*
673 *Anal. Mach. Intell.* **PAMI-1**, 306–307 (1979).
- 674 47. Capobianco, E. & Lio', P. Comorbidity: a multidimensional approach. *Trends Mol.*
675 *Med.* **19**, 515–521 (2013).
- 676 48. Pearl, J. Reverend bayes on inference engines: a distributed hierarchical
677 approach. in *Proceedings of the Second AAAI Conference on Artificial Intelligence*
678 133–136 (AAAI Press, 1982).
- 679 49. Pearl, J. *Causality: models, reasoning, and inference*. (Cambridge University
680 Press, 2013).
- 681 50. 1.1. Linear Models — scikit-learn 0.24.1 documentation. [https://scikit-](https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression)
682 [learn.org/stable/modules/linear_model.html#logistic-regression](https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression).
- 683 51. Pedregosa, F. *et al.* Scikit-learn: Machine Learning in Python. *J. Mach. Learn. Res.*
684 **12**, 2825–2830 (2011).

- 685 52. Hong, Y. On computing the distribution function for the Poisson binomial
686 distribution. *Comput. Stat. Data Anal.* **59**, 41–51 (2013).
- 687 53. Hart, J. F. *Computer approximations*. (Wiley, 1968).
- 688 54. *GNU scientific library: reference manual*. (Network Theory, 2009).
- 689 55. Behnel, S. *et al.* Cython: The Best of Both Worlds. *Comput. Sci. Eng.* **13**, 31–39
690 (2011).
- 691 56. Hunter, J. D. Matplotlib: A 2D Graphics Environment. *Comput. Sci. Eng.* **9**, 90–95
692 (2007).
- 693 57. Grinberg, M. *Flask web development: developing web applications with Python*.
694 (O'Reilly, 2018).
- 695 58. Bostock, M., Ogievetsky, V. & Heer, J. D³: Data-Driven Documents. *IEEE Trans.*
696 *Vis. Comput. Graph.* **17**, 2301–2309 (2011).

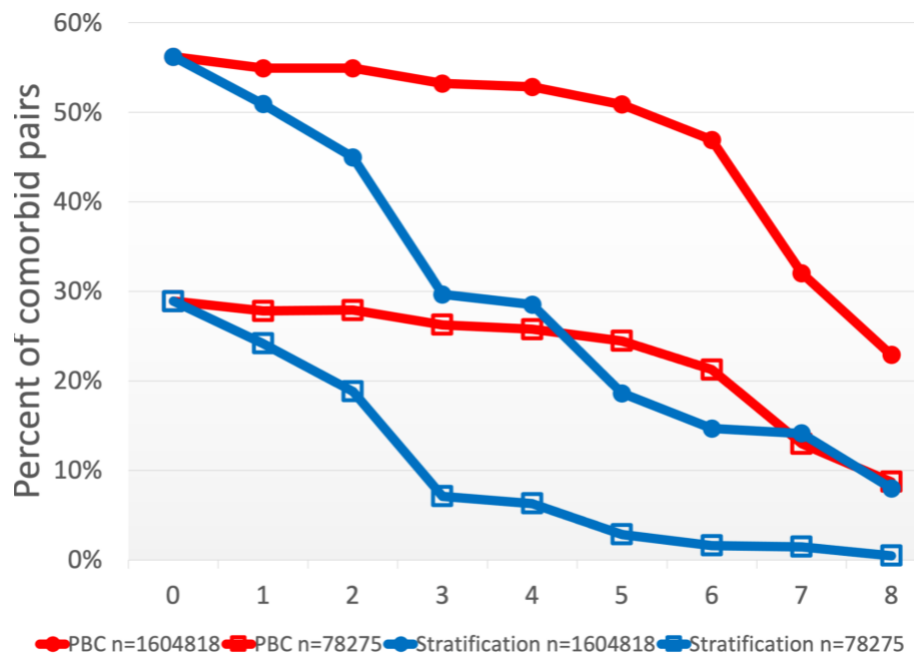
697 Display Items



698

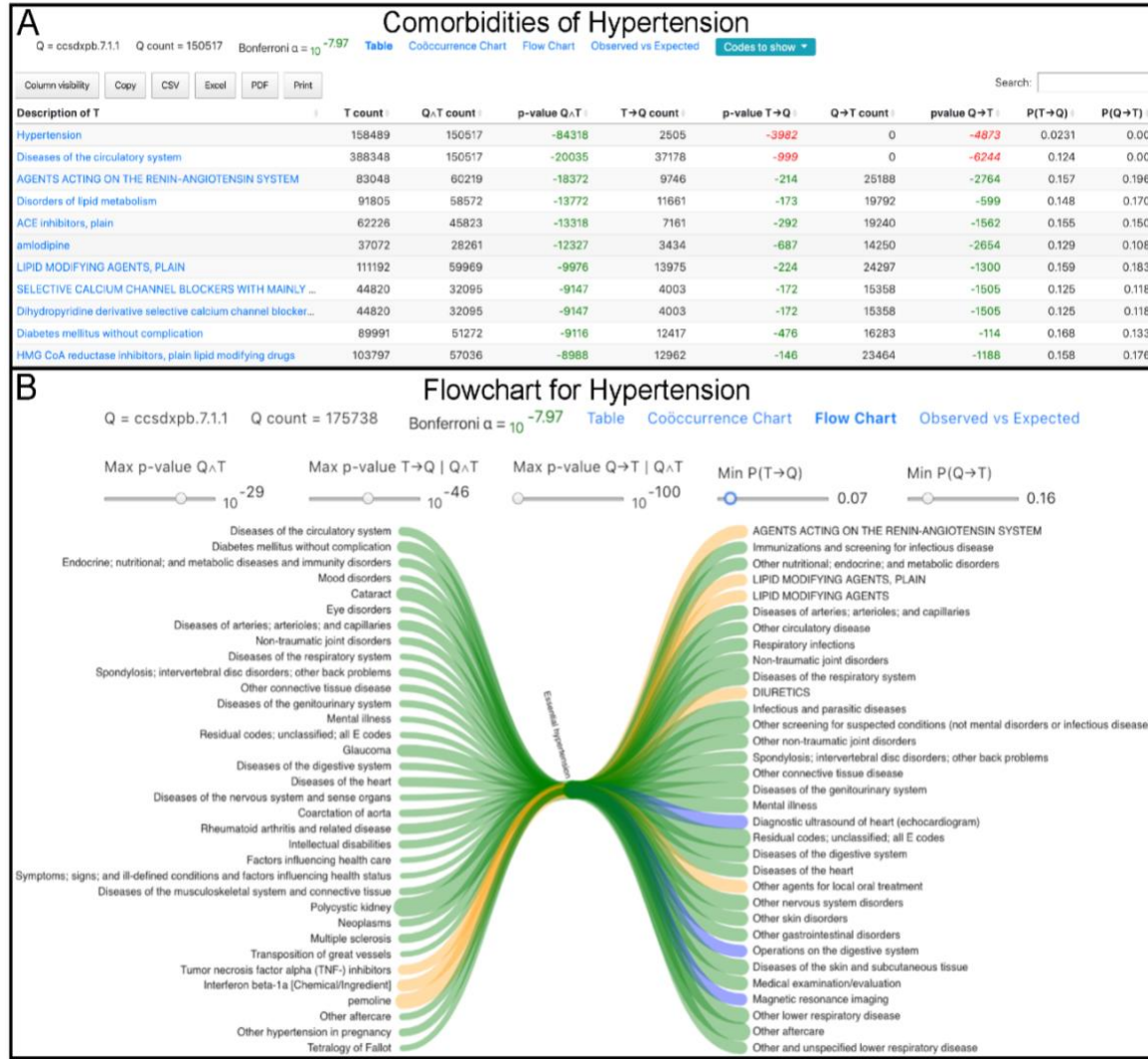
699 **Figure 1. Feature selection by L1 regularization.** Percent of medical term logistic
700 regression models that include each demographic feature. For example, EHR exposure,
701 i.e., the length and density of a person’s medical history, is an important predictor for
702 every medical diagnosis, procedure and medication.

703 **Figure 2. PBC maintains power for discovery by modeling the effects of**
704 **confounding variables.** We calculate the percent of significantly co-occurring pairs of
705 medical terms ($p < 1.08e-8$) using either PBC (blue lines) or stratification (red lines) for
706 two different sample sizes (Entire EHR corpus [1.6 million individuals, filled circles], and
707 a 78,275 patient sample [open squares]). Moving from left to right we introduce
708 additional features to the PBC approach and additional filters to the stratification
709 approach. The X-axis numbering corresponds to the following features/filters: 0: no
710 features; 1: race / African American; 2: sex / Female; 3: Age / 50-59; 4: Ethnicity /
711 nonhispanic; 5: Insurance / Commercial; 6: Span / at least 2 years; 7: Number of visits /
712 at least 3 visits; 8: Date of last visit / at least as recent as Jan 2018. The figure
713 highlights 2 important trends. First, the number of significant associations decreases as
714 a function of smaller datasets. Second, controlling for specific features markedly
715 reduces the number of recovered comorbidities for the stratification approach (an effect



716 further exacerbated by reducing the initial cohort size), while preserving significant

717 comorbidities using PBC.



719 **Figure 3. Screenshots from PBC web.** “Q” refers to the query term, in this case
720 hypertension. “T” refers to the term possibly comorbid with Q. **Panel A**, Code prefixes in
721 the first column can be deciphered as follows: ccs = “clinical classification system”, dx =
722 diagnosis, px = procedure, pb = provider billing (we omit hospital billing codes in this
723 figure), cui = RxNorm concept unique identifier. P-values that pass the Bonferroni
724 corrected significance threshold are colored green or red. Green indicates the
725 relationship occurs more often than expected. Red indicates less often than expected.
726 The last two columns represent flow rates which indicate the actual percent of patients
727 in our database that transit from one term to the next over time. **Panel B**, Terms that
728 significantly precede or follow hypertension (separated by at least 90 days) are shown
729 to the left and right of hypertension respectively. Green connections are diagnoses, blue
730 connections are procedures and orange connections are medications. The thickness of
731 the connection relates to the flow rate – the percent of patients that flow through the
732 given path.

733 **Table 1. PBC retains power as features are added; stratification loses power.**

734 Shown are 3 established comorbidities from the medical literature: concussion and
735 migraine (41), multiple myeloma and multiple sclerosis (42) and cancer of pancreas and
736 hypertension (43). Comorbidities passing a Bonferroni corrected alpha threshold of
737 **1.08e-8** are colored blue. As features are added to the stratification criteria, sample size
738 shrinks and statistical significance is lost. Rather than controlling for confounding
739 variables, PBC models their effects. Thus, for PBC sample size remains constant at
740 1,538,059 and statistical significance is preserved.

Stratification filters	PBC features	Concussion and migraine		Multiple myeloma and Multiple sclerosis		Cancer of pancreas and hypertension	
		P-values					
		χ^2	PBC	χ^2	PBC	χ^2	PBC
No filters (n=1538059)	none	1e-933	1e-933	1e-121	1e-121	1e-405	1e-405
female (n=794281)	+sex	1e-860	1e-937	1e-70	1e-135	1e-191	1e-400
+age 50-59 (n=69527)	+age	1e-126	1e-1031	1.6e-5	2.5e-65	7e-7	1e-239
+Caucasian (n=45782)	+ancestry	1e-92	1e-912	1.5e-5	2e-87	5e-7	1e-159
+nonhispanic (n=39897)	+ethnicity	1e-91	1e-886	4.8e-7	2e-87	6e-11	1e-61
+commercial (n=21148)	+insurance	1e-4	1e-859	0.32	3e-80	8e-4	5e-61
+3yr history (n=9243)	+span	.087	1e-522	0.41	1.6e-60	0.03	5e-78

742 **Table 2. PBC identifies comorbidities specific to underrepresented minorities**
743 **even when data is limited.** Here we compare the stratification-based approach with
744 PBC as regards ability to identify a known comorbidity: sickle cell anemia (SCA) paired
745 with malaise and fatigue. P-values passing a Benjamini-Hockberg corrected alpha
746 threshold of **1.0e-6** are colored blue. Without modeling the effects of confounding
747 variables, both approaches identify the association. But because ancestry is a key
748 determinant of risk for SCA, we need to control for this confounding variable to
749 determine whether malaise and fatigue is an actual symptom of sickle cell anemia or
750 whether the connection is being driven by ancestry. Filters applied under stratification
751 lead to samples too small to detect this association. In contrast, not only does PBC
752 detect the comorbidity, but the strength of the association increases as confounders
753 such as ancestry, ethnicity and age are included in the model.

Stratum filters	PBC features	Sickle Cell Anemia paired with Malaise and Fatigue		
		Stratum Pair count	χ^2 p-value	PBC p-value
No filters (n=477,070)	no features	19	4e-8	4e-8
Caucasian (n=276,496)	ancestry	5	0.002	2e-8
African American (n=7,035)	ancestry	9	2e-5	2e-8
+Nonhispanic (n=6,039)	+ethnicity	8	2e-4	1e-8
+Female (n=2,789)	+gender	7	1e-4	9e-8
+50-59 (n=302)	+age	3	0.004	3e-19

755 **Table 3. PBC is the only comorbidity search tool that takes a high-resolution**
 756 **approach.** CytoCom and comoR are no longer available; comoRbidity and
 757 Comorbidity4j failed to scale to datasets of our size. In addition to scaling to handle
 758 arbitrarily large EHR datasets, PBC models the effects of demographic information
 759 rather than relying on stratification.

Package	CytoCom	comoR	comoRbidity	Comorbidity4j	DTB*	PBC-web
Release date	2014	2014	2018	2019	2020	2021
Relative Risk [†]	✓	✓	✓	✓	✓	✓
ϕ -correlation [†]	✓	✓	✓	✓		✓
Comorbidity score [†]			✓	✓		✓
Odds ratio [†]			✓	✓		✓
Fisher's Exact test [†]		✓	✓	✓		✓
χ^2 or Binomial test					✓	✓
Poisson Binomial [‡]						✓
Inverse comorbidities						✓
Temporal directionality			✓	✓	✓	✓
Interactive network of comorbidities	✓				✓	✓
Arbitrarily large datasets					See note [§]	✓
Dealing with confounders	stratify	no support	stratify	stratify	matching	Logistic regression
Platform	Cytoscape	R package	R package	Java/website	Publicly available website	Publicly available website

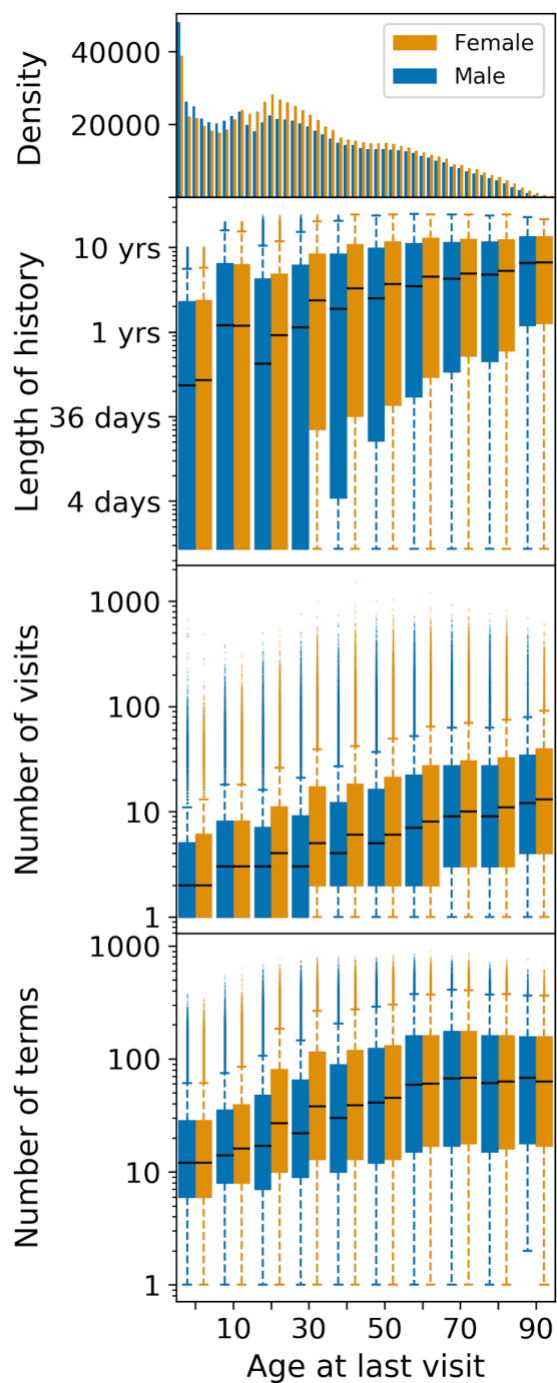
*Disease Trajectory Browser

[†]Statistics that rely on population incidence rate

[‡]Statistics that rely on per-person per-term probabilities

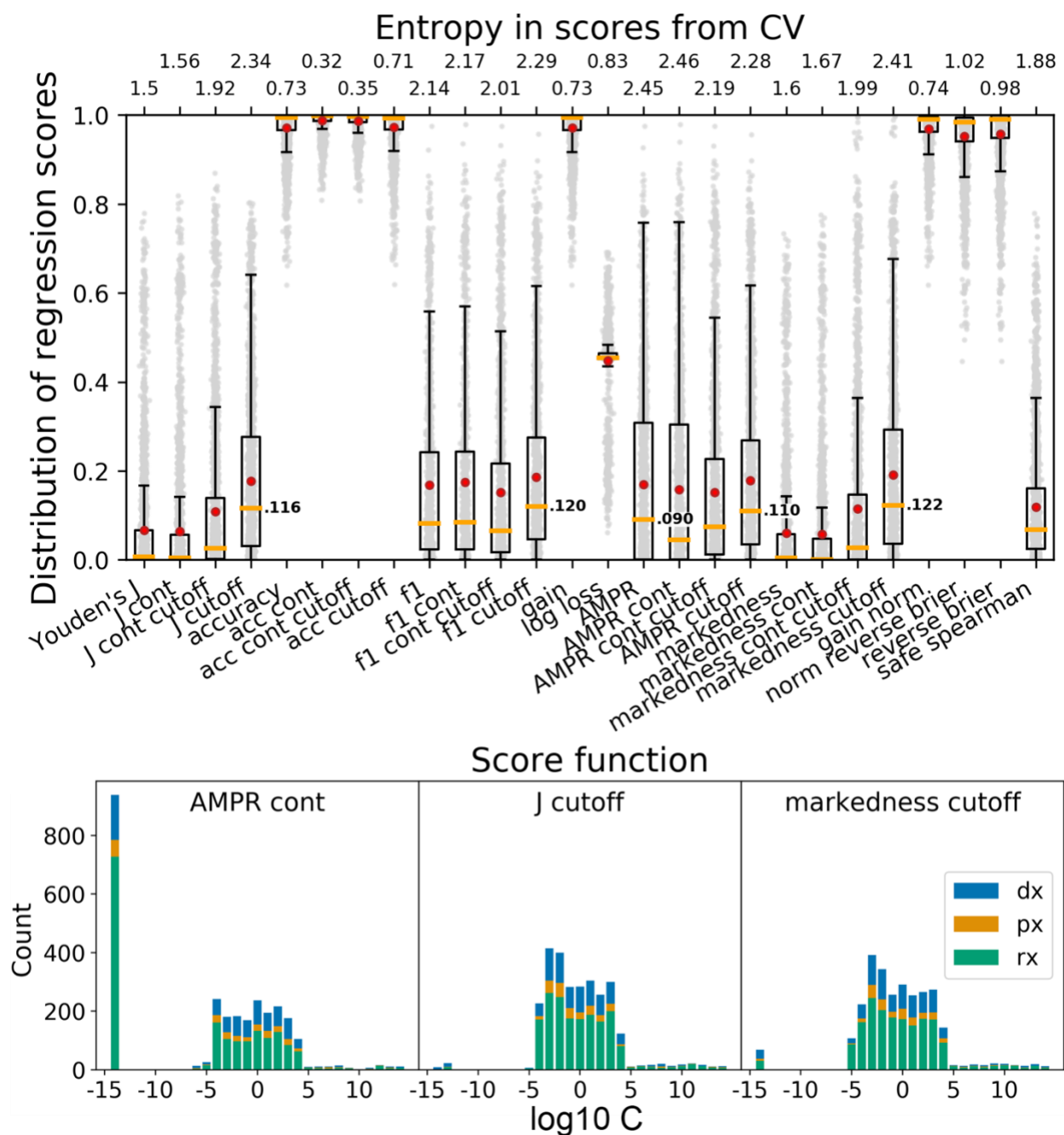
[§]Authors apply a pre-filtering step because calculating all pairs of comorbidities is too computationally demanding

760 Supplemental Display Items

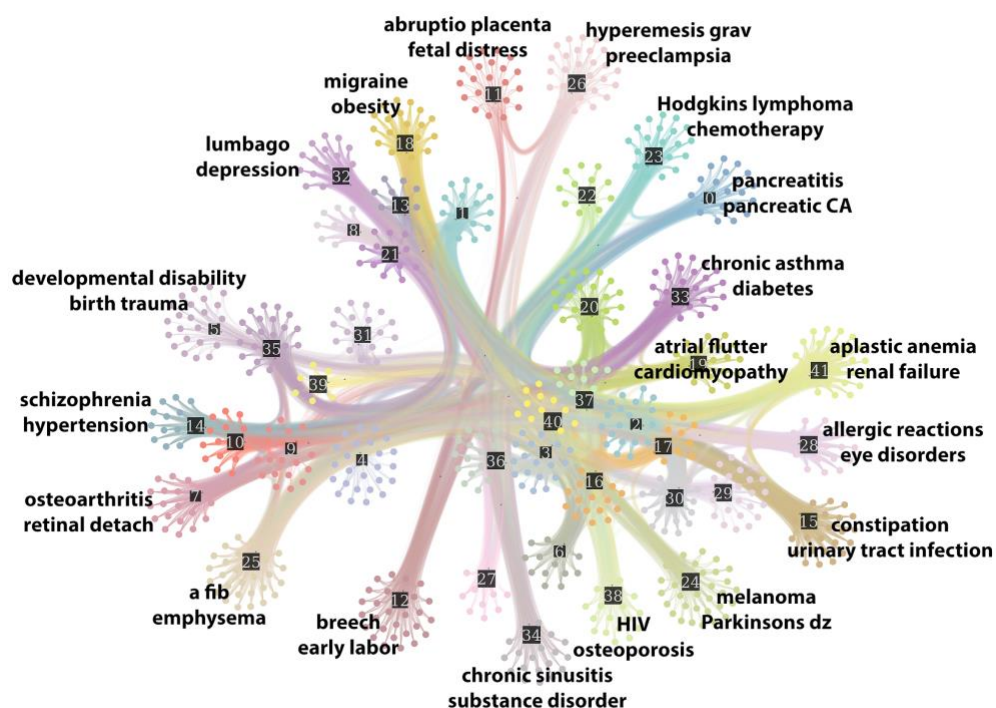


761

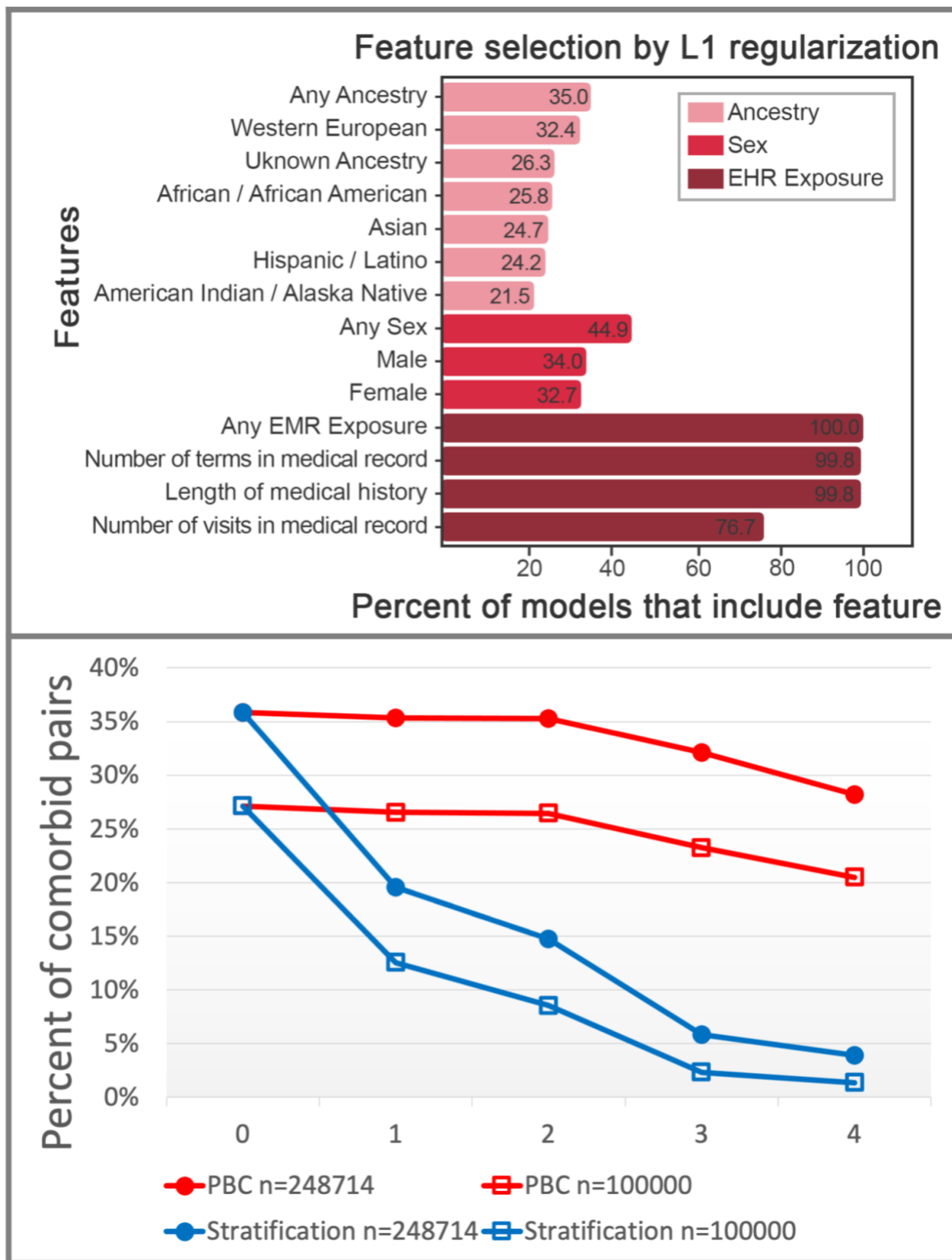
762 **Figure S1. University of Utah medical records binned by age-decade.** Boxplots
 763 show median (black line), 25th and 75th percentile (box ends), 95th and 5th percentile
 764 (whisker caps) and outliers. Number of terms (bottom panel) is a count of distinct
 765 diagnoses, procedures and medications found in each patient's medical history.



767 **Figure S2. J cutoff maximizes entropy and minimizes outliers.** Comparison of score
768 functions for logistic regression C-value optimization. For each score function, we
769 evaluated C-values ranging from 10^{-14} to 10^{14} . Top: For each of 3041 DX, PX, and RX
770 terms, we use cross validation to select the C-value that achieves the best score. Each
771 boxplot contains these 3041 best scores as evaluated with different score functions.
772 Bottom: Distribution of C-values for 3 score functions with high entropy. J_cutoff was
773 chosen for downstream analysis because it has high entropy and has a smooth C-value
774 distribution without the large outlier at C=-14.



775 **Figure S3.** Minimum description length of the comorbidity network discovered by the
776 PBC approach for diagnoses in the University of Utah EDW. Examples of significantly
777 associated medical conditions within each cluster are displayed. Citations supporting
778 these associations are listed in Suppl Table S5.



780 **Figure S4. Deployment of PBC on MIMIC-IV EHR data.** See Figure 1 legend for
781 description of top panel and Figure 2 legend for description of bottom panel. Bottom
782 panel, the X-axis ticks correspond to the addition of regression features (PBC) or
783 stratification criteria from left to right: 0 - no features, no stratification, 1- gender/female,
784 2 - ancestry/African American, 3 - length of medical history/at least 2 years, 4 - number
785 of visits/at least 3 visits. The MIMIC-IV results are very similar to the University of Utah
786 results, reinforcing a key message of this paper - that PBC retains the power to identify
787 comorbid relationships that are lost by stratification.

788 **Table S1. University of Utah patient demographics.** Total number of patients is
789 1,604,818.

Gender		Ethnicity	
Female	868758	Not hispanic	894511
Male	815855	Hispanic or latino	176944
Unknown	192	Unknown	613350
Ancestry			
Caucasian	959848	Asian	30952
Unknown/other	644805	African American	25074
Native Hawaiian or Pacific Islander	14445	American Indian or Alaska Native	9681
Financial Class (insurance)			
Commercial	912535	Medicaid/uninsured	379651
Medicare	159668	Other/Unknown	152,964

790

791 **Table S2. Overlap between comorbidities plotted in Figure 2.** Each line in Figure 2
792 consists of 9 points. For every two adjacent points in these lines, we calculate the size
793 of the intersection of the left point (A) and the right point (B) and divide by the size of B.
794 In set notation, we calculate $|A \cap B| / |B|$, i.e., the percent of B that is contained in A. As
795 seen in this table, the comorbidities discovered as features are added - are almost
796 entirely subsets of the model without the feature. Values of nan are present when B is
797 empty.

A	B	PBC	Binomial, female 50-59	PBC	Binomial, female 50-59	PBC	Stratification
0	1	100%	99%	100%	96%	100%	96%
1	2	99%	97%	98%	89%	98%	83%
2	3	98%	98%	97%	95%	96%	nan
3	4	98%	98%	97%	nan	96%	nan
4	5	99%	81%	98%	nan	96%	nan
5	6	99%	99%	98%	nan	96%	nan
6	7	99%	99%	99%	nan	97%	nan
7	8	99%	83%	97%	nan	95%	nan

798

799 **Table S3. PBC retains known breast cancer comorbidities lost by the stratification**
800 **approach.** Comorbidities passing a Bonferroni corrected alpha threshold of **1.08e-8** are
801 colored blue. Stratification by sex and age (considering only females in their 50s)
802 eliminates significant associations for most known comorbidities of breast cancer, with
803 the exception of endometriosis, a comorbidity of uncertain significance. PBC retains
804 statistical power while modelling the effects of confounding variables. Our approach
805 does not consider endometriosis as a comorbid condition of breast cancer.

Representative Comorbidities of Breast Cancer: p-values			
Potential Comorbidity	binomial all data	binomial, female 50-59	PBC
Hypertension	1e-1680	3.2e-8	1e-63
Osteoarthritis	1e-1491	2.5e-5	1e-101
Cancer of colon	1e-265	1.3e-6	1e-25
Melanomas of skin	1e-165	1.6e-8	1e-15
Cancer of kidney	-1e-108	5e-4	1.6e-13
Endometriosis	1e-106	1e-15	3.2e-6

806 **Table S4. Replication of Table 1 on MIMIC-IV data.** See legend for Table 1. MIMIC-IV
807 data does not include patient age, ethnicity or insurance type so these rows are omitted.
808 Only 1 patient in the MIMIC-IV data set has both multiple myeloma and multiple
809 sclerosis. For the other two known comorbidites - the trend is clear - PBC retains power
810 as additional features are added to the model. Stratification results in a loss of statistical
811 power.

Stratification filters	PBC features	Concussion and migraine		Multiple myeloma and Multiple sclerosis		Cancer of pancreas and hypertension	
		P-values					
		χ^2	PBC	χ^2	PBC	χ^2	PBC
No filters (n=248,714)	none	-2.46	-2.46	-0.00	-0.00	-420	-420
female (n=128676)	+sex	-2.43	-2.46	-0.00	-0.00	-183	-419
+Caucasian (n=78092)	+ancestry	-1.50	-3.02	-0.00	-0.00	-117	-391
+3yr history (n=10239)	+span	0.0	-2.38	-0.00	-0.00	-1.88	-331

812

813 **Table S5. Replication of Table 2 on MIMIC-IV data.** See legend for Table 2. While the
 814 dataset is too small to measure a significant comorbidity between Sickle Cell Anemia
 815 and Malaise and Fatigue using either method, we still see that PBC retains statistical
 816 power lost by stratification. For very small sample sizes, stratification is not an option.

Stratum filters	PBC features	Sickle Cell Anemia paired with Malaise and Fatigue		
		Stratum Pair count	χ^2 p-value	PBC p-value
No filters (n=477,070)	no features	9	-2.00	-2.00
Caucasian (n=276,496)	ancestry	0	0	-0.97
African American (n=7,035)	ancestry	8	-0.67	-0.97
+Female (n=2,789)	+gender	5	0	-1.5

817

818 **Table S5: Several co-occurrence and out-of-window p-values identified by PBC.**

819 Every p-value shown below passes a Bonferroni corrected alpha threshold of **1.08e-8**.

820 Each pair of terms exhibits a significant temporal trend as indicated by arrows.

T1	T2	Comorbid log10 p-value	Out-of-window Direction	Directional log10 p-value
Milrinone (1,676)	Heart Transplant (122)	1e-104	→	1e-32
T2D (97,484)	CKD (15,123)	1e-3374	→	1e-195

HIV (2,730)	Pneumocystis (203)	1e-120	→	1e-16
Obesity due to excess calories (26494)	Essential Hypertension (-1299)	1e-848	→	1e-1299

821 **Table S7: Citations corresponding to comorbidities discovered by minimum**
822 **description length clustering.** See Figure S3 for a graphical representation of MDL
823 clusters.

Term 1	Term 2	PubMed ID
abruptio placenta	fetal distress	26393335
hyperemesis gravidarum	preeclampsia	23360164
Hodgkins lymphoma	chemotherapy	26541251
pancreatitis	pancreatic cancer	30315287
Chronic asthma	diabetes	30489598
aplastic anemia	renal failure	17426071
allergic reactions	eye disorders	31343437
constipation	urinary tract infection	30212423
melanoma	Parkinson's disease	29991141
HIV	osteoporosis	25709813
chronic sinusitis	substance disorder	28812909
breech	early labor	31741046
atrial fibrillation	emphysema	25900353
osteoarthritis	retinal detachment	20462780
schizophrenia	hypertension	27855222
developmental disability	birth trauma	31240076
lumbago	depression	31703727

migraine	obesity	27358118
----------	---------	----------