

1 Topic modelling with ICD10-
2 informed priors identifies novel
3 genetic loci associated with
4 multimorbidities in UK Biobank

5 Yidong Zhang^{1,2,3}, Xilin Jiang^{1,4,5,6,7,8}, Alexander J
6 Mentzer^{1,5}, Gil McVean^{1*,#}, Gerton Lunter^{9,10*}

7
8 **Affiliations**

9 1 Big Data Institute, Li Ka Shing Centre for Health Information and Discovery, University of
10 Oxford, Oxford OX3 7LF, UK

11 2 CAMS China Oxford Institute, Nuffield Department of Medicine, University of Oxford,
12 Oxford OX3 7BN, UK

13 3 Department of Radiation Oncology, Peking Union Medical College Hospital, Chinese
14 Academy of Medical Sciences and Peking Union Medical College, Beijing 100006, China

15 4 Department of Statistics, University of Oxford, Oxford OX1 3LB, UK

16 5 Wellcome Centre for Human Genetics, Nuffield Department of Medicine, University of
17 Oxford, Oxford OX3 7BN, UK

18 6 Department of Epidemiology, Harvard T.H. Chan School of Public Health, Boston 02115,
19 MA, USA

20 7 British Heart Foundation Cardiovascular Epidemiology Unit, Department of Public Health
21 and Primary Care, University of Cambridge, Cambridge CB2 0SR, UK

22 8 Heart and Lung Research Institute, University of Cambridge, Cambridge CB2 0BB, UK

23 9 MRC Weatherall Institute of Molecular Medicine, John Radcliffe Hospital, University of
24 Oxford, Oxford OX3 9DS, UK

25 10 Department of Epidemiology, University Medical Center Groningen, University of
26 Groningen, Groningen 9700 RB, The Netherlands

27 *These authors jointly supervised the work: gil.mcvean@bdi.ox.ac.uk, g.a.lunter@umcg.nl

28 #Lead contact

29

30 Summary

31 Studies of disease incidence have identified thousands of genetic loci
32 associated with complex traits. However, many diseases occur in combinations
33 that can point to systemic dysregulation of underlying processes that affect
34 multiple traits. We have developed a data-driven method for identifying such
35 multimorbidities from routine healthcare data that combines topic modelling
36 through Bayesian binary non-negative matrix factorization with an informative
37 prior derived from the hierarchical ICD10 coding system. Through simulation
38 we show that the method, treeLFA, typically outperforms both Latent Dirichlet
39 Allocation (LDA) and topic modelling with uninformative priors in terms of
40 inference accuracy and generalisation to test data, and is robust to moderate
41 deviation between the prior and reality. By applying treeLFA to data from UK
42 Biobank we identify a range of multimorbidity clusters in the form of disease
43 topics ranging from well-established combinations relating to metabolic
44 syndrome, arthropathies and cancers, to other less well-known ones, and a
45 disease-free topic. Through genetic association analysis of inferred topic
46 weights (topic-GWAS) and single diseases we find that topic-GWAS typically
47 finds a much smaller, but only partially-overlapping, set of variants compared to
48 GWAS of constituent disease codes. We validate the genetic loci (only)
49 associated with topics through a range of approaches. Particularly, with the
50 construction of PRS for topics, we find that compared to LDA, treeLFA
51 achieves better prediction performance on independent test data. Overall, our
52 findings indicate that topic models are well suited to characterising
53 multimorbidity patterns, and different topic models have their own unique
54 strengths. Moreover, genetic analysis of multimorbidity patterns can provide
55 insight into the aetiology of complex traits that cannot be determined from the
56 analysis of constituent traits alone.

57

58 Key words

59 Multimorbidity, topic modelling, treeLFA, topic-GWAS, UK Biobank

60

61

62 Introduction

63 Multimorbidity, defined as the co-existence of multiple chronic conditions, is a major challenge
64 for modern healthcare systems. Its prevalence has increased because of a worldwide increase
65 in life expectancy¹⁻³, and it is associated with substantially lower quality of life^{3,4}, worse clinical
66 outcomes³, and increased healthcare expenditure⁵. The management of multimorbidity is
67 challenging given that most guidelines and research are still targeted at single diseases. As a
68 result, the negative impact of multimorbidity is often greater than the additive effects of
69 individual diseases⁶.

70 Several common multimorbidity patterns, such as a cluster composed of cardiovascular and
71 mental health disorders, and a musculoskeletal disease cluster, have been identified from
72 literature reviews^{3,7}. In recent years, the widespread adoption of electronic health records
73 (EHR) has enabled the systematic study of multimorbidity, and a variety of approaches have
74 been employed for this purpose, including factor analysis⁸, clustering⁹, graph or network
75 based methods^{10,11}, and statistical models such as latent class analysis¹²⁻¹⁴. These
76 approaches have both validated the previously identified multimorbidity patterns^{12,14,15} and,
77 through the inclusion of a wider range of diseases, identified additional multimorbidity patterns
78^{14,16}. In addition, downstream analyses enabled by these approaches have helped to identify
79 the clinical events and outcomes associated with specific multimorbidity patterns^{13,17}, which
80 may provide insights about early intervention and risk stratification for patients.

81 The existence of common multimorbidity patterns raises the question of their etiology. One
82 way to approach this question is to analyse multimorbidity patterns together with appropriate
83 -omics data to determine the biological pathways involved. These analyses have been made
84 possible with the establishment of biobanks linking individuals' biological samples and genetic
85 information to their EHR¹⁸⁻²⁰. A recent study investigating genome-wide association studies
86 (GWAS) of 439 common diseases recorded in UK biobank (UKB) hospital inpatient data found
87 that 46% multimorbidity disease pairs have evidence for shared genetics¹¹, suggesting that
88 this may be a fruitful approach.

89 Intrinsic to the study of multimorbidity is the joint analysis of multiple disease phenotypes. To
90 enable this, various multi-trait GWAS methods have been developed which promise to better
91 exploit the deep phenotype data available for individuals in biobanks. These methods can be
92 subdivided based on their analytical approaches. Univariate methods combine signals from
93 single-trait GWAS²¹⁻²⁷, while multivariate methods offer improved power by directly modelling
94 the individual level genotype and phenotype data²⁸⁻³³. Several of these methods use
95 transformations such as principal components analysis (PCA) on the original traits before
96 association analysis so that very large numbers of traits can be handled. Topic models such
97 as Latent Dirichlet Allocation (LDA) and Non-negative Matrix Factorization (NMF) were
98 dimension reduction algorithms developed to model word occurrence in text documents, and
99 have subsequently found application in biological studies to extract complex patterns from
100 high-dimensional data³⁴⁻³⁶. They can be used to find multimorbidity clusters from diagnostic
101 data, by viewing individuals as "documents" and diseases as "words". The topics learnt by

102 these models then are mathematical representations of groups of diseases that tend to co-
103 occur within the same individual³⁷. Earlier studies have shown that joining single diseases
104 into topics increases statistical power for genetic discovery, and helps to disentangle the
105 pleiotropic effects of several known genetic loci^{38–40}.

106 Despite these advances, existing methods all have limitations. First, diagnostic data is often
107 binary in nature, with zeros and ones representing the absence and presence of diseases, yet
108 topic models like LDA and NMF were designed for count data, while algorithms designed for
109 binary data⁴¹ have not found wide application in biomedical studies. Second, biobank data is
110 often sparse and, particularly for less common diseases, inclusion of prior domain knowledge
111 may well improve results. Domain knowledge has been used successfully in topic models^{42,43},
112 and, for example, medical ontologies like the ICD-10 disease classification system could serve
113 as prior for disease co-occurrence, as they encode the complex relationships of diseases as
114 a hierarchical structure which is amenable to mathematical analysis^{44–46}. Third, while for
115 statistical models, such as LDA, principled approaches exist for selecting the number of
116 clusters and optimising other hyperparameters^{47–51}, this is often not true for other methods,
117 and these choices can strongly impact the final results^{9,52}. In addition, methods not based on
118 statistical foundations typically lack estimates of uncertainty in the inferred clusters, which
119 makes interpretation difficult.

120 Here, we develop and validate an analytic framework for the study of multimorbidity using topic
121 models and multi-trait GWAS on biobank datasets. Central to our approach is “treeLFA” (latent
122 factor allocation with a tree-structured prior), a statistical model to identify multimorbidity
123 clusters of common diseases based on co-occurrence patterns and an informed prior derived
124 from a tree-structured disease ontology. Applying treeLFA to Hospital Episode Statistics (HES)
125 data extracted from UKB we gain insights about the relationships of diseases and their shared
126 genetic components. We identify multimorbidity clusters in the form of disease “topics” and
127 show that these agree with accepted medical understanding. Performing a series of GWAS
128 on the quantitative traits defined by individuals’ weights for these topics (topic-GWAS), we
129 show that the approach identifies novel loci that correlate in expected ways with several
130 genomic annotations. We validate the topic-GWAS results using test data, and show that
131 topic-GWAS can improve genetic risk prediction for multiple disorders, in particular immune
132 disorders, and those for which currently few associated loci are known.

133

134

135 Results

136 Overview of treeLFA

137 treeLFA is a topic model designed to identify multimorbidity clusters from binary disease
138 diagnosis data. We describe the model in terms of the associated generative process. To
139 generate data, first topic vectors ϕ_k containing disease probabilities, and a topic weight vector
140 θ_d for each individual d are sampled from a prior distribution. An individual's disease
141 probabilities are given by a mixture of different topic vectors, with the topic weights (θ_d) acting
142 as mixing proportions. The model can equivalently be defined by the likelihood for
143 observations, which involves factoring a latent matrix of disease probabilities:

$$144 \quad P(W|\theta, \Phi) = \prod_{d,s} \text{Bernoulli}(W_{d,s} | [\theta \Phi]_{d,s})$$

145 Here, W is the binary input matrix recording individuals' diagnosed diseases, with rows
146 representing individuals and columns representing disease codes; $P(W_{d,s}=1)$ is the
147 probability of disease s being 1 (diagnosed) for individual d . ϕ is the topic-disease matrix
148 (each row a topic and each column a disease), and θ the topic weights matrix (each row an
149 individual and each column a topic). The occurrence of disease s for individual d is modelled
150 with a Bernoulli distribution parameterized by the corresponding entry in the product of
151 matrices θ and ϕ : $[\theta \Phi]_{d,s}$. This model differs from LDA in three ways. First, LDA samples
152 diseases (or words) according to a multinomial distribution, so that diseases can occur
153 multiple times, while treeLFA only allows presence or absence. Second, LDA conditions on
154 the number of observed diseases, whereas for treeLFA the number of diseases is
155 informative. Third, treeLFA uses an informative prior on topic vectors ϕ_k guided by a tree-
156 structured ontology such as ICD-10 (Figure 1 in the analytic note). This prior has the
157 property that diseases that are closely related on the tree tend to have correlated
158 probabilities. Inference on the treeLFA model is performed using partially collapsed Gibbs
159 sampling⁵³, integrating out the topic weight variable. See the analytic note for more details,
160 including on hyperparameter optimization.

161

162 Validation of treeLFA; Comparison with related topic models

163 We assessed treeLFA's performance in a simulation experiment, comparing it to the same
164 model but without an informative tree prior (flatLFA; Fig 2a), and to LDA. We designed the
165 simulation to test the model with respect to the degree of multimorbidity in the data; the size
166 of the data; and the correctness of the prior. The degree of multimorbidity was governed by α ,
167 the concentration parameter of the Dirichlet prior for the topic weight variable θ , with large
168 values corresponding to the presence of several multimorbidity clusters in individuals, and
169 small values resulting in individuals mainly presenting diseases from a single cluster. To test
170 the influence of prior misspecification we used two sets of topics for simulation. In one set
171 ("correct tree prior") the active disease codes in topics were aligned with the tree structure of
172 disease codes, resulting in a high likelihood under the prior, which specifies that child nodes
173 on the tree tend to (though not exclusively) have the same activity as their parent nodes (Fig
174 2b,c). In the other setting ("incorrect tree prior") the pattern of active disease codes in topics

175 was possible but unlikely under the prior (Supplementary Fig 1a,b). For each set of topics we
176 considered four combinations of D and α , resulting in eight parameter combinations in total
177 (Supplementary Table 1). For each of these we generated 20 data sets. To evaluate the
178 performance of treeLFA, flatLFA and LDA we used two metrics: the accuracy of the inferred
179 topic ($\Delta\phi$, average absolute per-disease difference in probability between aligned true and
180 inferred topics; see Methods for details), and R_{pi} , the ratio of the average per-individual
181 predictive test likelihood for treeLFA and flatLFA.

182 Both metrics indicate that on data simulated using the correct tree prior, treeLFA performs
183 better than flatLFA, which does not have the benefit of an informative prior (Figures 2d-g,
184 Supplementary Table 2). This is most pronounced for small datasets with strong
185 multimorbidity (Figure 2d; $\Delta\phi$ 0.012 ± 0.004 (treeLFA) and 0.025 ± 0.010 (flatLFA); R_{pi}
186 1.003 ± 0.002), while for larger datasets, the two models show similar performance and the
187 prior has less influence (Figure 2e; $\Delta\phi$ 0.009 ± 0.006 (treeLFA) and 0.011 ± 0.005 (flatLFA)).
188 treeLFA outperforms LDA except for large data sets with weak multimorbidity (Figure 2g; $\Delta\phi$
189 0.0100 ± 0.0009 (treeLFA) and 0.0084 ± 0.0021 (LDA)) where both models gave accurate
190 inference. For simulation using incorrect tree priors, flatLFA gave results comparable to
191 simulation with correct priors, and the performance of flatLFA and treeLFA is similar across
192 the four parameter combinations (Supplementary Figure 1), indicating that treeLFA is robust
193 against prior misspecification. Overall, these results indicate that both treeLFA and flatLFA
194 give accurate results when sufficient training data is available, even if treeLFA's informative
195 prior is inaccurate; but when the tree prior is correct, treeLFA performs better than flatLFA,
196 particularly when training data is limited. Even in larger real-world data sets, low-frequency
197 topics will have limited training data, hence this suggests that treeLFA could add power to the
198 analysis of multimorbidity in biobank data.

199

200 Topics of ICD-10 codes inferred from UK Biobank data

201 To investigate the properties of treeLFA on real-world data, we built an exploratory data set
202 using the HES data in UKB from 502,413 individuals, consisting of the 100 most frequent
203 codes from chapters 1-13 of the ICD-10 coding system (top-100 UKB dataset, Supplementary
204 Table 3). We split these data randomly into training (80%) and testing (20%) datasets, and
205 trained treeLFA with an initial $K=11$ topics (There is a discussion of the optimal number of
206 topics below).

207 The inferred topics include an “empty” topic, in which all codes have near-zero probability of
208 occurring (Figure 3a, Supplementary Table 4). Its associated entry in the optimal Dirichlet prior
209 parameter α is very large (0.585) compared to that of other topics (0.016-0.06), indicating that
210 the empty topic is frequently assigned to an individuals' disease profile. The remaining topics
211 all contain active codes. Most topics are sparse (8 topics contain fewer than 10 high-probability
212 (>0.2) codes), but the model also infers dense topics, such as topics 8 and 10 which include
213 41 and 43 high-probability codes respectively. To assess whether codes tend to be specific to
214 a topic, we normalised their probabilities across topics to make the largest probability 1. We
215 found that, in general, codes are typically specific to topics: most (87/100) are active
216 (normalised probability >0.5) in 3 or fewer topics. However, some codes are active in many
217 topics, such as I10 (essential hypertension, active in 6 topics) and C44 (other malignant
218 neoplasms of skin, active in 8 topics) (Figure 3a), suggesting that they have both large

219 prevalence and a large number of multimorbidity partners belonging to different disease
220 clusters. The top disease codes (i.e. the five with the largest probabilities) in the 10 non-empty
221 topics are consistent with known disease mechanisms (Figure 3b), and are most frequently
222 drawn from two (6/10) of the ICD-10 chapters. Specifically, in Topic 5, codes E78 (disorders
223 of lipoprotein metabolism and other lipidemias) and I10 (essential hypertension) are
224 components of the metabolic syndrome⁵⁴, which is associated with increased risk for
225 cardiovascular diseases (CVD)⁵⁵, an association supported by the other three inferred top
226 codes for this topic (I20, I21 and I25, all heart diseases). Another example is Topic 11, whose
227 top codes include four spondylopathy subtypes, while the remaining one is G55 (nerve root
228 and plexus compressions), a common complication of intervertebral disk disorders.

229 In addition to defining topic vectors, the model also infers individuals' weights for all topics
230 (shown for 2,000 individuals in Figure 3c, Supplementary Table 5). As expected, most
231 individuals have substantial weight for the empty topic, and this weight was strongly and
232 negatively associated (Pearson correlation: -0.853) with the total number of diagnoses.
233 Individuals that were not diagnosed with any of the top-100 ICD-10 codes (629/2,000) have a
234 weight near 1 for the empty topic, while the majority of other individuals (1056/1371) have
235 large weight (>0.1) for less than two disease (non-empty) topic, as expected from the sparsity
236 of the data.

237 To compare the performance of treeLFA with flatLFA and LDA, we used the same input data
238 to train the flatLFA model with 11 topics and the LDA model with 10 topics (no empty topic
239 would be inferred by LDA, so it was trained with one fewer topic). Topics inferred by treeLFA
240 and flatLFA were almost identical (Supplementary Figure 2a), indicating that the input data
241 was large enough to make the impact of the informative prior minimal. Most topics inferred by
242 LDA also had a high level of similarity to the non-empty topics inferred by treeLFA, except for
243 two topics, for which the cosine similarities were 0.685 and 0.853 (Supplementary Figure 2b).
244 Overall, these results indicate that the three topic models captured the same multimorbidity
245 pattern from the top-100 UKB dataset.

246

247 GWAS on topic weights

248 We next investigated whether the quantitative traits defined by the topic weights can be used
249 to identify genetic variants that are associated with an individual's risk for developing
250 multimorbidities represented by the topics. We performed GWAS on individuals' weights for
251 the 11 topics inferred by treeLFA (topic-GWAS), and identified associations that reached
252 genome-wide significance ($p < 5 \times 10^{-8}$; non-lead SNPs with $r^2 > 0.1$ were removed). For
253 comparison, we also performed standard binary GWAS for the 100 ICD-10 codes and the 296
254 Phecodes mapped from these ICD-10 codes (see Methods for details).

255 We found 128 independent loci associated with at least one of the 11 topics, while 812
256 independent loci were associated with at least one of the 100 ICD-10 codes; 82 loci were
257 shared between the sets (Figure 4a). Phecode GWAS showed similar patterns
258 (Supplementary Figure 3a,b). Breaking this down by topic, we find that unique loci found by
259 topic-GWAS were highly non-randomly distributed (Figure 4b, Supplementary Table 6). Most
260 unique loci were associated with the empty topic (20/36), followed by Topic 8 (17/28) which
261 contains a large number (41) of high probability codes (>0.2) from Chapter 11 (Diseases of
262 the digestive system, 12 codes) and 13 (Diseases of the musculoskeletal system and

263 connective tissue, 13 codes). In contrast, four sparse topics showed no unique loci. Topics 5
264 (metabolic and heart diseases) and 6 (joint diseases) had many associated loci (50 and 17
265 respectively) and also a substantial number of unique loci (5 and 7 respectively), suggesting
266 that active codes in these topics include shared genetic components. The identification of
267 novel loci indicates that topic-GWAS provides additional power for discovery. For example,
268 Figure 4c (Supplementary Table 7) compares P-values of lead SNPs from the topic-GWAS
269 for Topic 5, and P-values for association of the same loci with the top five active codes in Topic
270 5 (E78, I10, I20, I21, I25) from the single code GWAS. For most topic-associated lead SNPs,
271 P-values given by topic-GWAS are smaller than those given by the corresponding single code
272 GWAS, indicating increased power for these examples. This also explains some loci uniquely
273 found by topic-GWAS, including some loci that show single-code P-values well below genome-
274 wide significance (see Figure 4d for two example loci). Despite the limited numbers of topic-
275 associated loci, the genomic control inflation factor and LD score regression (LDSC) indicate
276 that most topics are in fact highly polygenic traits, with the exception of the empty topic and
277 Topic 8, for which LDSC analysis suggests that uncontrolled confounding factors exist
278 (Supplementary Table 8; age, sex and the first 10 PCs were controlled for in topic-GWAS).

279 We next asked whether topic-GWAS simply identify associations with disease groups
280 (categories) represented by internal nodes of the ICD-10 or the Phecode ontology tree, which
281 correspond to expert-led disease clusters and provide a useful contrast to our data-driven
282 multimorbidity clusters. To answer this we performed GWAS on groups of ICD-10 codes or
283 Phecodes corresponding to internal nodes in the respective classification systems. We found
284 634 loci associated with the 68 internal ICD-10 codes and 296 loci associated with the 136
285 internal Phecodes. Of the 128 topic-associated loci, 41 were not associated with any of the
286 internal or terminal ICD-10 codes; and for Phecodes the corresponding number was 56
287 (Supplementary Figure 3c-f). This indicates that topic modelling provides insights into the
288 relationships of diseases beyond those provided by expert-driven disease groupings encoded
289 in ontologies. For example, Topic-8 has the majority of its active codes coming from Chapter
290 11 (Diseases of the digestive system) and 13 (Diseases of the musculoskeletal system and
291 connective tissue), and a similar multimorbidity cluster was also identified by a recent study
292 on UK Biobank ¹¹. Interestingly, this cluster has many unique loci found by topic-GWAS,
293 possibly indicating that these two categories of diseases share some underlying biology.

294 We then compared the topic-GWAS results for topics inferred by treeLFA, flatLFA and LDA.
295 As expected from the similarity of topics inferred by treeLFA and flatLFA, similar numbers of
296 associated loci were identified (128 and 126; Supplementary Figure 4a), compared to LDA,
297 which identified many fewer (65; Supplementary Figure 4b), of which 44 overlap with the
298 treeLFA loci. This difference in numbers of associated loci is mostly due to treeLFA's empty
299 topic (associated with 36 loci), which is not identified by LDA, and also due to differences in
300 the dense Topic 8 (treeLFA, 28 loci; LDA, 4 loci), and topics 5 and 6 (Supplementary Figure
301 4c). One reason for the relatively poor performance of LDA may be that LDA-derived topic
302 weights are negatively correlated with each other, as they must sum to one, while treeLFA's
303 topic weights are only negatively correlated with the empty topic weight, but are otherwise
304 almost independent of each other (Supplementary Figure 5).

305 **Validation of topic-GWAS results**

306 To exclude the possibility that the unique topic-GWAS associations were driven largely by

307 technical biases or population stratification, we validated the results in three ways. First we
308 considered overlap with previously identified loci reported in the GWAS catalogue⁵⁶. We find
309 that 114/128 (89.1%) of all topic-associated loci and 36/46 (78.3%) of unique associations
310 have records in the GWAS catalogue, and this overlap is consistent across topics
311 (Supplementary Figure 6a). Second, we looked at enrichment of topic-associated loci in
312 functional genomic regions. To do this we defined three groups of SNPs, including lead SNPs
313 for all loci associated with ICD-10 codes, a random selection of 10,000 GWAS tag SNPs
314 (controls) and topic-associated lead SNPs that were not found by single code GWAS, and
315 then compared the proportions of them that have different functional properties (two-proportion
316 Z-test, adjusted P-value<0.05, Bonferroni correction). We find that compared to random SNPs,
317 a significantly larger proportion of topic-associated SNPs are in genomic regions with strong
318 transcription activity (using chromHMM-predicted chromatin states as proxy⁵⁷). In addition,
319 the proportions of SNPs that are QTLs and have chromatin interactions (CI) in at least one
320 tissue in the first and third groups are similar (0.83 and 0.77 are eQTL, 0.90 and 0.98 have
321 CI), and larger than the corresponding proportions in the control group (0.50 and 0.65 for eQTL
322 and CI), indicating that loci associated with single codes and topics have comparable
323 functional properties which are different from those for controls. (Supplementary Figure 6b-d).

324 Third, we made use of the test data to validate topic-GWAS results. We reasoned that if topics
325 and their associated loci represent true biological processes, then topic-GWAS results should
326 enable us to predict the risks of individual diseases with an accuracy comparable to that
327 achieved using single code GWAS. To do this we first constructed PRS for topic weights using
328 topic-GWAS results on training data. We found that they all show significant association with
329 inferred topic weights on test data (Supplementary Table 9, see Methods for details). We then
330 used these PRS for topics to construct PRS for the 100 ICD-10 codes, by adding individuals'
331 PRS for the ten disease topics weighted by the probability of the ICD-10 code of interest in
332 each topic. For comparison, we also constructed PRS for all ICD-10 codes directly using the
333 single-code GWAS results in the standard way. Each pair of PRS for an ICD-10 code was
334 used to predict individuals' corresponding diagnosed disease in the test data, and their
335 performance was evaluated using the area under the receiver-operator curve (AUC) statistic.
336 For 65 ICD-10 codes, topic-PRS AUCs are larger than single-code PRS AUC (Figure 4f,
337 Supplementary Table 10), with increases ranging from 1% to 5%. This increase was seen
338 most for ICD-10 codes from chapters 5 (Mental and behavioural disorders, 75% (3/4) showing
339 increased AUC), 11 (Diseases of the digestive system, 86% (18/21)) and 13 (Diseases of the
340 musculoskeletal system and connective tissue, 70% (14/20)). By contrast, single-code PRS
341 performed well for codes that have a relatively large number of associated loci found by single
342 code GWAS (>10 associated loci; 18/22 disease codes show larger AUC for single-code PRS
343 than topic-PRS; Supplementary Figure 7a,b). Finally, to make an objective comparison of the
344 topic-GWAS for treeLFA and LDA, we constructed PRS for ICD-10 codes from LDA's topic-
345 GWAS using the same approach, and found that in the majority of cases (99/100) the PRS
346 based on treeLFA's results have larger AUC (Supplementary Figure 7c), indicating treeLFA's
347 topic-GWAS is more informative. Taken together, the three complementary approaches
348 indicate that topic-GWAS associations broadly represent true genetic associations with
349 biological phenotypes.

350 Inference and topic-GWAS results across models

351 Before applying treeLFA to a larger data set containing more diseases, we considered how to
352 select the number of topics (K), a fundamental problem for topic models. We trained treeLFA
353 models with different numbers of topics (K=2-20, 50, 100) on the top-100 data set, and found
354 that for larger K topic vectors were frequently duplicated, therefore we performed clustering
355 on the posterior samples of topics (see Methods for more details). The resulting topics always
356 included an empty topic, and as K increased the topics tended to become more sparse,
357 although some dense topics always remained (Figure 5a, Supplementary Table 11). As K
358 increased, topics tended to split into sub-topics, which we visualised in a tree by connecting
359 each topic to its most similar topic (measured by Pearson correlation of topic vectors) in the
360 layer above, and we observed that the topics split in a stable way (Figure 5b, Supplementary
361 Table 12). These observations indicate that topic-GWAS loci and associated effect sizes
362 should also be stably identified. We verified this for many loci (examples in Supplementary
363 Figure 8), and Figure 5b illustrates this for a single variant. On the top-100 UKB dataset, the
364 number of distinct topics remaining after clustering is saturated at 25-30 topics
365 (Supplementary Figure 9a). Similarly, the total number of topic-GWAS loci, the number of
366 unique such loci, and the predictive likelihood on the test data all began to saturate beyond
367 K=20 (Figure 5b; Supplementary Figure 9b). We do note that for models with K=50 or K=100,
368 we infer several near-empty topics after clustering, which are unlikely to be stable
369 multimorbidity patterns and are challenging to interpret. Taken together, these results indicate
370 that selecting a sufficiently large value for K, combined with post-hoc clustering of topics, is a
371 computationally efficient strategy for producing a stable and comprehensive set of topics.

372 Results on a larger UKB dataset

373 We next defined a larger data set consisting of the 436 ICD-10 codes from chapters 1-14 with
374 a prevalence exceeding 0.001 in UKB (top-436 UKB dataset, Supplementary Table 13), and
375 again randomly split this 80/20 into training and testing datasets. Training treeLFA/flatLFA
376 models with 100 topics, we identified about 40 distinct topics after clustering of posterior
377 samples. Therefore, we kept 40 topics for both models (for the convenience of an objective
378 comparison of predictive likelihood), and collapsed the remaining near-empty topics into the
379 empty topic.

380 Since the inference results (topics) given by different treeLFA/flatLFA chains were not exactly
381 the same, we used the result given by the chain with the largest predictive likelihood on the
382 test data for the downstream analyses. Among the 40 inferred topics, 29 were found by all
383 three treeLFA chains, and five were found by two treeLFA chains, suggesting most topics
384 were stably identified from the data (Figure 6a, Supplementary Table 14). The 40 topics again
385 include several dense topics, with topics 1-5 including more than 40 active codes (having a
386 normalised probability>0.5 in a topic). The set includes many sparse topics, with most
387 including active codes enriched (Fisher exact test, adjusted P values<0.05, FDR corrected)
388 for 1-2 ICD-10 chapters, and again a single empty topic (Figure 6a-b; Supplementary Table
389 14,15). The top active codes in topics (defined as having an unnormalized probability>0.3) are
390 shown in Table 1, where topics are annotated based on the categories of these top active
391 codes. For most topics, their top active codes represent similar diseases, such as diseases

392 affecting the same physiological system or having the same pathological mechanism.
393 Comparing topics identified by treeLFA and flatLFA, we found that 32 topics were identified
394 by both models (cosine similarity>0.9; Supplementary Figure 10a), while the remaining topics
395 have substantial differences. Overall, the predictive likelihood of treeLFA chains was better
396 than that of flatLFA, and has a smaller range (Supplementary Figure 10b), indicating that the
397 tree-based prior is indeed helpful in extracting meaningful patterns from the data.

398 We then performed topic-GWAS on the 40 treeLFA and flatLFA topics we kept. We found 278
399 treeLFA and 260 flatFLA genome-wide significant loci, with the majority (207) found in both
400 sets and associated with corresponding topics (Supplementary Figure 10c-d, Supplementary
401 Table 16). We also performed single-code GWAS on the 436 ICD-10 codes and found 1,093
402 associated loci, among them 198 were also associated with treeLFA topics. Lead SNPs for
403 loci only associated with treeLFA topics (80 unique loci) had smaller effect sizes (median
404 absolute effect size: 0.021) compared to loci supported by both topics and single codes (0.024)
405 (Figure 6d, Supplementary Table 17), indicating that topic-GWAS enabled the discovery of
406 variants with small effects on multiple related diseases. Unique loci were not uniformly
407 distributed across topics; as in the top-100 dataset, many were associated with the empty
408 topic (21). Topics 3 and 30 also have large proportions of unique loci (81.5% and 83.3%), and
409 most of their active codes are from Chapter 13 (Diseases of the musculoskeletal system and
410 connective tissue) (Figure 6e, Supplementary Table 18). Other topics that are associated with
411 substantial numbers of unique loci are shown in the Supplementary Table 19.

412 We validated the topic-associated loci for the larger dataset using the same approach as used
413 with the top-100 dataset. Overall, 89.2% (248/278) of topic-associated loci and 78.9% (63/80)
414 of unique associations have records in the GWAS catalogue. The functional validation results
415 were also similar to that on the top-100 dataset, with unique topic-associated loci and single
416 code associated loci exhibiting similar profiles (Supplementary Figure 11). The two types of
417 PRS for single codes were also constructed using single code and topic-GWAS results. For
418 130 in 436 (30%) ICD-10 codes, PRS based on topic-GWAS resulted in larger AUC on the
419 test data. These codes were mainly from Chapter 3 (Diseases of the blood and blood-forming
420 organs and certain disorders involving the immune mechanism, 5/12), Chapter 9 (Diseases of
421 the circulatory system, 21/48) and Chapter 13 (Diseases of the musculoskeletal system and
422 connective tissue, 24/52). In contrast, for most (93 of 109) ICD-10 codes with more than 5
423 associated loci, PRS based on single code GWAS resulted in better performance
424 (Supplementary Table 20). To compare treeLFA and flatLFA, we also constructed PRS for
425 single codes based on flatLFA results, and compared their AUC on the test data to that of
426 treeLFA. For 231 in 436 codes (53%), PRS based on treeLFA showed better performance.
427 Supplementary Figure 10e compares the density plots for the AUC of PRS for all codes given
428 by the two methods, where treeFLA shows a minor advantage.

429

430 Discussion

431 Multimorbidity is a major challenge for today's healthcare systems, yet our understanding of it
432 remains limited⁵⁸. The establishment of biobanks linked to electronic health records presents
433 an opportunity for a more systematic study of multimorbidity, and highlights the need for
434 reliable and powerful analytic tools to enable the identification of major multimorbidity clusters
435 and downstream analyses of paired phenotype and -omics data.

436 Here we developed treeLFA, a topic model inspired by Latent Dirichlet Allocation (LDA) that
437 admits a prior for topics constructed on existing tree-structured medical ontologies. We
438 compared it to flatLFA and LDA on both simulated and UKB data, and found that the prior was
439 effective at extracting relevant topics from limited input data, such as data involving rare
440 diseases. We also found that the novel model structure better fits the binary input data,
441 resulting in the identification of an empty ("healthy") topic, and ensuring that topic weights for
442 the remaining disease topics were largely uncorrelated, improving power for downstream
443 topic-GWAS. We implemented algorithms to optimise hyperparameters of the model, and
444 developed a computationally efficient approach to determine the number of meaningful topics
445 to be inferred.

446 By applying treeLFA to HES data for 436 common diseases recorded in UKB, we identified 40
447 topics reflecting combinations of diseases that tend to co-occur. These topics varied in density
448 and include a single empty topic, many sparse topics that each include a small number of
449 active disease codes, and several dense topics. We found that the inferred topics were largely
450 consistent with the current disease classification system (ICD-10), yet treeLFA also combines
451 diseases distant on the tree structure into the same topic, indicating the utility of supplementing
452 expert-led knowledge systems with data-driven methods. By comparing the topics inferred by
453 treeLFA and LDA, and topics inferred by different treeLFA models, we found that
454 multimorbidity clusters were consistently inferred, indicating the existence of stable
455 multimorbidity clusters, and that topic modelling of cross-sectional diagnosis data is an
456 informative method of finding them.

457 Most inferred topics likely reflect underlying aetiology, as indicated by the fact that the large
458 majority (34/40) show genome-wide significant associations with genetic markers, while 20
459 topics are associated with 80 novel loci that do not reach genome wide significance in GWAS
460 for single ICD-10 codes, and show evidence of functionality using multiple methods. The active
461 ICD-10 codes in the topics with the most novel associations are mainly from Chapter 13
462 (Diseases of the musculoskeletal system and connective tissue), with substantial contributions
463 also from Chapters 4 (Endocrine, nutritional and metabolic diseases), 9 (Diseases of the
464 circulatory system) and 14 (Diseases of genito-urinary system), suggesting that diseases in
465 these chapters share genetic risk factors.

466 With topic-GWAS results, we explored constructing PRS for a single code as the sum of PRS
467 for all topics weighted by the probabilities of the code in topics. We found that for certain codes,
468 especially codes with very few GWAS hits from Chapter 13, this new type of PRS outperforms
469 the standard PRS based on single code GWAS results. This improvement in prediction might
470 result from better estimation of effect sizes of variants by topic-GWAS. This is because
471 although treeLFA factorises the input matrix in a linear way, it achieves a dimension reduction
472 by mapping from the disease space to the topic space. This results in fewer traits and therefore

473 more ‘cases’ for each one, which is especially beneficial for the study of the highly polygenic
474 traits (topic weights), where there are a large number of SNPs with small effects.

475 In contrast to topic-GWAS, single-code GWAS on the 463 common diseases resulted in 1,093
476 significant associations, of which the vast majority (895) were not associated with any topic.
477 Taken together, these genetic analyses indicate that the majority of genetic associations are
478 driven through links to individual diseases. Meanwhile, most multimorbidity clusters also have
479 genetic bases, which are mainly composed of pleiotropic genetic variants affecting risks of
480 multiple active diseases in the cluster. Besides, there are also a substantial number of
481 associations for topics that are difficult to identify by single code GWAS due to the lack of
482 power. From a biological perspective, they might reflect the complex connections between
483 upstream pathways that are distant to individual diseases, or variants with direct yet small
484 effects across multiple diseases. Overall, these observations are a helpful starting point for
485 our pursuit of a deeper understanding of the mechanisms underlying multimorbidity clusters.

486 We consistently identified an “empty” topic, which reflects individuals’ overall disease burden,
487 since their weights for the “empty” topic are negatively correlated with their total number of
488 diagnosed diseases. Perhaps surprisingly we found that this topic showed many genetic
489 associations, many of which (21/35) had not been identified before. An enrichment analysis
490 for topic-associated genes (Methods) indicated several lifestyle-associated factors, such as
491 gym attendance and religious observance, suggesting that this topic may be related to
492 individuals’ health behaviour, which in part is genetically determined. In contrast, for most
493 sparse disease topics, the enriched gene sets are usually directly related to the active codes
494 in the topic (Supplementary Table 21).

495 There have been many studies aiming at identifying multimorbidity patterns using various
496 methods ^{12,13,16,17,59}. Most of these studies focus on dozens of diseases that in addition varied
497 from study to study, making comparisons difficult. We compared the topics identified by
498 treeLFA with multimorbidity networks (which are interpretable at the level of genetic loci) found
499 by a recent study using the HES data for 439 common diseases in UKB (433 of these are
500 included in the top-436 UKB dataset in our study) ¹¹. For most of the disease networks, there
501 are specific corresponding disease topics identified by treeLFA (Supplementary Table 22).
502 However, the overlap of active codes in treeLFA inferred topics and the disease networks is
503 limited, suggesting significant differences in the details of the inference results. This
504 discrepancy could be caused by the fundamental differences between the two methods, since
505 treeLFA analyses all diseases simultaneously, while multimorbidity networks were constructed
506 based on pairs of diseases that tend to co-occur. One limitation of topic models is that they
507 cannot determine the relationships between active diseases in the same topic. In contrast, an
508 advantage of topic models is that they make direct use of individual level data for the genetic
509 analyses, and allows for making predictions on the test data, which provides an objective way
510 to compare different methods.

511 Our work represents real progress in the understanding of multimorbidity, yet also reveals
512 important and unsolved challenges. For instance, while we showed that taken together the
513 novel genetic associations likely represent true biology, we have not performed individual
514 replication of the findings in this study in independent data sets and this may be challenging
515 unless the data sources and methods of data collection are comparable to UK Biobank. The

516 problem of inferring disease topics that are stable, tractable and biologically meaningful across
517 geographies and healthcare systems represents a major challenge for future research.

518 Acknowledgements

519 This work uses data provided by patients and collected by the NHS as part of their care and
520 Support. Funded by the Chinese Academy of Medical Sciences (CAMS) China Oxford
521 Institute (COI) and Studentship from the China Scholarship Council (CSC) (YZ).
522 Computation used the Oxford Biomedical Research Computing (BMRC) facility, a joint
523 development between the Wellcome Centre for Human Genetics and the Big Data Institute
524 supported by Health Data Research UK and the NIHR Oxford Biomedical Research Centre.
525 We thank Chris Holmes for the discussion.

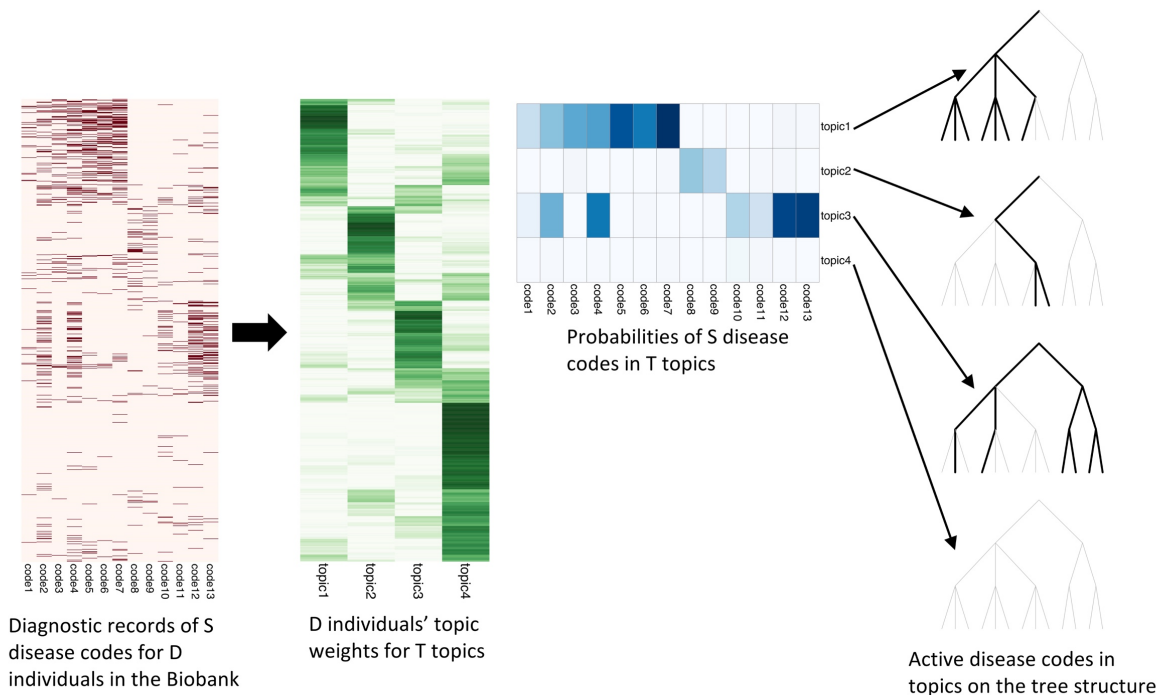
526 Author Contributions

527 Y.Z.: conceptualization, data curation, formal analysis, methodology, software, validation,
528 visualisation, writing-original draft; X.J.: conceptualization, data curation, methodology,
529 software, writing: review&editing; A.J.M: conceptualization, project administration,
530 supervision, writing: review&editing; G.L: conceptualization, investigation, methodology,
531 project administration, resources, supervision, writing-original draft, writing: review&editing;
532 GM: conceptualization, funding acquisition, methodology, project administration, resources,
533 supervision, writing: review&editing;

534 Declaration of Interests

535 G.M. is a director of and shareholder in Genomics PLC, and a partner in Peptide Groove LLP.
536 G.L. is a shareholder in Genomics PLC. The other authors declare no competing financial
537 interests.
538

539 Figures

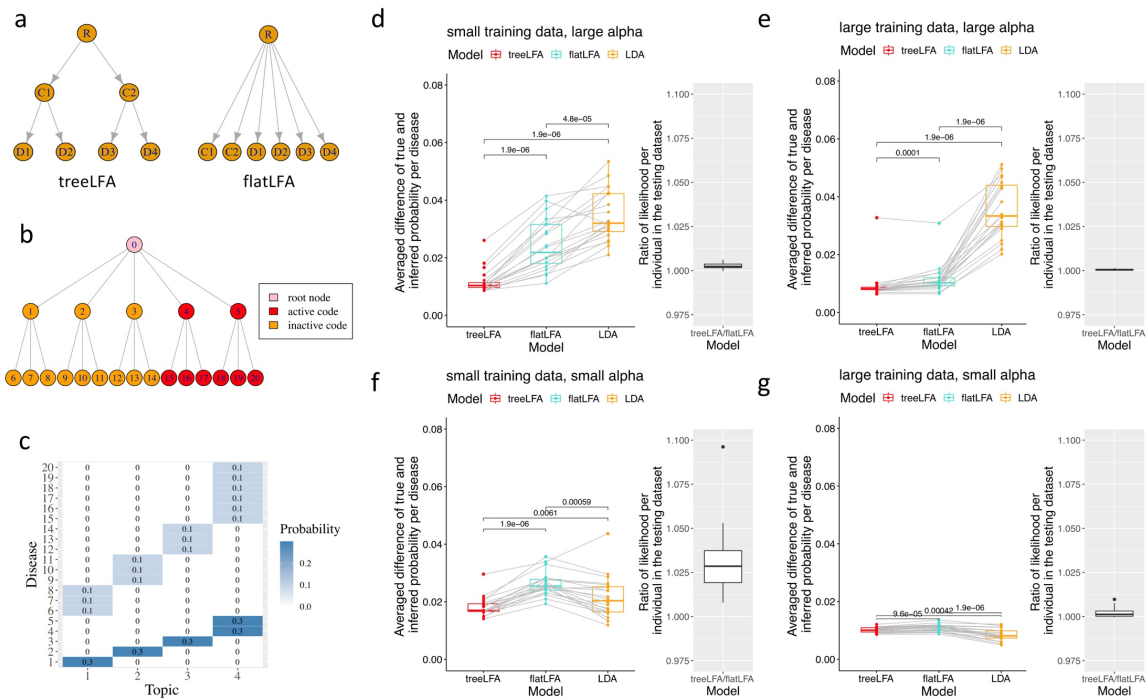


540

541 **Figure 1 Schematic for topic modelling of the diagnosis data in UKB with treeLFA.**

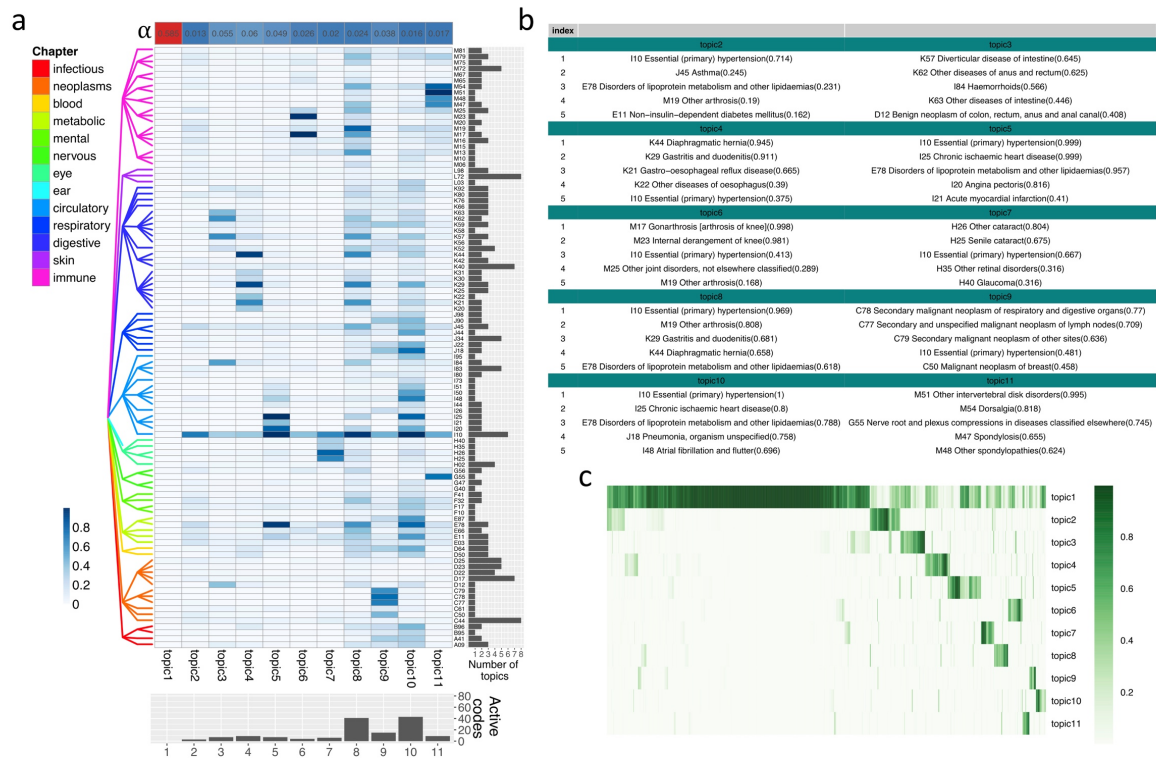
542 The presence and absence of S disease codes for D individuals in the biobank is modelled
 543 with $D \times S$ Bernoulli distributions. The matrix of Bernoulli probabilities is factored into the
 544 product of a topic matrix and a topic weight matrix. Individuals' weights for topics are
 545 modelled with categorical distributions with a Dirichlet prior. Each topic is composed of S
 546 probability variables with Beta priors to parameterize the Bernoulli distributions for disease
 547 codes. Disease codes can be either active or inactive in topics. Active disease codes have
 548 large probability while inactive ones have near zero probability. A prior for topics (specifies
 549 the likelihood of different disease codes to be active in topics) is constructed on the tree
 550 structure of disease codes specified by a medical ontology (such as the ICD-10 coding
 551 system). The path from the root node to active leaf nodes (corresponding to all active
 552 disease codes in a topic) are highlighted on a three-layered tree structure for 13 disease
 553 codes in 4 topics.

554



555
556 **Figure 2 Comparison of three related topic models (treeLFA, flatLFA and LDA) on**
557 **simulated datasets.**

558 a, The informative and non-informative tree structures used by treeLFA and flatLFA.
559 b, The tree structure of the 20 disease codes used for simulation. Red nodes on the tree correspond to
560 all active disease codes in Topic 4 used for simulation (in Figure 2c).
561 c, The four topics used for simulation. The heatmap shows the probabilities of diseases in topics. Each
562 row corresponds to a disease code, each column corresponds to a topic. Inactive disease codes in
563 topics have zero probability. All active disease codes in each of the first three topics come from the
564 same branch of the tree in Figure 2b. Active codes in the last topic come from the last two branches of
565 the tree in Figure 2b.
566 d-g, Comparison of three topic models on simulated datasets. The performance of three topic models
567 (treeLFA, flatLFA and LDA) on four groups of simulated datasets are shown. The four groups of datasets
568 were generated using the same topics (Figure 2c), and different values for D (number of individuals in
569 the training dataset) and α (the concentration parameter of the Dirichlet prior for individuals' topic
570 weights). For each combination of D and α , 20 datasets (including both training and testing datasets)
571 were simulated. Inference accuracy of topic models is evaluated using the averaged per disease
572 difference between true and inferred probability of all diseases in all 4 topics (box plots). Each dot in a
573 box plot is the result of one model on one dataset, and dots for different models on the same dataset
574 are connected with grey lines. For treeLFA and flatLFA, the predictive likelihood on the testing datasets
575 were calculated using topics inferred on the training data. Each dot in the point plot represents the
576 treeLFA to flatLFA ratio of per individual averaged predictive likelihood for one dataset. d, Results on
577 datasets simulated using D=2500 and $\alpha=1$. e, Results on datasets simulated using D=5000 and $\alpha=1$. f,
578 Results on datasets simulated using D=300 and $\alpha=0.1$. g, Results on datasets simulated using D=1000
579 and $\alpha=0.1$.
580



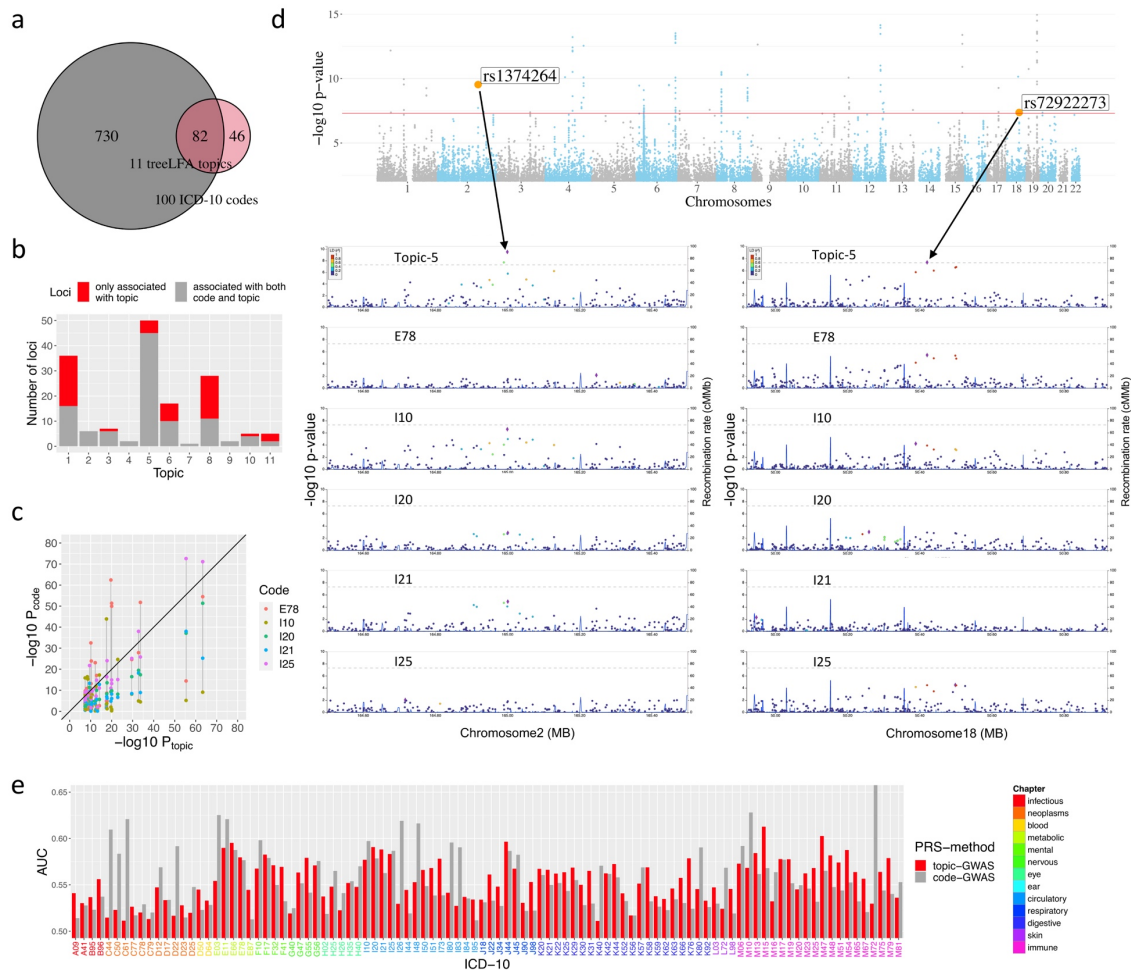
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595

Figure 3 Inference results given by treeLFA on the top 100 UKB dataset.

a, 11 topics inferred by treeLFA from the top-100 UKB dataset. The heatmap shows the probabilities of 100 ICD-10 codes in the 11 treeLFA topics, in which each row is an ICD-10 code and each column is a topic. Topics are arranged in a descending order of their corresponding entries in the optimised α vector (the single-row heatmap on the top). The tree structure for the 100 ICD-10 codes is shown to the left of the heatmap. Codes from different chapters of the ICD-10 coding system are colored differently. The barplot below the heatmap shows the numbers of ICD-10 codes with a probability of at least 0.2 in topics. The barplot on the right side of the heatmap shows the number of topics in which an ICD-10 code is active (with a normalised probability of at least 0.5).

b, The top 5 codes with the largest probability in the 10 non-empty treeLFA topics (topics 2-11 Figure 3a). Numbers in the brackets show the probabilities of disease codes in topics.

c, Inferred weights for the 11 topics for 2000 random individuals. Each row in the heatmap is a topic, and each column is an individual.



596

597 **Figure 4 topic-GWAS result for the 11 topics inferred by treeLFA.**

598 a, The total numbers of significant loci found by topic-GWAS for the 11 treeLFA topics and single code

599 GWAS for the 100 ICD-10 codes, and the overlap of these two sets of loci.

600 b, The numbers of significant loci found by both single code/topic-GWAS, and the numbers of loci only

601 found by topic-GWAS for the 11 treeLFA topics.

602 c, Comparison of P-values given by topic-GWAS for all lead SNPs for Topic 5 and P-values for the

603 same SNPs given by single code GWAS for the top 5 active codes (E78, I10, I20, I21, I25) in Topic 5.

604 d, The Manhattan plot for Topic 5, and the regional Manhattan plots for single code/topic-GWAS results

605 for two example lead SNPs of Topic 5.

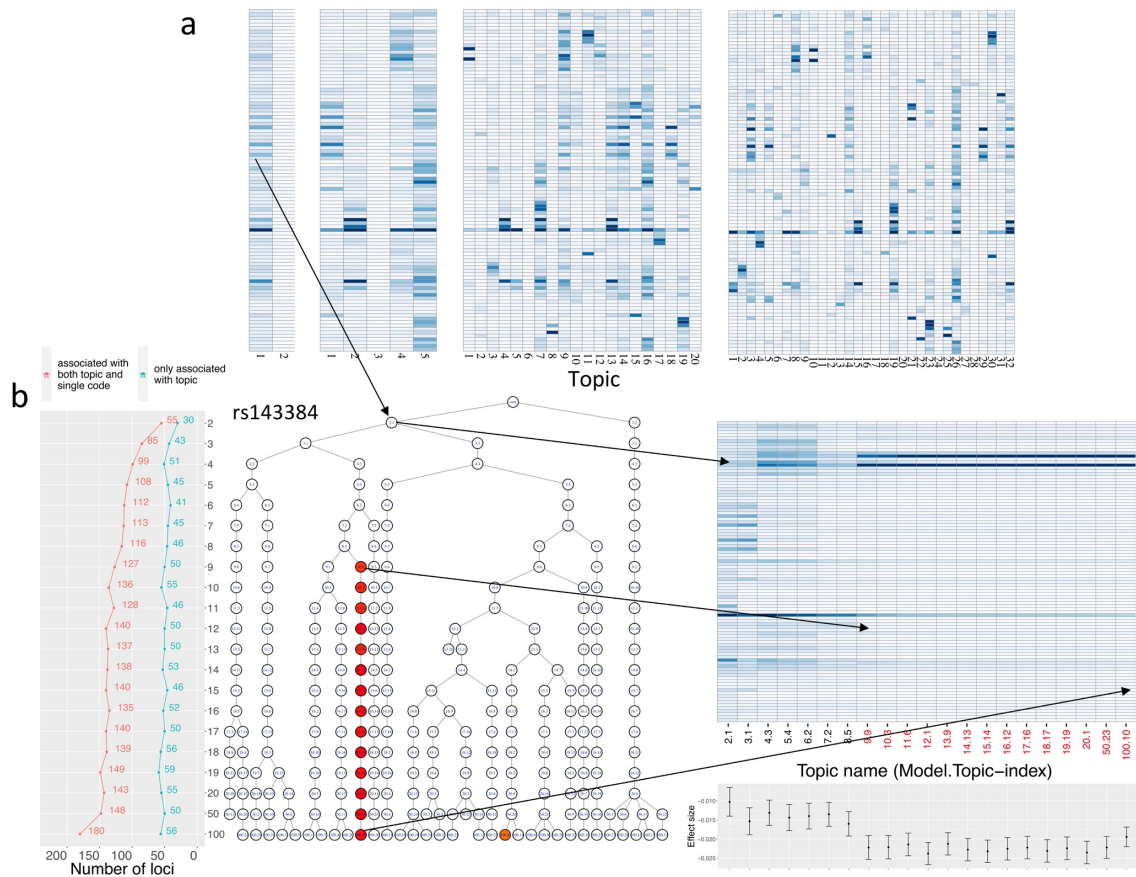
606 e, Comparison of two types of PRS for the 100 ICD-10 codes. One type of PRS is directly constructed

607 using the single code GWAS results. Another type of PRS for ICD-10 codes is constructed as the sum

608 of PRS for topics weighted by the probabilities of an ICD-10 code in all topics. The AUC of these two

609 types of PRS on the test dataset are plotted.

610

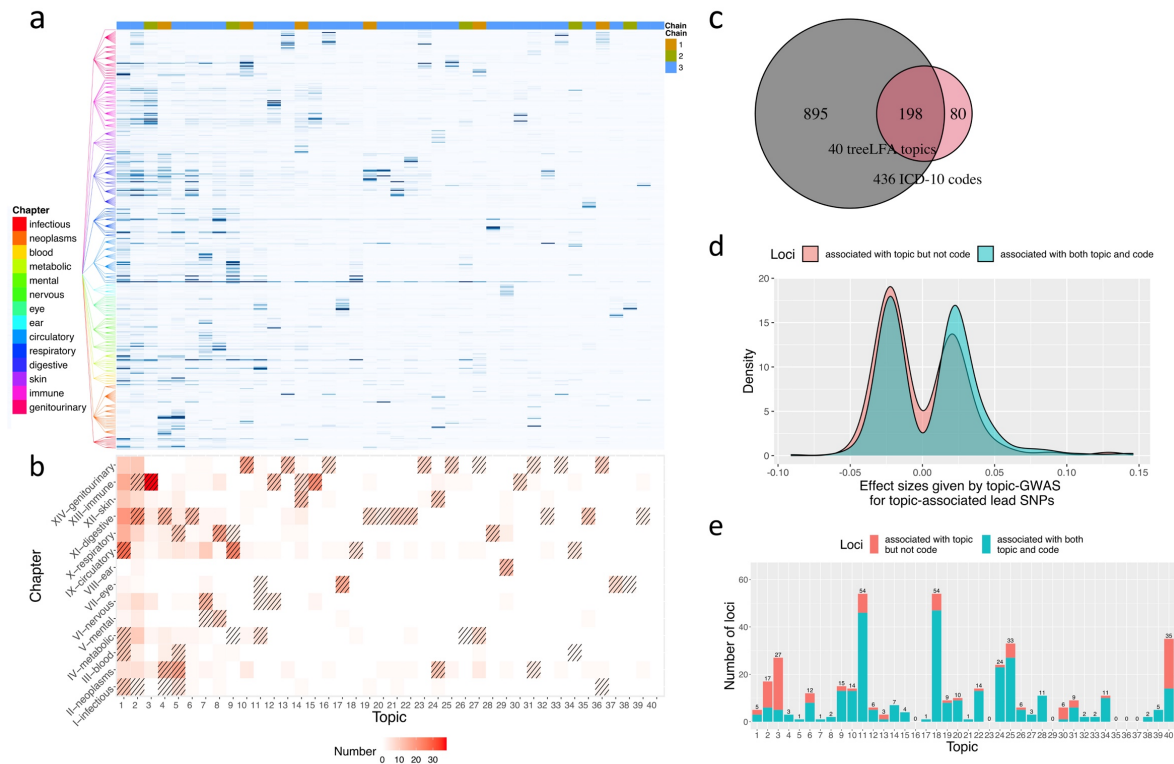


611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630

Figure 5 Inference and topic-GWAS results across treeLFA models set with different numbers of topics.

a, Averaged inferred topics given by the treeLFA models set with 2, 5, 20 and 100 topics. For each treeLFA model, ten Gibbs chains were trained, and 50 posterior samples of topics were collected from each chain. Posterior samples of topics from different chains are mixed and clustered, and the mean topic vectors are then calculated for each cluster.

b, Relationship of topics inferred by different models. All topics inferred by all treeLFA models are organised into a tree structure. Each node on the tree is a topic inferred by a model, and all nodes on the same layer (level) of the tree are all topics inferred by the same treeLFA model (set with a certain number of topics). Each topic in the tree is connected to its most similar topic (measured with the Pearson correlation) in the layer above. Topics associated with the SNP rs143384 are colored according to the $-\log_{10}(P\text{-value})$ for the SNP in the corresponding linear regression. Most of the associated topics are in the same branch of the tree, so all topics in this branch are plotted in the heatmap on the right side of the tree, with names of topics (model.topic-index) associated with rs143384 highlighted in red. In the barplot below the heatmap, effect sizes and standard errors of rs143384 given by topic-GWAS for the above topics are plotted. The line plot to the left of the tree shows the total numbers of topic-associated loci and the numbers of topic-associated loci that are not found by single code GWAS for different treeLFA models.



631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647

Figure 6 Inference and topic-GWAS results for the 40 topics inferred by treeLFA from the top-436 UKB dataset.

a, The 40 topics inferred by the treeLFA model set with 100 topics. Topics are ordered according to their density (the sum of the probability of all codes in the topic). The tree structure of the 436 ICD-10 codes is plotted to the left of the heatmap. The colour bar on top shows for each topic the number of treeLFA chains which inferred it.

b, The numbers of active codes (with a normalised probability of at least 0.5) in different topics coming from different ICD-10 chapters. Enriched chapters among active codes in topics are highlighted with shades in the cells (Fisher exact test, FDR<0.05).

c, The total numbers of loci associated with the 40 treeLFA topics and the 436 ICD-10 codes, and the overlap of these two sets of loci.

d, Distribution of effect sizes given by topic-GWAS for lead SNPs associated with only topic and lead SNPs associated with both topic and single code.

e, The total numbers of loci associated with different topics (in the same order as topics in Figure 6a), and the numbers of topic-associated loci that are not found by single code GWAS (red).

648 Tables

649 **Table 1 Names and top active codes of the 40 topics inferred by treeLFA on the** 650 **top-436 UKB dataset.**

651 Top ICD-10 codes (with an unnormalized probability of at least 0.3) in the 40 topics inferred by treeLFA
652 on the top-436 UKB dataset are shown. Topics are named using the ICD-10 chapters which make major
653 contributions to their top active codes (different chapters are connected with “-” in the names). Words
654 after colons in the names of topics give further summary of the top active codes from a specific chapter.
655 Topic 1 and 2 are not named since they are very dense, and there is no obvious pattern in their top
656 active codes.

Topic	Name	Top active codes
1	-	-
2	-	-
3	Musculoskeletal:joint-1	M19,I10,M17,M25,M13,M54,K21,E78,M47,M23,K44,E66,M15,M79,J45,K29,M16,M75,K57,G56,F32,E11,M51,M65
4	Neoplasm-Digestive	C78,C18,C77,I10,K56,D64,K63,C79,K59,K57,K66,C20,J90,K91,N39,K62,J18,K52,E87,N17,D12,B96,K43,K29,E86,A41
5	Neoplasm-Respiratory-Blood	C79,C78,D70,C77,J18,A41,J90,C34,I10,D64,J22,K59,C50,E87,N17,N39,A09,C80,J98,K52
6	Digestive-Circulatory	I10,I25,E78,K29,I20,K44,K57,K21,E11,K62,K63,I84,D12,K22,D50,J45,D64,K20,K92,J44,K31,I48
7	Nervous:cerebrovascular	I10,I63,E78,G81,I67,N39,I69,G40,F32,I48,J18,K59,E87,H53,B96,G45,J22,G93
8	Respiratory-Mental	J44,J18,F17,F32,J45,I10,F10,J22,F41,J43,J98,J47,J96,E87
9	Circulatory:heart-Metabolic-1	I10,I25,I48,I50,E78,I51,I20,I44,I34,I21,I08,I47,I35,J90,J18,I49,I42,E11
10	Urinary:male-1	I10,N32,N40,N39,E78,C67,N30,N35,N13,E11,N20,C61,K57
11	Metabolic:diabetes-Eye	E11,I10,E78,E10,H26,H36,E14,I25,H25,E66
12	Musculoskeletal:spine	M51,M54,G55,M47,M48,I10,M79,M43,M50
13	Reproductive:female-1	D25,N83,N92,N80,N73,N94
14	Skin:infection	L03,I10,L40,B95,L02

15	Musculoskeletal:joint-2	M19,M20,M25,I10,M75,M16
16	Reproductive:female-2	N95,N84,N85,I10,D25,N81
17	Eye-1	H26,H25,I10,H35,H40,H33,H52
18	Circulatory:heart-Metabolic-2	I25,I10,E78,I20,I21
19	Digestive:lower GI-1	K52,K51,K62,K63,K50,K57,I84,A09
20	Digestive:lower GI-2	K57,D12,K63,K62,I84,I10
21	Digestive:upper GI	K44,K29,K21,K22,K20,K31,I10
22	Digestive:hepatobiliary	K80,I10,K81,K82,K85
23	Urinary:female	N39,N81,N32
24	Neoplasm:skin-Skin	C44,D22,L57,L98,L82,I10,C43
25	Urinary:male-2	N40,N32,I10,C61,K40
26	Metabolic	I10,E78,I48,E11
27	Renal:kidney	N20,I10,N23,N13
28	Respiratory:upper	J34,J33,J32,J45
29	Ear:hearing	H72,H91,H90,H65,I10
30	Musculoskeletal:knee	M23,M17,I10
31	Neoplasm:breast	C50,C77,D05
32	Digestive:lower GI-3	I84,K62,K57
33	Reproductive:female-4	N92,D25,N84

34	Circulatory:peripheral vascular	I80,M79,I83,I26
35	Digestive:teeth	K02,K08,K01,K04
36	Reproductive:female-5	N87,N84,N95
37	Eye:appendicle	H02,H04
38	Eye-2	H26,H33
39	Digestive:hernia	K40
40	Healthy	Empty

657

658 **Methods**

659 Full technical details for treeLFA are given in the analytical note in the supplemental data.

660 **Validation of treeLFA with simulated data**

661 The simulation study served two purposes. Firstly, we aimed to verify that treeLFA can
662 accurately infer latent topics from the diagnosis data encoded as binary variables. Secondly,
663 we aimed to compare the performance of the three topic models (treeLFA, flatLFA and LDA)
664 discussed in the next section, such that the influence of model structure and treeLFA's
665 informative prior for topics on the inference can be assessed.

666 **Overview of the three related topic models**

667 1. treeLFA

668 treeLFA ("latent factor allocation with a tree structured prior for topics") is a topic model
669 designed for binary diagnosis data in biobanks based on the Bayesian mean-parameterized
670 binary non-negative matrix factorization⁴¹. It models the presence and absence of S disease
671 codes for D individuals with $D \times S$ Bernoulli distributions, and the matrix of Bernoulli probability
672 for all disease variables is factored into the topic matrix (ϕ) and the topic weight matrix (θ).
673 The loading of disease code s in topic t (ϕ_{ts}) is its corresponding Bernoulli probability in the
674 topic, and the Bernoulli probability of the disease variable for individual d and disease code s
675 is a mixture of the Bernoulli probability of code s in all topics, with the mixing coefficients
676 specified by the topic weights for individual d. treeLFA also incorporates an informative prior
677 for topics constructed by running a Markov process on a tree structure of individual words,
678 which assumes that disease codes on the same subtree are likely to have similar Bernoulli
679 probability in the same topics (see the details in the analytic note).

680 2. flatLFA

681 flatLFA has the same model space as treeLFA, with the only difference being that flatLFA
682 uses a non-informative prior for topics, which is constructed on a tree structure where all
683 nodes (representing all disease codes in a topic) are placed directly under the common root
684 node. By comparing the performance treeLFA and flatLFA, the contribution of treeLFA's
685 informative prior for topics to the inference can be assessed.

686 3. LDA

687 Latent Dirichlet Allocation (LDA⁶⁰)'s model configuration is different from treeLFA. LDA only
688 models the disease codes that are diagnosed for individuals with categorical distributions.
689 For LDA, each topic is a categorical distribution (or a Multinomial distribution if the input data
690 is viewed as a count matrix) across the S disease codes, and a Dirichlet prior distribution is
691 used to generate these topics. By contrast, for treeLFA each topic is a sequence of Bernoulli
692 probability for the S disease codes, and S Beta distributions are used to generate the topics.

693 **Description of simulated data**

694 We simulate multiple data sets using different topics and hyperparameters to assess the
695 performance of the three topic models (treeLFA, flatLFA and LDA) in different situations. We

696 simulate the input data sets in two steps. Firstly, we build a tree structure for 20 disease
697 codes (Figure 2B). The tree structure has three layers. The first layer is the root node; the
698 second layer contains five nodes, and each of them has three children nodes in the third
699 layer. Secondly, we generate topics of disease codes using the tree structure above. A
700 Markov process on the tree structure is used to do this (see the analytic note), and the
701 rationale for choosing the parameter of the Markov process is explained below.

702 The Markov process chooses active disease codes (disease codes having large probability
703 in a topic) for each topic by generating binary indicator variables for disease codes (with 1
704 represents active codes, and 0 represent inactive codes in a topic) (see the analytic note).
705 The Markov process has two transition probabilities, ρ_{01} and ρ_{11} , which control the sparsity of
706 topics ($\rho_{01} = P(I = 1 | I^{\text{parent}} = 0)$, $\rho_{11} = P(I = 1 | I^{\text{parent}} = 1)$). Small values for both ρ_{01} and ρ_{11} give
707 rise to sparse topics (because most codes in a topic will be inactive), while large values for
708 both generate dense topics. ρ_{11} also controls the clustering of active codes in topics. With a
709 large ρ_{11} , most children codes of an active parent code will be active. As a result, active codes
710 in a topic will gather on the same branch of the tree. By contrast, if ρ_{11} is small, active codes
711 will spread across the entire tree. For our simulation, we use topics resembling those
712 generated with small ρ_{01} and large ρ_{11} . This reflects our belief that in the real world most topics
713 of disease codes should be sparse (thus we chose small ρ_{01}), and that active disease codes
714 in the same topic tend to come from the same subtree (thus we chose large ρ_{11}).

715 For our simulation, we construct two sets of topics manually. The first set of topics are likely
716 to be generated using a Markov process with small ρ_{01} and large ρ_{11} (hyperparameter setting
717 used for inference on the simulated dataset), while the second set of topics are unlikely to be
718 generated by this Markov process. The first set of topics are used to test if the tree structure
719 of codes improves inference accuracy, and the second set of topics are to test the robustness
720 of treeLFA's inference when the tree structure of codes is wrongly specified. We manually
721 specified the topics to ensure that they are completely distinct from each other and have
722 strong patterns with respect to the clustering of active codes.

723 Figure 2B shows the first set of topics. For the first three topics, all codes on one branch of
724 the tree are active, and the remaining codes are inactive. In the last topic, all codes from two
725 branches of the tree are active. These topics are likely to be generated using a Markov
726 process with small ρ_{01} and large ρ_{11} , since a parent code and all its children codes are always
727 in the same state (either active or inactive). In the second simulation setting, active diseases
728 in topics are not generated according to their adjacency on the tree (Supplementary Figure
729 1B). We construct these topics by switching a fraction of active codes between topics in the
730 first simulation setting. As a result, active parent codes always have inactive children codes,
731 and inactive parent codes always have active children codes.

732 We simulate disease data using the topics described above and the generative process of
733 treeLFA. For each dataset, we split the data into training and testing data of the same size.
734 To evaluate the topic models in different situations in each topic setting, we use four
735 combinations of two hyperparameters to simulate data: α (the concentration parameter of the
736 Dirichlet prior for topic weights θ) and D (the number of individuals in the training dataset). A
737 large value for α means that most individuals will have large topic weights spread across
738 topics. By contrast, a small value for α will make topic weights for each individual more

739 concentrated on a single topic. A large α makes the inference difficult, since most individuals
740 are a mixture of multiple topics. Therefore, data sets simulated using large α require larger D
741 for accurate inference. For each hyperparameter and topic setting, we simulate 20 datasets.
742 Supplementary Table 1 summarises the hyperparameter and topic settings.

743 Implementation of the inference procedure

744 LDA is implemented using the R package “topicmodels”, and collapsed Gibbs sampling is
745 used to do the inference. treeLFA and flatLFA are implemented from scratch by us using the
746 R package “RcppParallel” and “Rcpp”.

747 For the hyperparameters of treeLFA and flatLFA, we provide the true value of α to train the
748 three models. Beta priors for the probability of active and inactive disease codes in topics (ϕ)
749 are Beta(2,4) and Beta(0.3,80). Beta priors for the transition probability (ρ_{01} and ρ_{11}) of the
750 Markov process are Beta(4.8,20) and Beta(20,4.8) (treeLFA). For flatLFA only ρ_{01} will be used
751 on the tree, and its prior is Beta(7,20), resulting in approximately the same expected number
752 of active codes in topics as treeLFA. For LDA we try a few different values (0.01,0.1,1) for η ,
753 the concentration parameter of the Dirichlet prior for topics. We find that with $\eta = 0.01$ LDA
754 has the best performance evaluated by the inference accuracy.

755 For the initialization of hidden variables for treeLFA and flatLFA, we initialised all indicator
756 variables (l) as 0, and then simulate all probability variables (ϕ) using the Beta prior for
757 inactive disease codes. Topic assignment variables (Z) are randomly sampled for all
758 individuals.

759 For each simulation scenario, ten Gibbs chains were sampled, and 20 posterior samples of
760 hidden variables were collected from each Gibbs chain with an interval of 100 iterations after
761 15,000 burn-in iterations.

762 Evaluation and comparison of topic models on simulated datasets

763 Two metrics are used to evaluate the models in simulations. The first metric is the inference
764 accuracy, measured with the averaged per disease difference between true and inferred

765 topic loadings: $\Delta\phi = \frac{\sum_{k=1}^K \sum_{s=1}^S |\phi_{ks}^{\text{true}} - \phi_{ks}^{\text{infer}}|}{K \times S}$, where ϕ_{ks} is the Bernoulli probability of disease

766 code s in topic k . To reorder the inferred topics as the true topics, we match each inferred
767 topic to the true topic that has the highest cosine similarity, in a greedy procedure (i.e. once
768 a true topic is matched, it is removed from the matching of the next inferred topic). The
769 pairwise t-test is used to test for statistical difference between two different models. The
770 second metric is the predictive likelihood on the test data (see the analytic note for more
771 details). For each posterior sample of topics, 200 Monte-carlo samples of topic weight θ are
772 used to approximate the predictive likelihood. A sensitivity analysis is done to ensure the
773 number of samples for θ is enough to have a stable estimate of the predictive likelihood.

774 Inference on the top-100 UK Biobank dataset

775 The input data for treeLFA

776 The input to treeLFA includes the diagnosis data for individuals in the UK Biobank, and the
777 tree structure for disease codes. The diagnosis dataset is constructed from the Hospital
778 Episode Statistics (HES) data in the UK Biobank, which is coded using the five-layered
779 hierarchical ICD-10 billing system. The first layer of the ICD-10 tree structure is the root node;
780 the second layer is composed of chapters of diseases coded using capital English letters; the
781 third layer contains blocks of disease categories; the fourth layer contains single disease
782 categories; and lastly, the bottom layer contains sub-categories of diseases, which can be,
783 for instance, the same disease occurring at different sites of human body, or subtypes of the
784 same disease. In UK Biobank, most of the diagnosed diseases are encoded using codes on
785 the bottom layer (fifth layer) of the tree. We use the fourth layer of encoding as diagnoses,
786 where we replace all diagnoses with their parental code in the fourth layer.

787 The top 100 most frequent ICD-10 codes in UK Biobank from the first 13 chapters of the ICD-
788 10 coding system are chosen to construct the dataset. This selection of chapters provides a
789 balance between breadth of phenotype and depth within any one chapter so that the potential
790 benefits of treeLFA can be explored. The diagnosis data is a binary matrix, with each row
791 represents an individual, and each column a disease code. Zeros and ones in the matrix are
792 used to represent the absence and presence of diagnosed ICD-10 codes for individuals. If
793 an individual is diagnosed with the same disease code several times, we keep only one
794 record to avoid bias of repeated diagnoses. The full dataset is randomly split into a training
795 dataset and a testing dataset, containing the diagnosis record for 80% and 20% individuals.

796 The tree structure of disease codes is encoded in a table with 2 columns: the first column
797 contains all the ICD-10 codes on the tree, and the second column records the parent codes
798 of the corresponding codes in the first column (Supplementary Table 3).

799 Implementation of treeLFA

800 Training strategy for treeLFA

801 The Gibbs-EM algorithm is firstly used to optimise α in two stages. In the first stage we run
802 2,000 iterations of the Gibbs-EM algorithm. In the E-step of each iteration, we run the Gibbs
803 sampler for treeLFA for 20 iterations and collect one posterior sample of Z (19 burn-in Gibbs
804 sampling iterations before the collection of the posterior sample). In the M-step, α is optimised
805 using this single posterior sample of Z collected in the E-step. In the second stage we continue
806 to run the Gibbs-EM algorithm for 200 iterations. In the E-step of each iteration, we run the
807 Gibbs sampler for treeLFA for 200 iterations and collect ten posterior samples of Z in total (19
808 burn-in Gibbs sampling iterations before the collection of each posterior sample). The reason
809 to have two stages of training is to balance the computational speed with the inference
810 accuracy. In the first stage, we optimise α more frequently and quickly get close to its optimal
811 value. In the second stage, α is more accurately optimised based on multiple posterior
812 samples of Z .

813 After optimising α , we use collapsed Gibbs sampler to simulate posterior distributions of all
814 hidden variables (Z, I, ϕ, ρ), with α fixed at and hidden variables initialised at the values
815 provided by the last iteration of the Gibbs-EM algorithm. 5,000 burn-in iterations are run before
816 the collection of posterior samples of hidden variables. For each topic model, ten Gibbs chains
817 are constructed, and 50 posterior samples are collected with an interval of 100 iterations from
818 each chain.

819 **Choices of hyperparameters and initialization of hidden variables**

820 To shorten the training with the Gibbs-EM algorithm, we initialise α as $(1, 0.1, \dots, 0.1)$. The first
821 entry in α is much larger than the others, and it corresponds to the empty topic that will always
822 be inferred from real-world diagnosis data. We also initialise α in other ways, such as using
823 $(1, \dots, 1)$, and we find that model converge to the same results regardless of the ways of
824 initialization of α , and the optimised α is usually more close to $(1, 0.1, \dots, 0.1)$ than other
825 choices. For topic assignment variable Z , we assign the empty topic (topic 1, corresponds to
826 the first entry in α) to all disease variables for individuals without any diagnosed disease
827 codes. For individuals with at least one diagnosed disease code, all topics are randomly
828 assigned to all disease variables. For topics, all indicator variable I are initialised as 0, and
829 probability variable ϕ are randomly sampled from $\text{Beta}(1, 5,000,000)$. $\text{Beta}(0.3, 80)$ and
830 $\text{Beta}(2, 4)$ are used as the prior for ϕ of inactive and active codes. $\text{Beta}(3, 20)$ and $\text{Beta}(3, 3)$
831 are used as the prior for transition probability ρ_{01} and ρ_{11} of the Markov process on the tree.
832 The hyperparameters for flatLFA are set in the same way as treeLFA. For LDA, the
833 concentration parameters of the Dirichlet priors for topic weights (α) and topics (η) are both
834 initialised as a vector of 0.1. α is not optimised since we find that this has negligible influence
835 on the inference result (inferred topics) and downstream analyses (topic-GWAS).

836 **Post-processing of inference result**

837 **Approximation of topic weights**

838 The topic weight variable θ is integrated out during the collapsed Gibbs sampling, therefore
839 their posterior samples need to be approximated using posterior samples of Z and α . θ can
840 be computed as in Griffiths and Steyvers⁵³: $\theta_{dt} = \frac{N_{dt} + \alpha_t}{N_d + \sum_{k=1}^K \alpha_k}$, where N_{dt} is the total number of
841 disease variables assigned with topic t for individual d , and N_d is the total number of disease
842 variables.

843 **Combining inference results from different Gibbs chains**

844 To combine the inference results given by different Gibbs chains, the “identifiability” issue
845 needs to be addressed, since the order of topics in different posterior samples from different
846 chains may not be the same.

847 We combine all posterior samples of topics from all chains together, and cluster topics before
848 taking the average within each cluster. To cluster topics from all samples, we firstly construct
849 a shared nearest neighbour (SNN) graph using the R package “scran”⁶¹. With the SNN graph,
850 we use the “Louvain” algorithm⁶², a community detection algorithm implemented in the R
851 package “igraph”, to assign topics into clusters. After clustering, similar topics coming from
852 different chains or posterior samples will be put into the same cluster. In addition to topics

853 (ϕ), we also assign posterior samples of other hidden variables (l , ρ and α) to the
854 corresponding clusters according to the clustering result for topics.

855 The Louvain algorithm doesn't allow us to directly specify the total number of clusters
856 (communities) to be found. Instead, the number of clusters is decided by the hyperparameter
857 k (the number of nearest neighbours to consider) for the construction of the SNN graph. A
858 large k will result in a small number of clusters, while a small k gives rise to a large number
859 of clusters, though some clusters might be alike. Empirically, we choose $k = N_{ch} \times \frac{N_{ps}}{2}$, where
860 N_{ch} is the number of Gibbs chains for the same treeLFA model, and N_{ps} is the total number of
861 posterior samples taken from each chain. This choice is to balance the total number of
862 clusters found by the algorithm and the uniqueness of different clusters.

863 **Post-processing inference results from models with a large number of topics**

864 For models set with a very large number of topics which far exceed the actual number needed
865 to explain the data (for instance, the treeLFA models set with 50 or 100 topics), multiple near-
866 empty topics (topics with few active codes having very small probability) will be inferred.
867 Although the small differences between these near-empty topics are not meaningful, they are
868 usually assigned to different clusters by the Louvain algorithm. To collapse these near-empty
869 topics into a single empty topic, we further apply hierarchical clustering on topics averaged
870 from different clusters given by the Louvain algorithm. During the hierarchical clustering,
871 similar topics are kept being combined until all the remaining topics are significantly different
872 from each other. By visualising the inferred topics using heatmap, one can roughly decide
873 the number of distinct meaningful topics (topics that are not empty or near-empty) to keep,
874 and then set the number of clusters (topics) to keep for the hierarchical clustering.

875 **Genetic analyses**

876 **topic-GWAS and single code GWAS**

877 To find genetic variants influencing individuals' risks for topics of diseases, we perform GWAS
878 using inferred topic weights as continuous traits (topic-GWAS). Since topic weights are real
879 numbers in the range of 0 to 1, the basic assumptions of linear regression do not hold. We
880 apply a logit transformation on topic weights to address this issue before fitting the standard
881 linear model for GWAS. We validate that using logit transformation gives better results than
882 using rank based inverse normal transformation and using no transformation on topic
883 weights. The validation is done by comparing the number of significant loci found by different
884 methods, and the predictive performance of PRS for single codes based on topic-GWAS
885 results (see section below). For topic-GWAS, we only include common SNPs (SNPs with a
886 minor allele frequency (MAF) larger than 0.01 in the UK Biobank) and individuals who self-
887 report having British ancestry in the training dataset (343,006 individuals in total). Sex, age
888 and the first ten principal components (PCs) of genomic variation are controlled for.

889 For comparison, we also perform GWAS (logistic regression) using the presence and
890 absence of single ICD-10 codes as binary traits (single code GWAS). The inclusion criterion
891 for individuals, SNPs and covariates are the same as topic-GWAS. In addition to ICD-10
892 codes, we also use terminal Phecodes mapped from the top 100 ICD-10 codes as traits for

893 single code GWAS. Phecodes are defined by systematically grouping terminal ICD-10 codes
894 into more applicable medical terms based on the judgements of clinicians and researchers
895 ⁶³, which reduces the granularity of terminal ICD-10 codes. Similar to ICD-10 codes, there is
896 also a hierarchical coding system for Phecodes. To map the ICD-10 codes used in the top-
897 100 UKB dataset to phecodes, we firstly extract all terminal ICD-10 codes (on the fifth layer
898 of the ICD-10 tree) that are children codes of the 100 level-4 ICD-10 codes, and then retrieve
899 their corresponding Phecodes according to the Phecode map
900 (<https://phewascatalog.org/phecodes>) In total, there are 296 terminal Phecodes mapped from
901 the 100 ICD-10 codes.

902 **Inflation in P-values given by topic-GWAS**

903 Inflation in P-values are observed for the topic-GWAS results given by all three topic models.
904 The inflation can either be resulted from true polygenicity of the traits (topic weights), or
905 stratification in the population. To differentiate these two possibilities, we carry out the LD
906 score regression (LDSC) ^{64,65} using the summary statistics of topic-GWAS for all topics. Pre-
907 computed LD scores (based on 1000 Genomes European data) are downloaded and used
908 in the analyses as recommended ⁶⁵. The genomic control inflation factor λ_{GC} and the intercept
909 of LDSC are output by the algorithm, and compared with each other. A large λ_{GC} and small
910 intercept for the same trait suggest true polygenicity causing the inflation in P-values, while
911 large values for both λ_{GC} and intercept suggest stratification in the population.

912 **Processing GWAS results**

913 To define genomic loci from significant SNPs ($P < 5 \times 10^{-8}$) found by GWAS, we use the
914 clumping function implemented in PLINK-1.9. $r^2 > 0.1$ is used as the threshold for clumping
915 SNPs in linkage disequilibrium (LD). We define a loci to be an association for both topic-
916 GWAS and single code GWAS as follows: the significant lead SNP found by one GWAS
917 method can be clumped with a significant lead SNP found by the other method.

918 **GWAS on internal disease codes on the tree**

919 In addition to grouping disease codes via topic modelling, we also group disease codes
920 completely following the medical ontologies (ICD-10 and Phecode systems). In other words,
921 we use internal codes (such as blocks of categories of diseases and chapters of disease,
922 corresponding to the nodes in the third and second layers of the ICD-10 coding system) of
923 the two disease classification systems as binary traits for single code GWAS. For instance, if
924 both disease codes A and B are under a common parent code C on the tree, then C will be
925 used as the trait for GWAS, and individuals who are diagnosed with either A or B will be used
926 as cases for the single code GWAS for code C. For the 100 ICD-10 codes there are 68
927 internal codes, and for the 296 Phecodes there are 136 internal codes.

928 **Comparison of topic-GWAS results for the three topic models**

929 In addition to topics inferred by treeLFA, topic-GWAS for flatLFA and LDA inferred topics are
930 also performed. For LDA, only individuals with at least one diagnosed disease code are used
931 as input for inference. For topic-GWAS, there are two options to deal with the individuals
932 without any diagnosis. We can either exclude them or include them and give them small

933 random weights for all disease topics. We experiment with both methods, and find that
934 excluding the completely healthy individuals results in a larger power for topic-GWAS.

935 Validation of topic-associated loci

936 **Validation using the GWAS Catalogue**

937 We check the GWAS Catalogue⁵⁶ to see if topic-associated loci were also found by previous
938 GWAS as significant. We download the full GWAS Catalogue⁶⁶, and clump all SNPs in it to
939 topic-associated lead SNPs ($r^2 > 0.5$ as threshold). If a topic-associated lead SNP found by us
940 can be clumped, it means that a SNP in LD with it was found by a previous GWAS as
941 significant.

942 **Validation using functional genomic resources**

943 Integrated analysis of GWAS results and functional genomic datasets has gained popularity
944 in recent years^{56,67}. Checking the enrichment of genomic annotations among topic-
945 associated loci (lead SNPs) is another angle of validation. We obtain various genomic
946 annotations for topic-associated loci using the software “FUMA”^{68,69}. Since most topics only
947 have a small number of associated loci, we combine all loci (lead SNPs) that are associated
948 with at least one topic and perform analyses on them as a whole. Meanwhile, we also perform
949 the same analyses on all single code associated lead SNPs and 10,000 random SNPs
950 sampled from all SNPs used in the GWAS (the distribution of their MAF are matched to all
951 topic-associated SNPs) for comparison. The assumption made here is that if topic-GWAS
952 find true associations, then the significant SNPs should have an enrichment profile that is
953 similar to single code associated SNPs (positive control) and different from randomly selected
954 SNPs (negative control).

955 Three types of functional annotations are used for the validation of topic-associated lead
956 SNPs. Firstly, the three groups of loci (lead SNPs) are annotated using the 15-core chromatin
957 states predicted by the chromHMM algorithm⁵⁷. Since the predicted chromatin states in the
958 127 available types of tissues are different, for each genomic locus we use the smallest
959 chromatin state across all tissues. Secondly, we calculate the proportions of lead SNPs in
960 the three groups that are eQTL (expression quantitative trait loci) in different tissues using
961 the eQTL mapping function implemented in FUMA, based on the GTEx8 dataset^{70,71}. Thirdly,
962 we calculate the proportions of lead SNPs having chromatin interactions with other genomic
963 regions in different tissues, based on the HiC data from the GSE87112 dataset⁷⁰. The default
964 setting of FUMA for parameters is used in all the analyses. The two proportion z-test is used
965 to test for significant differences between the proportions of two groups.

966 **Genetic risk prediction based on topic-GWAS results**

967 **PRS for topics**

968 Another way to validate topic-GWAS results is to use them for prediction tasks on the test
969 data. Because individual variants' effects on traits of interest are usually small, polygenic risk
970 scores (PRS) are constructed to aggregate the effects of tens of thousands of variants. With
971 topic-GWAS carried out on the training data, PRS for topic weights (traits of topic-GWAS) are

972 constructed using the software “PRSSice-2”⁷², which uses a “C+T” (clumping and
973 thresholding) method. No threshold for P-values is manually set for the inclusion of SNPs.
974 We use the test data to evaluate PRS for topics constructed on the training data. Topic
975 weights for individuals in the testing dataset are inferred by running the Gibbs sampler for
976 treeLFA on them, with ϕ and α fixed at values learnt from the training data (averaged from all
977 posterior samples from all Gibbs chains). Ten Gibbs chains are simulated to infer topic
978 weights for individuals in the testing dataset, and 50 posterior samples are collected from
979 each chain, and their average is used in the subsequent analyses. With inferred topic
980 weights, linear models are fit to evaluate the associations of PRS for topics and the
981 corresponding topic weights, using the logit transformed topic weights as response variables,
982 and PRS for topics as independent variables. The heritabilities of topic weights are estimated
983 using LDSC as a reference.

984 **PRS for single codes based on topic-GWAS results**

985 To evaluate topic-GWAS results using single code GWAS results as reference, and to
986 compare the topic-GWAS results for different topics models (such as treeLFA and LDA)
987 under a common criterion, we construct two types of PRS for single ICD-10 codes using
988 single code and topic-GWAS results, respectively. PRS based on single code GWAS are
989 constructed in the standard way. As for PRS based on topic-GWAS, for code s we extract its
990 probabilities in all topics (ϕ_{ts}), and calculate an individual’s PRS for it as: $PRS_s =$
991 $\sum_{t=1}^T (PRS_t \times \phi_{ts})$, where PRS_t is the individual’s PRS for topic t (constructed using the topic-
992 GWAS result for topic t). The area under the receiver-operator curve (AUC) is used to
993 evaluate the predictive performance of PRS on the test data.

994 **Analyses on the larger UKB dataset**

995 **The input data**

996 The larger UKB dataset (top-436 dataset) is constructed in the same way as the top-100 UKB
997 dataset, and contains the diagnostic records of the top 436 most frequent ICD-10 codes from
998 the first 14 chapters of the ICD-10 coding system for all individuals in UKB. These codes are
999 all the ones in UK Biobank with a prevalence of at least 0.001 at the date of selection
1000 (continued data collection means that prevalence will tend to increase over time),
1001 corresponding to approximately 500 cases. The prevalence threshold of 0.001 is chosen both
1002 for computational reasons (this is roughly the limit of what can be performed using available
1003 computing resources) and because there must be sufficient occurrences of diseases from
1004 which to discover multi-morbidity clusters. As with the top-100 dataset, we partition the full
1005 top-436 dataset into training (80%) and testing (20%) datasets. The top-436 and the top-100
1006 datasets use different partitions for the training and testing datasets.

1007 **Inference on the top-436 UKB dataset**

1008 Training strategy for the top-436 dataset

1009 The top-436 dataset is more than three times larger than the top-100 dataset, increasing the
1010 computational requirements for training topic models. On the top-100 dataset, treeLFA

1011 models with different numbers of topics are trained and compared. We find that when we set
1012 an excess number of topics for the model, both inferred topics and topic-GWAS results are
1013 stable across different models (Figure 5). Therefore, on the top-436 dataset, instead of
1014 training many models with different numbers of topics, we train treeLFA and flatLFA models
1015 with 100 topics, and cluster and collapse the inferred topics to combine all near-empty topics
1016 into a single one.

1017 For the optimization of α , the two-stage training strategy with the GibbsEM algorithm is used
1018 again. 1,500 iterations are run in the first stage (with a single posterior sample of Z collected
1019 in the E-step), and 350 iterations are run in the second stage (with 10 posterior samples of Z
1020 collected in the E-step). 50 posterior samples of hidden variables are collected during the last
1021 50 iterations for Gibbs-EM (with an interval of 200 iterations for the Gibbs sampling). For both
1022 treeLFA and flatLFA, three Gibbs chains are simulated.

1023 Choice of hyperparameters and initialization of hidden variables

1024 α is initialised as (1,0.1,...,0.1). Beta(0.1,3000) and Beta(1.2,3) are used as the prior for ϕ of
1025 inactive and active codes to account for diseases with small prevalence. The rest hidden
1026 variables and hyperparameters are set in the same way as for the top-100 dataset.

1027 Processing inference result

1028 We find that different Gibbs chains for treeLFA and flatLFA give slightly different inference
1029 results on the top-436 UKB dataset, while different posterior samples from the same chain
1030 have a very high level of consistency. Considering the variability among the inference results
1031 given by different chains, instead of clustering posterior samples of topics from all chains
1032 altogether, we cluster posterior samples from different chains separately. With the averaged
1033 ϕ and α for different chains, we calculate their predictive likelihood on the test data, and for
1034 both treeLFA and flatLFA we retain the chain which has the largest predictive likelihood, and
1035 use its inference result as the input for downstream analyses. For each topic inferred by the
1036 chain with the largest predictive likelihood, we check the inference results of the other chains,
1037 and annotate the topic with the number of chains that infer them to give a reference of its
1038 reliability.

1039 **Genetic analyses**

1040 topic-GWAS for the larger UKB dataset

1041 Since the top-436 dataset is much larger than the top100 dataset, to increase the inference
1042 accuracy for topic weights, after the training with Gibbs-EM algorithm we use Gibbs sampling
1043 to re-estimate individuals' topic weights, which is observed to increase the power of topic-
1044 GWAS. ϕ and α are fixed at values averaged from all posterior samples from the chain with
1045 the largest predictive likelihood. As a result, there is no longer an identifiability issue, so the
1046 results given by different chains (for the re-estimation of topic weights) can be combined
1047 directly. For both treeLFA and flatLFA, ten Gibbs chains are used to re-estimate topic
1048 weights, and 50 posterior samples are collected from each chain. Topic weights averaged
1049 from these chains are used as the input for topic-GWAS.

1050 Gene-set enrichment analysis for topic-associated SNPs

1051 The software FUMA can find genes that are close to the significant SNPs found by GWAS
1052 on the genome (the physical mapping function of FUMA). With the mapped genes, further
1053 analyses can be performed. Gene-set enrichment analysis (GSEA) tests for the enrichment
1054 of different gene sets among a group of genes. We choose genes that are associated with
1055 different traits in the GWAS catalogue as the reference gene sets to carry out GSEA for
1056 genes mapped from topic associated SNPs. By doing this, we can summarise the major
1057 associations of topic-associated SNPs found by previous GWAS. The default setting for
1058 FUMA is used in all the analyses in this section.

1059 Resource availability

1060 Lead Contact

1061 Further information and requests for resources should be directed to and will be fulfilled by
1062 the lead contact, Gil McVean (gil.mcvean@bdi.ox.ac.uk).

1063 Materials availability

1064 The key inference results (inferred topics) of the topic models and the topic-GWAS results
1065 are included in the supplementary material of the paper. The remaining results will be made
1066 available via the UK Biobank data return and linked to UK Biobank application
1067 number: 12788.

1068 Data and Code Availability

1069 This research has been conducted using the UK Biobank Resource: application number
1070 12788. The genotype data used for GWAS in this study comes from data field “22418” in UK
1071 Biobank. The diagnosis data used for topic modelling comes from data fields “41202” and
1072 “41204” in UK Biobank.

1073 The code for the treeLFA algorithm and a demo for using it on example data is available at:
1074 <https://github.com/zhangyd10/treeLFA-demo> or <https://doi.org/10.5281/zenodo.7420615>.

1075

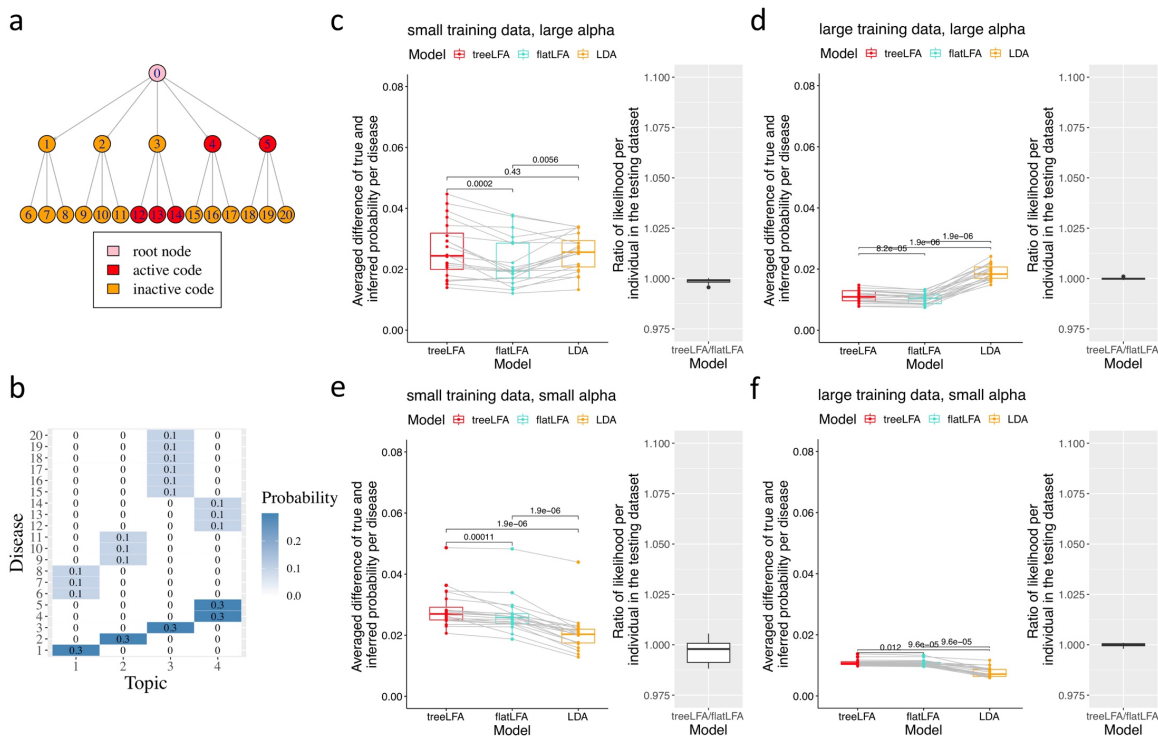
1076

1077 **Supplementary Information**

1078 **Description of Supplementary Data**

1079 Supplementary information includes 11 figures, 39 tables, and the analytical note.

1080 **Supplementary Figures**

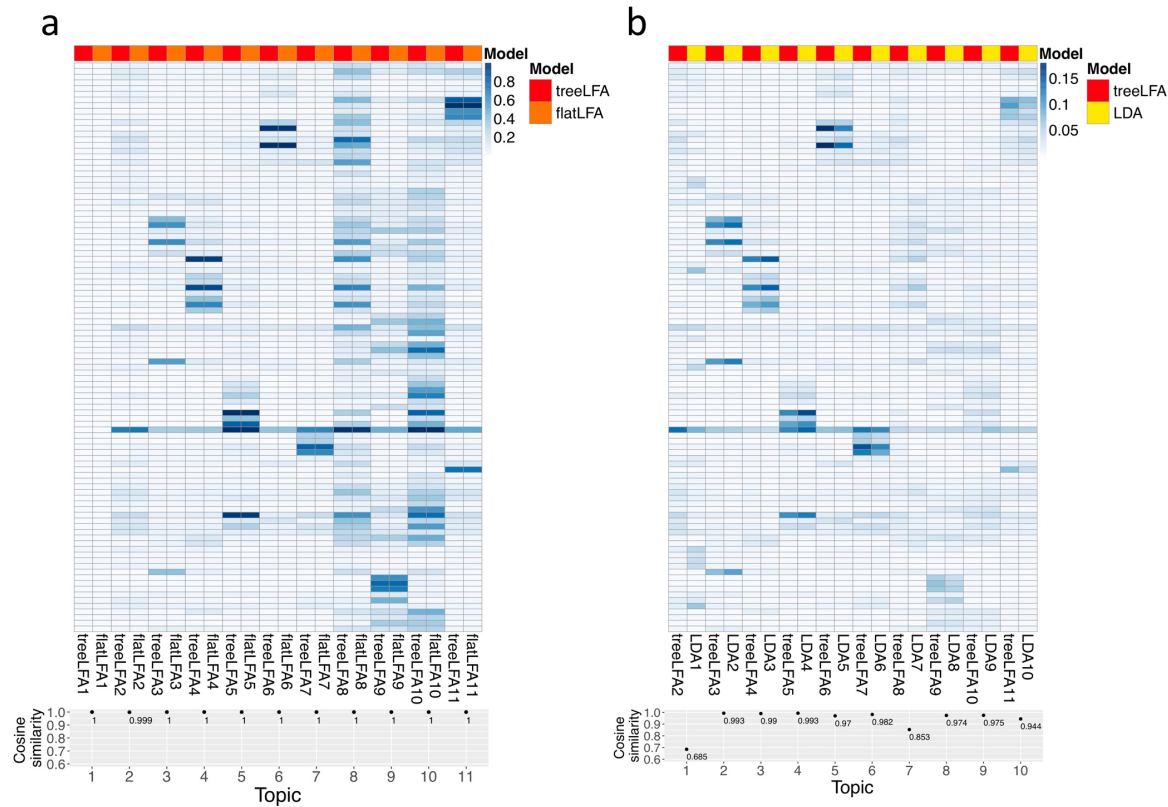


1081

1082 **Supplementary Figure 1 Comparison of three related topic models (treeLFA, flatLFA**
 1083 **and LDA) on simulated datasets.**

1084 a, The tree structure of 20 diseases. Red nodes correspond to the active codes in Topic 4 in panel b.
 1085 b, The four topics used for simulation. Active codes in these topics are unlikely to be generated by a
 1086 Markov process with small probability of transforming from inactive to active and large probability of
 1087 staying active while going from the parent node to its children nodes, since an active parent code always
 1088 has inactive children codes, while active children codes always have inactive parent code.
 1089 c-f Comparison of three topic models on simulated datasets. The Parameter setting and metrics are the
 1090 same as those in Figure 2. c, Results on datasets simulated using $D=2500$ and $\alpha=1$. d, Results on
 1091 datasets simulated using $D=5000$ and $\alpha=1$. e, Results on datasets simulated using $D=300$ and $\alpha=0.1$.
 1092 f, Results on datasets simulated using $D=1000$ and $\alpha=0.1$. The numeric results are in Supplementary
 1093 Table 23.

1094

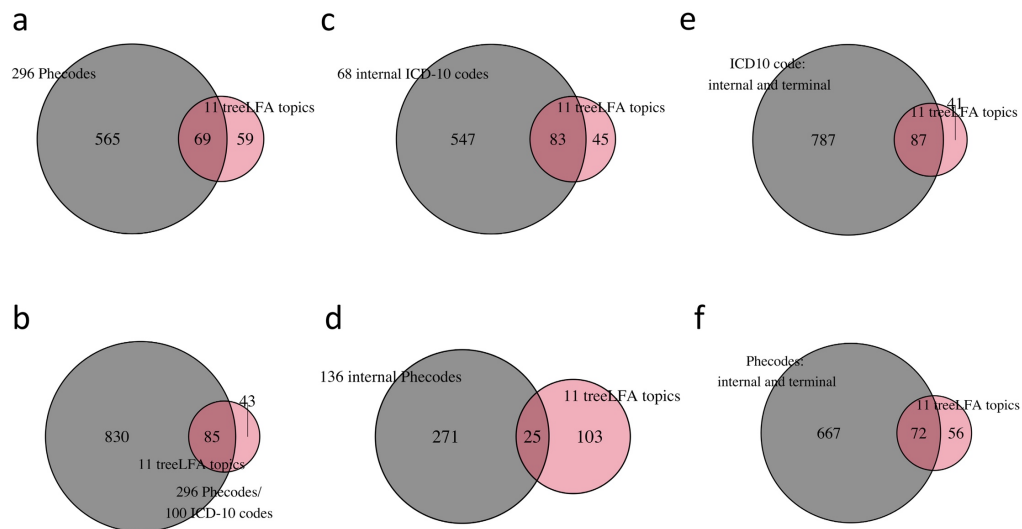


1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107

Supplementary Figure 2 Comparison of topics inferred by three topic models on the top-100 UKB dataset.

a, Comparison of the 11 topics inferred by treeLFA and flatLFA. The same topics inferred by the two models are placed next to each other. Cosine similarity was used to measure the similarity of topics inferred by the two models (point plot below the heatmap). The numeric results are in Supplementary Table 4,24.

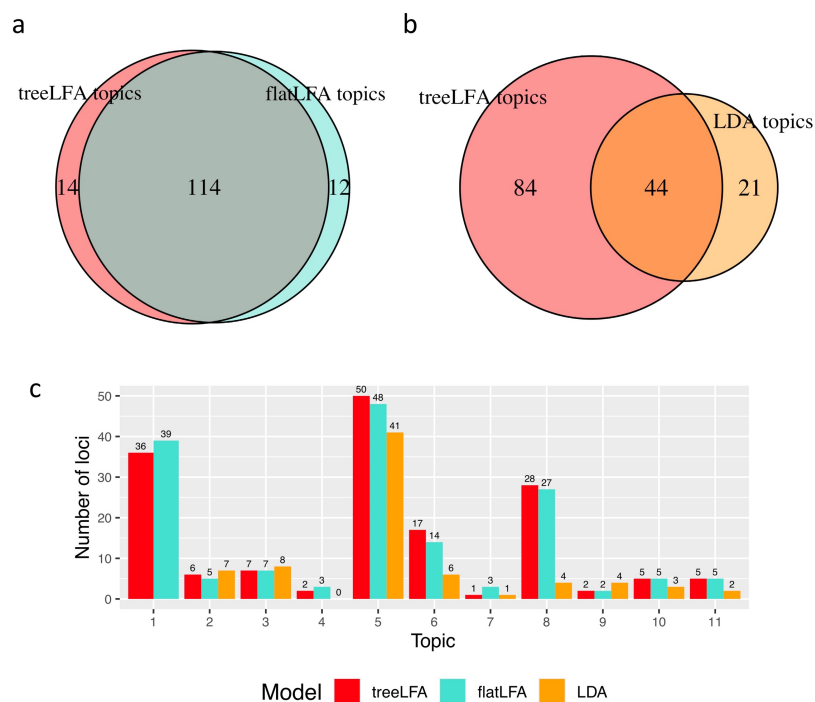
b, Comparison of the 10 topics inferred by LDA and the 10 non-empty topics inferred by treeLFA. Topics inferred by treeLFA are normalised such that probabilities of the 100 ICD-10 codes add up to 1 in any topic. The same topics inferred by the two models are placed next to each other. Cosine similarity was used to measure the similarity of topics inferred by the two models (point plot below the heatmap). The numeric results are in Supplementary Table 4,25.



1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122

Supplementary Figure 3 Overlap of significant loci found by GWAS for different traits.

- a, The total numbers of loci associated with any of the 296 Phecodes mapped from the 100 ICD-10 codes and any of the 11 treeLFA topics, and their overlap.
- b, The total numbers of loci associated with any of the 296 Phecodes or the 100 ICD-10 codes and any of the 11 treeLFA topics, and their overlap.
- c, The total numbers of loci associated with any of the 68 internal ICD-10 codes and any of the 11 treeLFA topics, and their overlap.
- d, The total numbers of loci associated with any of the 136 internal Phecodes and any of the 11 treeLFA topics, and their overlap.
- e, The total numbers of loci associated with any of the internal or terminal ICD-10 codes and any of the 11 treeLFA topics, and their overlap.
- f, The total numbers of loci associated with any of the internal or terminal Phecodes and any of the 11 treeLFA topics, and their overlap.



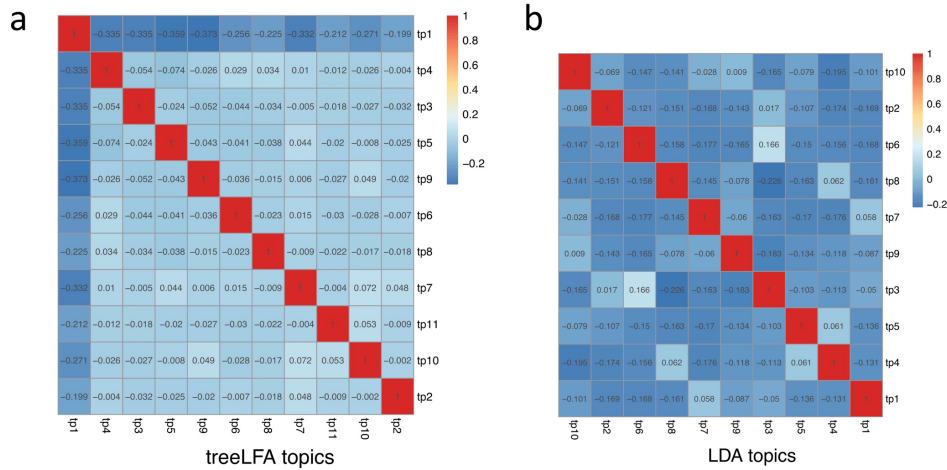
1123
 1124
 1125
 1126
 1127
 1128
 1129
 1130
 1131
 1132
 1133

Supplementary Figure 4 Comparison of the topic-GWAS results for three topic models.

a, The total numbers of loci associated with any of the 11 topics inferred by treeLFA and flatLFA, and their overlap.

b, The total numbers of loci associated with any of the 11 topics inferred by treeLFA and any of the 10 topics inferred by LDA, and their overlap.

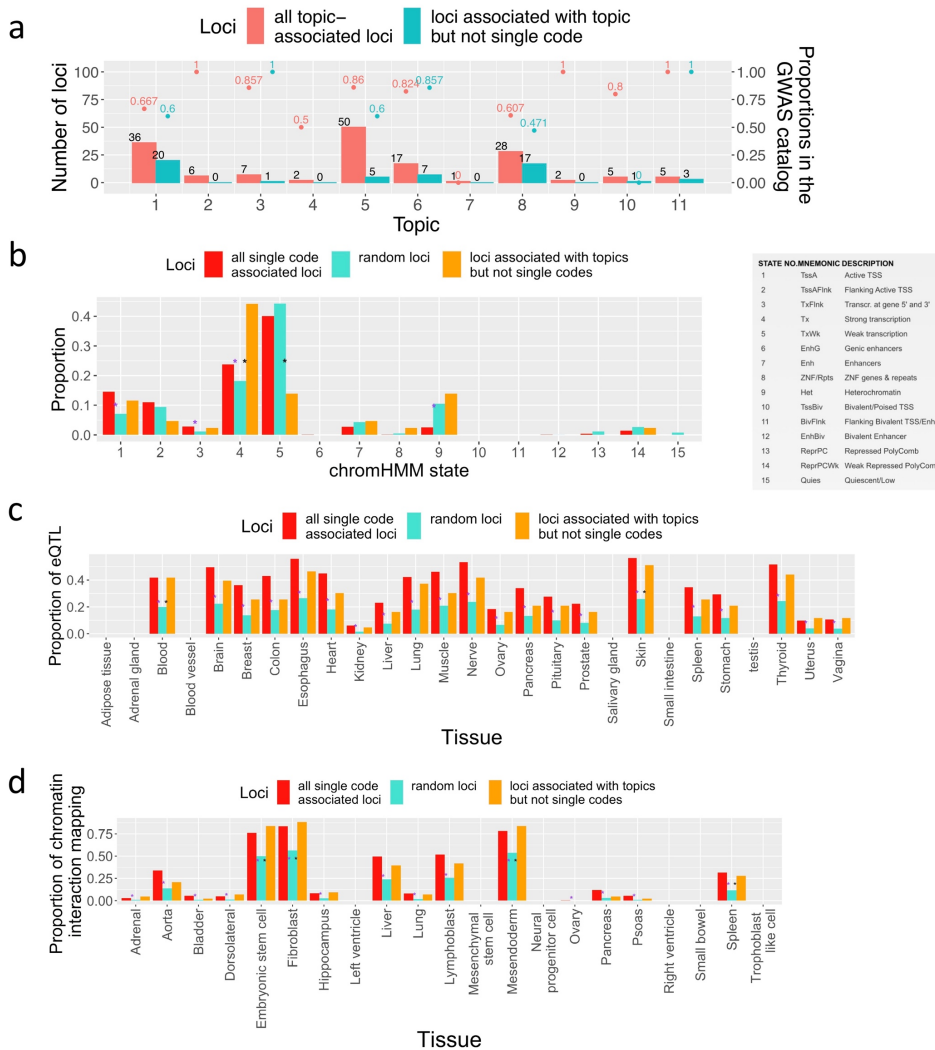
c, The numbers of loci associated with each of the 11 topics inferred by treeLFA and flatLFA, and each of the 10 topics inferred by LDA. The first topic is the empty topic, and is only inferred by treeLFA and flatLFA.



1134
1135
1136
1137
1138
1139
1140
1141

Supplementary Figure 5 Correlation of topic weights for topics inferred by treeLFA and LDA.

a, The correlation matrix for individuals' weights for the 11 topics inferred by treeLFA. The first topic is the empty topic.
b, The correlation matrix for individuals' weights for the 10 topics inferred by LDA. The matrix uses the same colour scheme as the matrix in panel a.



1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162

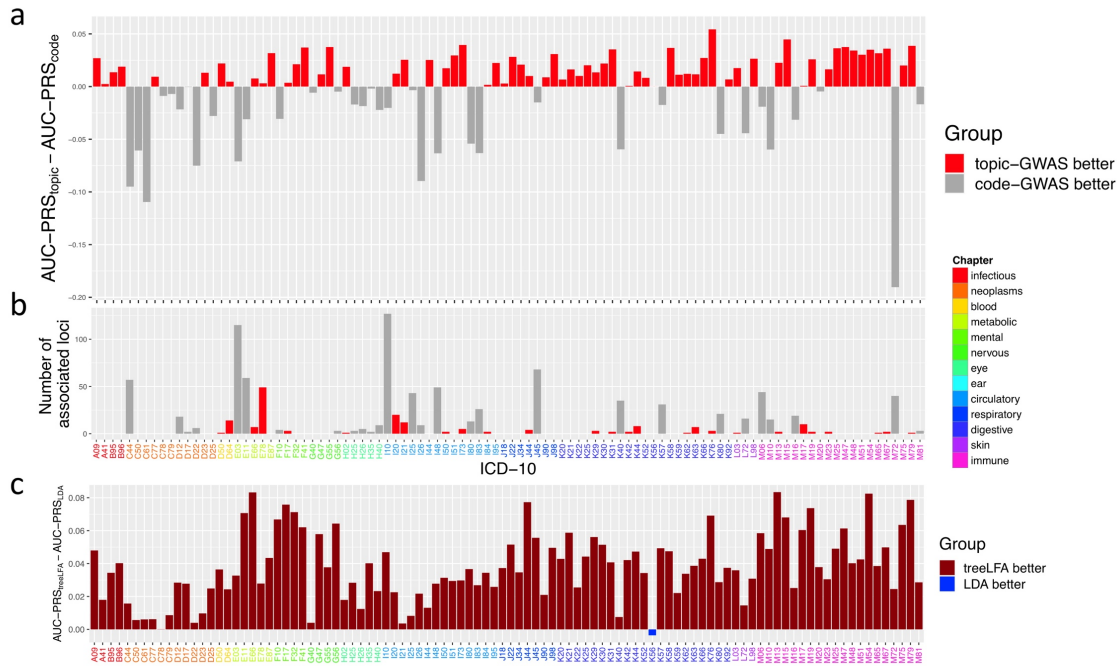
Supplementary Figure 6 Validation of topic associated loci.

a, The proportions of topic-associated loci recorded in the GWAS catalogue. Red bars show the results for all topic-associated loci, and green bars show the results for loci that are associated with topics but not any single code.

b, Proportions of three groups of SNPs that are in different chromHMM states. Meanings of different chromHMM states are shown in the right table. The first group contains lead SNPs associated with at least one ICD-10 code; the second group contains 10,000 random SNPs whose allele frequency is matched with that of topic-associated lead SNPs; The third group contains lead SNPs associated with at least one of the 11 topics but not any single code. The proportions of the first and third groups are compared with the second group respectively, and significant differences in proportions (two-proportion Z-test, adjusted P-value<0.05, Bonferroni correction) are marked with asterisks between the corresponding bars (purple asterisks between red and green bars mean significant differences between the first and second groups, black asterisks between blue and green bars mean significant differences between the third and second groups). The numeric results are in Supplementary Table 26.

c, Proportions of the three groups of SNPs in panel b that are eQTL in different tissues. The numeric results are in Supplementary Table 27.

d, Proportions of the three groups of SNPs in panel b that have chromatin interaction with genes in different tissues. The numeric results are in Supplementary Table 28.



1163

1164

Supplementary Figure 7 PRS for ICD-10 codes based on topic-GWAS results.

1165 a, Comparison of the AUC of two types of PRS for the 100 ICD-10 codes on the test data. One type of

1166 PRS is constructed using the topic-GWAS result for treeLFA, and the other type of PRS is constructed

1167 using single code GWAS result. Bars are colored according to the relative performance of the two

1168 types of PRS. The numeric results are in Supplementary Table 10.

1169 b, The numbers of loci associated with each of the 100 ICD-10 codes. Bars are colored the same way

1170 as in panel a. Codes are colored according to the ICD-10 chapters they belong to. The numeric

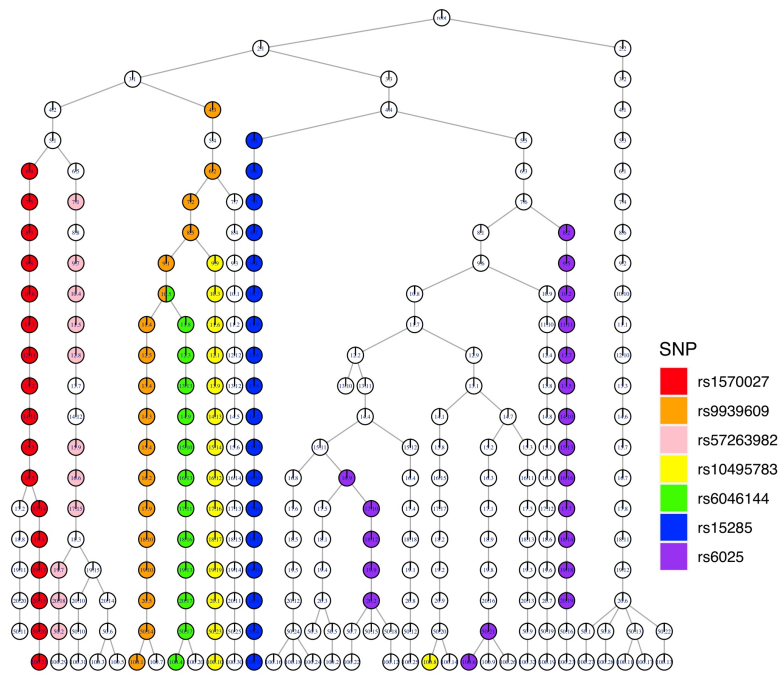
1171 results are in Supplementary Table 29.

1172 c, Comparison of the AUC of PRS constructed using the topic-GWAS results for treeLFA and LDA.

1173 Bars are colored according to the relative performance of the two types of PRS. The numeric results

1174 are in Supplementary Table 30.

1175



1176

1177

1178

Supplementary Figure 8 Association of SNPs and topics from different treeLFA models.

1179

The associations of a few example SNPs and topics inferred by different treeLFA models are

1180

visualised on the tree structure of topics. The tree structure of topics is the same as the one in Figure

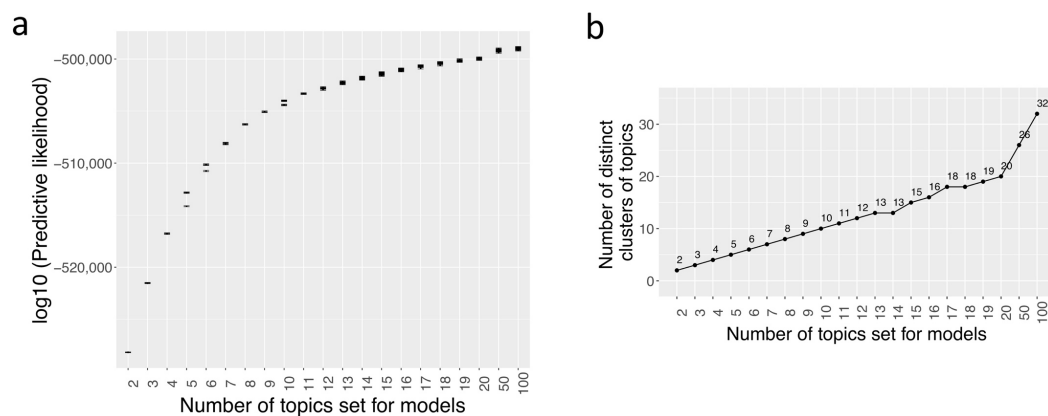
1181

5b. Topics significantly associated with 7 different SNPs are highlighted with different colours on the

1182

tree structure. The numeric results are in Supplementary Table 31.

1183

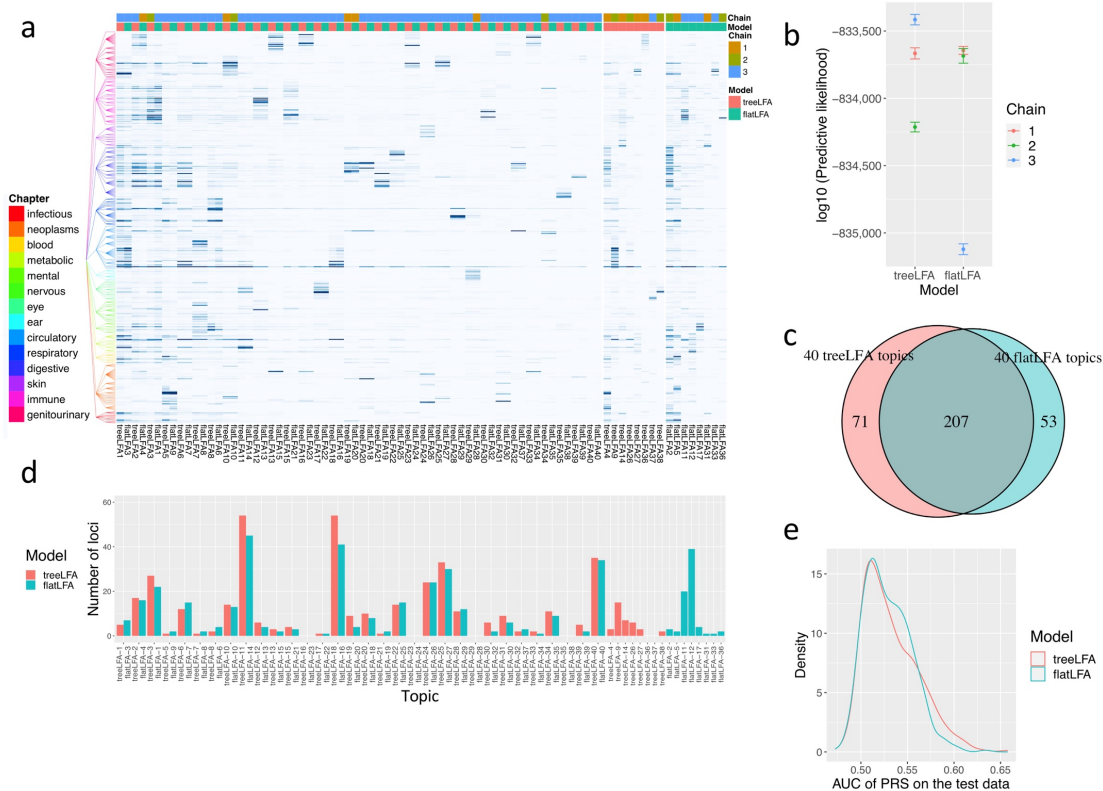


1184
1185
1186
1187
1188
1189
1190
1191
1192

Supplementary Figure 9 Summary statistics for models with different numbers of topics.

a, The predictive likelihood on the test data for the ten Gibbs chains for treeLFA models set with different numbers of topics. For each chain, the standard deviation of the likelihood calculated using different posterior samples of topics were shown. The numeric results are in Supplementary Table 32.

b, Numbers of distinct topics remained after clustering for treeLFA models set with different numbers of topics.



1193
1194
1195
1196
1197
1198
1199
1200
1201
1202
1203
1204
1205
1206
1207
1208
1209
1210
1211
1212

Supplementary figure 10 Comparison of the inference and topic-GWAS results for treeLFA and flatLFA on the top-436 UKB dataset.

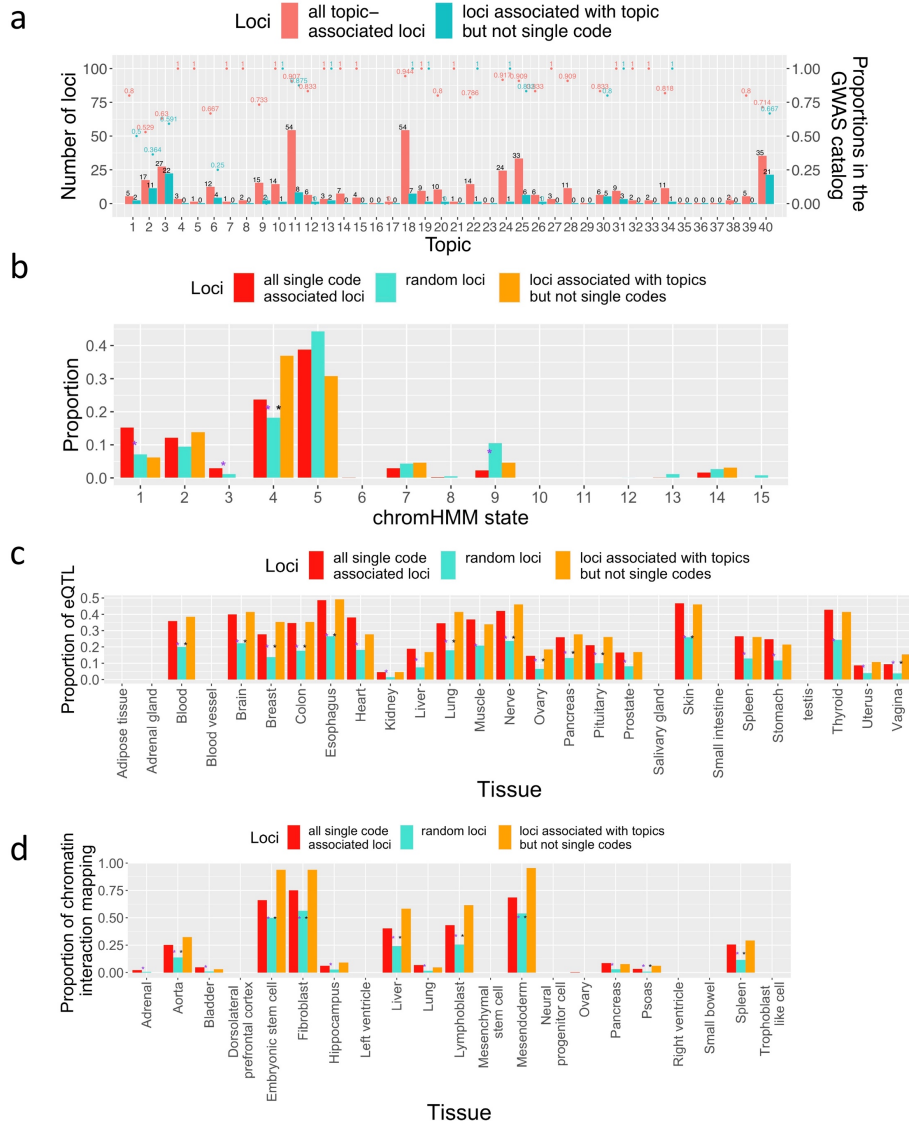
a, 40 topics inferred by treeLFA and flatLFA models set with 100 topics. The same topics inferred by treeLFA and flatLFA are placed next to each other. Topics inferred by both models are shown first, followed by topics inferred by only one model. For each model, the inferred topics are numbered according to their density. The tree structure of the 436 ICD-10 codes is shown to the left of the heatmap, and codes from different ICD-10 chapters are colored differently. For each topic, the number of chains that inferred it is shown with the colour bar on top of the heatmap. The numeric results are in Supplementary Table 33.

b, The predictive likelihood on the test data for the three treeLFA and flatLFA chains. The calculation of predictive likelihood was repeated ten times for each chain to get the standard deviation. The numeric results are in Supplementary Table 34.

c, The total numbers of loci associated with any of the topics inferred by treeLFA and flatLFA, and their overlap.

d, The numbers of loci associated with each of the treeLFA and flatLFA topics. Topics have the same order as those in panel a. The numeric results are in Supplementary Table 35.

e, Density plots for the AUC of PRS for the 436 ICD-10 codes on the test data based on the topic-GWAS results for treeLFA and flatLFA. The numeric results are in Supplementary Table 36.



1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224

Supplementary figure 11 Validation of topic associated loci for the top-436 UKB dataset.

All settings are the same as those in Supplementary Figure 6.

a, The proportions of topic-associated loci recorded in the GWAS catalogue.

b, Proportions of the three groups of SNPs that are in different chromHMM states. The numeric results are in Supplementary Table 37.

c, Proportions of the three groups of SNPs in panel b that are eQTL in different tissues. The numeric results are in Supplementary Table 38.

d, Proportions of the three groups of SNPs in panel b that have chromatin interaction with genes in different tissues. The numeric results are in Supplementary Table 39.

1225

1226 Supplementary Tables

1227 **Supplementary Table 1. Parameter setting for all simulated datasets.** 8 groups of
1228 datasets are simulated using different combinations of three parameters: the concentration
1229 parameter (α) of the Dirichlet prior for topic weight variable θ , the number of individuals in
1230 the training dataset (D), and the tree structure of diseases. Correct tree structure means that
1231 the topics used for simulation are likely to be constructed by running a Markov process with
1232 small ρ_{01} and large ρ_{11} (see Methods) on the tree structure.

1233 **Supplementary Table 2. Comparison of 3 topic models on the 4 groups of datasets**
1234 **simulated using correct tree structure of diseases.** Each group contains 20 datasets
1235 simulated using the same topics and hyperparameters (Supplementary Table 1).
1236 Phi_diff_ave: the averaged per disease difference in probability (Phi) between true and
1237 inferred topics. I_diff_ave: the averaged per disease difference in the values of indicator
1238 variables (I) between true and inferred topics. Predictive likelihood: predictive likelihood on
1239 the corresponding test data, which is calculated using the topics inferred from the training
1240 data.

1241 **Supplementary Table 3. Tree structure of the top-100 most frequent ICD-10 codes from**
1242 **chapters 1-13 of the ICD-10 coding system in UK Biobank.** In the first column both internal
1243 and terminal nodes of the tree structure of the top-100 most frequent ICD-10 codes are
1244 listed (ordered according to the layers of nodes on the tree and the names of nodes). In the
1245 second column the parent nodes of the nodes in the first column are shown.

1246 **Supplementary Table 4. The 11 topics inferred by treeLFA on the top-100 UKB dataset.**
1247 The probability of the 100 ICD-10 codes in the 11 topics inferred by the treeLFA model set
1248 with 11 topics on the top-100 UKB dataset is shown.

1249 **Supplementary Table 5. Topic weights for the 11 topics inferred by treeLFA for 2,000**
1250 **randomly selected individuals in UKB.**

1251 **Supplementary Table 6. Number of loci associated with individuals' weights for the 11**
1252 **topics inferred by treeLFA on the top-100 UKB dataset.**

1253 **Supplementary Table 7. Comparison of $-\log_{10}$ (P-value) for lead SNPs of Topic 5 given**
1254 **by single code GWAS for the top five active codes in Topic 5 and topic-GWAS for**
1255 **Topic 5.**

1256 **Supplementary Table 8. Inflation in P-values given by topic-GWAS for the 11 topics**
1257 **inferred by treeLFA on the top-100 UKB dataset.** Both the intercept and λ_{GC} are given by
1258 LDSC software (see Methods). λ_{GC} : genomic control inflation factor.

1259 **Supplementary Table 9. Performance of PRS for topics constructed using topic-**
1260 **GWAS results on the test data.** R2: phenotypic variance of topic weights explained by the
1261 PRS for topics. P: P-value of model fit on the testing data. Heritability: heritabilities of topic
1262 weights calculated using the LD score regression and summary statistics of topic-GWAS.
1263 NUM_SNP: numbers of SNPs used by the software "PRSice-2" to construct the PRS for
1264 topics. Threshold: threshold for P-values of SNPs used by "PRSice-2" to construct the PRS
1265 for topics.

- 1266 **Supplementary Table 10. Comparison of the performance of two types of PRS for the**
1267 **100 ICD10 codes.** One type of PRS for ICD-10 codes is constructed using topic-GWAS
1268 results (see Methods), and the other type of PRS is constructed using single code GWAS
1269 results. AUC: the area under the receiver-operator curve of the PRS on the test data.
- 1270 **Supplementary Table 11. Topics inferred by treeLFA models set with different**
1271 **numbers of topics on the top-100 UKB dataset.** The probability of the 100 ICD-10 codes
1272 in topics inferred by treeLFA models set with different numbers of topics on the top-100 UKB
1273 dataset is shown.
- 1274 **Supplementary Table 12. topic-GWAS results for the example SNP "rs143384" and all**
1275 **topics inferred by all treeLFA models on the top-100 UKB dataset.** Topics are named
1276 using the number of topics set for the corresponding treeLFA model, and the index of the
1277 inferred topic for that model. For instance, topic "10.10" means the 10th topic inferred by the
1278 treeLFA model set with 10 topics. P: the P-value of the SNP given by the linear regression
1279 for the topic weight. BETA: estimated effect size of the SNP. SE: standard error of the
1280 estimated effect size.
- 1281 **Supplementary Table 13. Tree structure of the top-436 most frequent ICD-10 codes**
1282 **from chapters 1-14 of the ICD-10 coding system in UKB.** In the first column both internal
1283 and terminal nodes of the tree structure of the top-436 most frequent ICD-10 codes are
1284 listed (ordered according to the layers of nodes on the tree and the names of nodes). In the
1285 second column the parent nodes of the nodes in the first column are shown.
- 1286 **Supplementary Table 14. The 40 topics inferred by treeLFA on the top-436 UKB**
1287 **dataset.** The probability of the 436 ICD-10 codes in the 40 topics inferred by treeLFA on the
1288 top-436 UKB dataset is shown (see Methods).
- 1289 **Supplementary Table 15. Numbers of active codes from different ICD10 chapters in**
1290 **the 40 topics inferred by treeLFA on the top-436 UKB dataset (see Methods).**
- 1291 **Supplementary Table 16. Significant SNPs found by topic-GWAS for the 40 treeLFA**
1292 **inferred topics.** CHR: chromosome. BP: physical position (base pair) of SNP. BETA:
1293 estimated effect size (regression coefficient). SE: standard error of the estimated effect size.
1294 Topic: the index of the inferred topic for the corresponding treeLFA model. Model: the
1295 number of topics set for the treeLFA model.
- 1296 **Supplementary Table 17. Effect sizes of lead SNPs associated with the 40 treeLFA**
1297 **inferred topics.**
- 1298 **Supplementary Table 18. Numbers of associated loci for the 40 topics inferred by**
1299 **treeLFA on the top-436 UKB dataset.**
- 1300 **Supplementary Table 19. treeLFA inferred topics with a substantial number of novel**
1301 **loci found by topic-GWAS.** Dominant ICD10 chapter: the ICD-10 chapter which contributes
1302 the majority of active codes in the corresponding topic.
- 1303 **Supplementary Table 20. Comparison of the performance of two types of PRS for the**
1304 **436 ICD10 codes.** One type of PRS for ICD-10 codes is constructed using topic-GWAS
1305 results (see Methods), and the other type of PRS is constructed using single code GWAS
1306 results. AUC: the area under the receiver-operator curve of the PRS on the test data.
- 1307 **Supplementary Table 21. Enriched gene sets in the GWAS catalogue among genes**
1308 **associated with the 40 treeLFA topics.** GeneSet: names of gene sets in the GWAS

1309 catalogue. N: the total numbers of genes in the gene sets. n: the numbers of genes in the
1310 gene sets that are associated with the treeLFA topics. P-value: P-values of the enrichment
1311 analyses. adjusted-P: adjusted P-values of the enrichment analyses. Genes: genes in the
1312 enriched gene sets that are associated with the topics. This table is the output of the “SNP-
1313 to-gene” function of the software “FUMA”

1314 **Supplementary Table 22. The correspondence between treeLFA inferred topics and**
1315 **loci level genetically interpretable multimorbidity networks found by the study “A**
1316 **global overview of genetically interpretable multimorbidities among common**
1317 **diseases in the UK Biobank”.** The information of the 9 loci level genetically interpretable
1318 multimorbidity networks found by a previous study using the diagnostic data in UKB is
1319 shown, together with their corresponding topics inferred by treeLFA on the top-436 UKB
1320 dataset. The first 5 columns come from the supplementary results of the paper "A global
1321 overview of genetically interpretable multimorbidities among common diseases in the UK
1322 Biobank".

1323 **Supplementary Table 23. Comparison of 3 topic models on the 4 groups of datasets**
1324 **simulated using wrong tree structure of diseases.** Each group contains 20 datasets
1325 simulated using the same topics and hyperparameters (Supplementary Table 1).
1326 Phi_diff_ave: the averaged per disease difference in probability (Phi) between true and
1327 inferred topics. I_diff_ave: the averaged per disease difference in the values of indicator
1328 variables (I) between true and inferred topics. Predictive likelihood: predictive likelihood on
1329 the corresponding test data, which is calculated using the topics inferred from the training
1330 data.

1331 **Supplementary Table 24. The 11 topics inferred by flatLFA on the top100 UKB dataset.**
1332 The probability of the 100 ICD-10 codes in the 11 topics inferred by the flatLFA model set
1333 with 11 topics on the top-100 UKB dataset is shown.

1334 **Supplementary Table 25. The 10 topics inferred by LDA on the top100 UKB dataset.**
1335 The probability of the 100 ICD-10 codes in the 10 topics inferred by the LDA model set with
1336 10 topics on the top-100 UKB dataset is shown.

1337 **Supplementary Table 26. Proportions of SNPs in different chromHMM states for the 3**
1338 **groups of SNPs on the top-100 UKB dataset.** P(random/single code): P-values for the
1339 comparison of proportions for random/single code-associated SNPs. P(random/topic): P-
1340 values for the comparison of proportions for random/topic-associated SNPs.

1341 **Supplementary Table 27. Proportions of SNPs that are eQTL in different tissues for**
1342 **the 3 groups of SNPs on the top-100 UKB dataset.** P(random/single code): P-values for
1343 the comparison of proportions for random/single code-associated SNPs. P(random/topic): P-
1344 values for the comparison of proportions for random/topic-associated SNPs.

1345 **Supplementary Table 28. Proportions of SNPs with chromatin mapping in different**
1346 **tissues for the 3 groups of SNPs on the top-100 UKB dataset.** P(random/single code): P-
1347 values for the comparison of proportions for random/single code-associated SNPs.
1348 P(random/topic): P-values for the comparison of proportions for random/topic-associated
1349 SNPs.

1350 **Supplementary Table 29. Numbers of significant loci given by single code GWAS for**
1351 **the 100 ICD10 codes.**

1352 **Supplementary Table 30. Comparison of the performance of PRS for the 100 ICD-10**

- 1353 **codes based on topic-GWAS results for treeLFA/LDA inferred topics on the top-100**
1354 **UKB dataset.** AUC: the area under the receiver-operator curve of the PRS on the test data.
- 1355 **Supplementary Table 31. topic-GWAS results for 7 example SNPs and all topics**
1356 **inferred by treeLFA models set with different numbers of topics on the top-100 UKB**
1357 **dataset.** P-values given by topic-GWAS for the examples SNPs are shown. Topics (in the
1358 first column) are named using the number of topics set for the corresponding treeLFA model,
1359 and the index of the inferred topic for that model. For instance, topic “10.10” means the 10th
1360 topic inferred by the treeLFA model set with 10 topics.
- 1361 **Supplementary Table 32. Predictive likelihood for treeLFA models set with different**
1362 **numbers of topics on the top-100 UKB dataset.** For each treeLFA model, 10 Gibbs
1363 chains are trained, and 50 posterior samples of topics are collected from each chain. The
1364 predictive likelihood on the test data is calculated using each posterior sample of topics, and
1365 standard error of predictive likelihood is calculated for each chain.
- 1366 **Supplementary Table 33. The 40 topics inferred by treeLFA and flatLFA on the top-436**
1367 **UKB dataset.** The probability of the 436 ICD-10 codes in the 40 topics inferred by treeLFA
1368 and flatLFA is shown.
- 1369 **Supplementary Table 34. Predictive likelihood on the test data for three treeLFA and**
1370 **flatLFA chains on the top-436 dataset.** For each chain, the predictive likelihood is
1371 calculated 10 times.
- 1372 **Supplementary Table 35. Numbers of associated loci for the 40 treeLFA and flatLFA**
1373 **topics on the top-436 UKB dataset.** Topics are named according to their order in
1374 Supplementary Table 33.
- 1375 **Supplementary Table 36. Comparison of the performance of PRS for the 436 ICD10**
1376 **codes based on topic-GWAS results for treeLFA and flatLFA topics.** AUC: the area
1377 under the receiver-operator curve of the PRS on the test data.
- 1378 **Supplementary Table 37. Proportions of SNPs in different chromHMM states for the 3**
1379 **groups of SNPs on the top-436 UKB dataset.** P(random/single code): P-values for the
1380 comparison of proportions for random/single code-associated SNPs. P(random/topic): P-
1381 values for the comparison of proportions for random/topic-associated SNPs.
- 1382 **Supplementary Table 38. Proportions of SNPs that are eQTL in different tissues for**
1383 **the 3 groups of SNPs on the top-436 UKB dataset.** P(random/single code): P-values for
1384 the comparison of proportions for random/single code-associated SNPs. P(random/topic): P-
1385 values for the comparison of proportions for random/topic-associated SNPs.
- 1386 **Supplementary Table 39. Proportions of SNPs with chromatin mapping in different**
1387 **tissues for the 3 groups of SNPs on the top-436 UKB dataset.** P(random/single code): P-
1388 values for the comparison of proportions for random/single code-associated SNPs.
1389 P(random/topic): P-values for the comparison of proportions for random/topic-associated
1390 SNPs.

1391

1392

1393 References

- 1394 1. Garin, N., Koyanagi, A., Chatterji, S., Tyrovolas, S., Olaya, B., Leonardi, M. et al. Global
1395 Multimorbidity Patterns: A Cross-Sectional, Population-Based, Multi-Country Study. *J.*
1396 *Gerontol. A Biol. Sci. Med. Sci.* **71**, 205–214 (2016).
- 1397 2. Fortin, M., Stewart, M., Poitras, M.-E., Almirall, J. & Maddocks, H. A Systematic Review
1398 of Prevalence Studies on Multimorbidity: Toward a More Uniform Methodology. *The*
1399 *Annals of Family Medicine* vol. 10 142–151 (2012).
- 1400 3. Violan, C., Foguet-Boreu, Q., Flores-Mateo, G., Salisbury, C., Blom, J., Freitag, M. et al.
1401 Prevalence, determinants and patterns of multimorbidity in primary care: a systematic
1402 review of observational studies. *PLoS One* **9**, e102149 (2014).
- 1403 4. Ryan, A., Wallace, E., O'Hara, P. & Smith, S. M. Multimorbidity and functional decline in
1404 community-dwelling adults: a systematic review. *Health Qual. Life Outcomes* **13**, 168
1405 (2015).
- 1406 5. Mair, F. S. & May, C. R. Thinking about the burden of treatment. *BMJ* **349**, (2014).
- 1407 6. Van Wilder, L., Devleeschauwer, B., Clays, E., De Buyser, S., Van der Heyden, J.,
1408 Charafeddine, R. et al. The impact of multimorbidity patterns on health-related quality of
1409 life in the general population: results of the Belgian Health Interview Survey. *Qual. Life*
1410 *Res.* **31**, 551–565 (2022).
- 1411 7. Prados-Torres, A., Calderón-Larrañaga, A., Hanco-Saavedra, J., Poblador-Plou, B. &
1412 van den Akker, M. Multimorbidity patterns: a systematic review. *J. Clin. Epidemiol.* **67**,
1413 254–266 (2014).
- 1414 8. Holden, L., Scuffham, P. A., Hilton, M. F., Muspratt, A., Ng, S.-K. & Whiteford, H. A.
1415 Patterns of multimorbidity in working Australians. *Popul. Health Metr.* **9**, 15 (2011).
- 1416 9. Ng, S. K., Holden, L. & Sun, J. Identifying comorbidity patterns of health conditions via
1417 cluster analysis of pairwise concordance statistics. *Stat. Med.* **31**, 3393–3405 (2012).
- 1418 10. Guisado-Clavero, M., Roso-Llorach, A., López-Jimenez, T., Pons-Vigués, M., Foguet-

- 1419 Boreu, Q., Muñoz, M. A. et al. Multimorbidity patterns in the elderly: a prospective cohort
1420 study with cluster analysis. *BMC Geriatr.* **18**, 16 (2018).
- 1421 11. Dong, G., Feng, J., Sun, F., Chen, J. & Zhao, X.-M. A global overview of genetically
1422 interpretable multimorbidities among common diseases in the UK Biobank. *Genome*
1423 *Med.* **13**, 110 (2021).
- 1424 12. Shang, X., Zhang, X., Huang, Y., Zhu, Z., Zhang, X., Liu, J. et al. Association of a wide
1425 range of individual chronic diseases and their multimorbidity with brain volumes in the
1426 UK Biobank: A cross-sectional study. *EClinicalMedicine* **47**, 101413 (2022).
- 1427 13. Ronaldson, A., Arias de la Torre, J., Bendayan, R., Yadegarfar, M. E., Rhead, R.,
1428 Douiri, A. et al. Physical multimorbidity, depressive symptoms, and social participation in
1429 adults over 50 years of age: findings from the English Longitudinal Study of Ageing.
1430 *Aging Ment. Health* 1–11 (2022) doi:10.1080/13607863.2021.2017847.
- 1431 14. Schramm, S., Møller, S. P., Tolstrup, J. S. & Laursen, B. Effects of individual and
1432 parental educational levels on multimorbidity classes: a register-based longitudinal
1433 study in a Danish population. *BMJ Open* vol. 12 e053274 (2022).
- 1434 15. Rajoo, S. S., Wee, Z. J., Lee, P. S. S., Wong, F. Y. & Lee, E. S. A Systematic Review of
1435 the Patterns of Associative Multimorbidity in Asia. *Biomed Res. Int.* **2021**, 6621785
1436 (2021).
- 1437 16. Zemedikun, D. T., Gray, L. J., Khunti, K., Davies, M. J. & Dhalwani, N. N. Patterns of
1438 Multimorbidity in Middle-Aged and Older Adults: An Analysis of the UK Biobank Data.
1439 *Mayo Clin. Proc.* **93**, 857–866 (2018).
- 1440 17. Ronaldson, A., Arias de la Torre, J., Prina, M., Armstrong, D., Das-Munshi, J., Hatch, S.
1441 et al. Associations between physical multimorbidity patterns and common mental health
1442 disorders in middle-aged adults: A prospective analysis using data from the UK
1443 Biobank. *Lancet Reg Health Eur* **8**, 100149 (2021).
- 1444 18. Bycroft, C., Freeman, C., Petkova, D., Band, G., Elliott, L. T., Sharp, K. et al. The UK
1445 Biobank resource with deep phenotyping and genomic data. *Nature* **562**, 203–209

- 1446 (2018).
- 1447 19. Chen, Z., Chen, J., Collins, R., Guo, Y., Peto, R., Wu, F. et al. China Kadoorie Biobank
1448 of 0.5 million people: survey methods, baseline characteristics and long-term follow-up.
1449 *International Journal of Epidemiology* vol. 40 1652–1666 (2011).
- 1450 20. Kanai, M., Akiyama, M., Takahashi, A., Matoba, N., Momozawa, Y., Ikeda, M. et al.
1451 Genetic analysis of quantitative traits in the Japanese population links cell types to
1452 complex human diseases. *Nat. Genet.* **50**, 390–400 (2018).
- 1453 21. Cotsapas, C., Voight, B. F., Rossin, E., Lage, K., Neale, B. M., Wallace, C. et al.
1454 Pervasive sharing of genetic effects in autoimmune disease. *PLoS Genet.* **7**, e1002254
1455 (2011).
- 1456 22. Bhattacharjee, S., Rajaraman, P., Jacobs, K. B., Wheeler, W. A., Melin, B. S., Hartge,
1457 P. et al. A subset-based approach improves power and interpretation for the combined
1458 analysis of genetic association studies of heterogeneous traits. *Am. J. Hum. Genet.* **90**,
1459 821–835 (2012).
- 1460 23. Zhu, X., Feng, T., Tayo, B. O., Liang, J., Young, J. H., Franceschini, N. et al. Meta-
1461 analysis of correlated traits via summary statistics from GWASs with an application in
1462 hypertension. *Am. J. Hum. Genet.* **96**, 21–36 (2015).
- 1463 24. van der Sluis, S., Posthuma, D. & Dolan, C. V. TATES: efficient multivariate genotype-
1464 phenotype analysis for genome-wide association studies. *PLoS Genet.* **9**, e1003235
1465 (2013).
- 1466 25. Trochet, H., Pirinen, M., Band, G., Jostins, L., McVean, G. & Spencer, C. C. A.
1467 Bayesian meta-analysis across genome-wide association studies of diverse
1468 phenotypes. *Genet. Epidemiol.* **43**, 532–547 (2019).
- 1469 26. Majumdar, A., Haldar, T., Bhattacharya, S. & Witte, J. S. An efficient Bayesian meta-
1470 analysis approach for studying cross-phenotype genetic associations. *PLoS Genet.* **14**,
1471 e1007139 (2018).
- 1472 27. Turley, P., Walters, R. K., Maghziyan, O., Okbay, A., Lee, J. J., Fontana, M. A. et al.

- 1473 Multi-trait analysis of genome-wide association summary statistics using MTAG. *Nat.*
1474 *Genet.* **50**, 229–237 (2018).
- 1475 28. O'Reilly, P. F., Hoggart, C. J., Pomyen, Y., Calboli, F. C. F., Elliott, P., Jarvelin, M.-R. et
1476 al. MultiPhen: joint model of multiple phenotypes can increase discovery in GWAS.
1477 *PLoS One* **7**, e34861 (2012).
- 1478 29. Jostins, L. & McVean, G. Trinculo: Bayesian and frequentist multinomial logistic
1479 regression for genome-wide association studies of multi-category phenotypes.
1480 *Bioinformatics* **32**, 1898–1900 (2016).
- 1481 30. Hartley, S. W. & Sebastiani, P. PleioGRiP: genetic risk prediction with pleiotropy.
1482 *Bioinformatics* **29**, 1086–1088 (2013).
- 1483 31. Stephens, M. A unified framework for association analysis with multiple related
1484 phenotypes. *PLoS One* **8**, e65245 (2013).
- 1485 32. Joo, J. W. J., Kang, E. Y., Org, E., Furlotte, N., Parks, B., Hormozdiari, F. et al. Efficient
1486 and Accurate Multiple-Phenotype Regression Method for High Dimensional Data
1487 Considering Population Structure. *Genetics* **204**, 1379–1390 (2016).
- 1488 33. Liu, J., Pei, Y., Papasian, C. J. & Deng, H.-W. Bivariate association analyses for the
1489 mixture of continuous and binary traits with the use of extended generalized estimating
1490 equations. *Genet. Epidemiol.* **33**, 217–227 (2009).
- 1491 34. Liu, L., Tang, L., Dong, W., Yao, S. & Zhou, W. An overview of topic modeling and its
1492 current applications in bioinformatics. *Springerplus* **5**, 1608 (2016).
- 1493 35. Lee, M., Liu, Z., Kelly, R. & Tong, W. Of text and gene--using text mining methods to
1494 uncover hidden knowledge in toxicogenomics. *BMC Syst. Biol.* **8**, 93 (2014).
- 1495 36. Bicego, M., Lovato, P., Perina, A., Fasoli, M., Delledonne, M., Pezzotti, M. et al.
1496 Investigating topic models' capabilities in expression microarray data classification.
1497 *IEEE/ACM Trans. Comput. Biol. Bioinform.* **9**, 1831–1836 (2012).
- 1498 37. McCoy, T. H., Castro, V. M., Snapper, L. A., Hart, K. L. & Perlis, R. H. Efficient genome-
1499 wide association in biobanks using topic modeling identifies multiple novel disease loci.

- 1500 *Mol. Med.* **23**, 285–294 (2017).
- 1501 38. McCoy, T. H., Castro, V. M., Snapper, L., Hart, K., Januzzi, J. L., Huffman, J. C. et al.
1502 Polygenic loading for major depression is associated with specific medical comorbidity.
1503 *Transl. Psychiatry* **7**, e1238 (2017).
- 1504 39. McCoy, T. H., Jr, Pellegrini, A. M. & Perlis, R. H. Using phenome-wide association to
1505 investigate the function of a schizophrenia risk locus at SLC39A8. *Transl. Psychiatry* **9**,
1506 45 (2019).
- 1507 40. Zhao, J., Feng, Q., Wu, P., Warner, J. L., Denny, J. C. & Wei, W.-Q. Using topic
1508 modeling via non-negative matrix factorization to identify relationships between genetic
1509 variants and disease phenotypes: A case study of Lipoprotein(a) (LPA). *PLoS One* **14**,
1510 e0212112 (2019).
- 1511 41. Lumbreras, A., Filstroff, L. & Févotte, C. Bayesian mean-parameterized nonnegative
1512 binary matrix factorization. *Data Min. Knowl. Discov.* **34**, 1898–1935 (2020).
- 1513 42. Andrzejewski, D., Zhu, X. & Craven, M. Incorporating Domain Knowledge into Topic
1514 Modeling via Dirichlet Forest Priors. *Proc. Int. Conf. Mach. Learn.* **382**, 25–32 (2009).
- 1515 43. Hu, Y., Boyd-Graber, J., Satinoff, B. & Smith, A. Interactive topic modeling. *Mach.*
1516 *Learn.* **95**, 423–469 (2014).
- 1517 44. Li, C., Rana, S., Phung, D. & Venkatesh, S. Hierarchical Bayesian nonparametric
1518 models for knowledge discovery from electronic medical records. *Knowledge-Based*
1519 *Systems* **99**, 168–182 (2016).
- 1520 45. Cortes, A., Dendrou, C. A., Motyer, A., Jostins, L., Vukcevic, D., Dilthey, A. et al.
1521 Bayesian analysis of genetic association across tree-structured routine healthcare data
1522 in the UK Biobank. *Nat. Genet.* **49**, 1311–1318 (2017).
- 1523 46. Choi, E., Bahadori, M. T., Song, L., Stewart, W. F. & Sun, J. GRAM: Graph-based
1524 Attention Model for Healthcare Representation Learning. *KDD* **2017**, 787–795 (2017).
- 1525 47. Cao, J., Xia, T., Li, J., Zhang, Y. & Tang, S. A density-based method for adaptive LDA
1526 model selection. *Neurocomputing* vol. 72 1775–1781 (2009).

- 1527 48. Wallach, H. M., Murray, I., Salakhutdinov, R. & Mimno, D. Evaluation methods for topic
1528 models. *Proceedings of the 26th Annual International Conference on Machine Learning*
1529 - *ICML '09* (2009) doi:10.1145/1553374.1553515.
- 1530 49. Carrasco, M. V., Manolopoulou, I., O'Sullivan, J., Prior, R. & Musolesi, M. Posterior
1531 summaries of grocery retail topic models: Evaluation, interpretability and credibility.
1532 *Journal of the Royal Statistical Society: Series C (Applied Statistics)* (2022)
1533 doi:10.1111/rssc.12546.
- 1534 50. Wallach, H., Mimno, D. & McCallum, A. Rethinking LDA: Why Priors Matter. in
1535 *Advances in Neural Information Processing Systems* (eds. Bengio, Y., Schuurmans, D.,
1536 Lafferty, J., Williams, C. & Culotta, A.) vol. 22 (Curran Associates, Inc., 2009).
- 1537 51. Minka, T. Estimating a Dirichlet distribution.
1538 <https://vismod.media.mit.edu/pub/tpminka/papers/minka-dirichlet.ps.gz>.
- 1539 52. Islam, M. M., Valderas, J. M., Yen, L., Dawda, P., Jowsey, T. & McRae, I. S.
1540 Multimorbidity and comorbidity of chronic diseases among the senior Australians:
1541 prevalence and patterns. *PLoS One* **9**, e83783 (2014).
- 1542 53. Griffiths, T. L. & Steyvers, M. Finding scientific topics. *Proceedings of the National*
1543 *Academy of Sciences* vol. 101 5228–5235 (2004).
- 1544 54. Eckel, R. H., Grundy, S. M. & Zimmet, P. Z. The metabolic syndrome. *Lancet* **365**,
1545 1415–1428 (2005).
- 1546 55. Cornier, M.-A., Dabelea, D., Hernandez, T. L., Lindstrom, R. C., Steig, A. J., Stob, N. R.
1547 et al. The metabolic syndrome. *Endocr. Rev.* **29**, 777–822 (2008).
- 1548 56. MacArthur, J., Bowler, E., Cerezo, M., Gil, L., Hall, P., Hastings, E. et al. The new
1549 NHGRI-EBI Catalog of published genome-wide association studies (GWAS Catalog).
1550 *Nucleic Acids Res.* **45**, D896–D901 (2017).
- 1551 57. Ernst, J. & Kellis, M. ChromHMM: automating chromatin-state discovery and
1552 characterization. *Nat. Methods* **9**, 215–216 (2012).
- 1553 58. Ng, S. K., Tawiah, R., Sawyer, M. & Scuffham, P. Patterns of multimorbid health

- 1554 conditions: a systematic review of analytical methods and comparison analysis. *Int. J.*
1555 *Epidemiol.* **47**, 1687–1704 (2018).
- 1556 59. Bisquera, A., Gulliford, M., Dodhia, H., Ledwaba-Chapman, L., Durbaba, S., Soley-Bori,
1557 M. et al. Identifying longitudinal clusters of multimorbidity in an urban setting: A
1558 population-based cross-sectional study. *The Lancet Regional Health - Europe* vol. 3
1559 100047 (2021).
- 1560 60. Blei, D. M., Ng, A. Y. & Jordan, M. I. Latent dirichlet allocation. *the Journal of machine*
1561 *Learning research* **3**, 993–1022 (2003).
- 1562 61. Xu, C. & Su, Z. Identification of cell types from single-cell transcriptomes using a novel
1563 clustering method. *Bioinformatics* **31**, 1974–1980 (2015).
- 1564 62. Que, X., Checconi, F., Petrini, F. & Gunnels, J. A. Scalable Community Detection with
1565 the Louvain Algorithm. in *2015 IEEE International Parallel and Distributed Processing*
1566 *Symposium* 28–37 (2015). doi:10.1109/IPDPS.2015.59.
- 1567 63. Bastarache, L. Using Phecodes for Research with the Electronic Health Record: From
1568 PheWAS to PheRS. *Annu Rev Biomed Data Sci* **4**, 1–19 (2021).
- 1569 64. Bulik-Sullivan, B. K., Loh, P.-R., Finucane, H. K., Ripke, S., Yang, J., Schizophrenia
1570 Working Group of the Psychiatric Genomics Consortium et al. LD Score regression
1571 distinguishes confounding from polygenicity in genome-wide association studies. *Nat.*
1572 *Genet.* **47**, 291–295 (2015).
- 1573 65. Heritability and Genetic Correlation. *GitHub* <https://github.com/bulik/ldsc>.
- 1574 66. Burdett, T., Hastings, E., Welter, D., SPOT, EMBL-EBI & NHGRI. GWAS Catalog.
1575 <https://www.ebi.ac.uk/gwas/>.
- 1576 67. Cano-Gamez, E. & Trynka, G. From GWAS to Function: Using Functional Genomics to
1577 Identify the Mechanisms Underlying Complex Diseases. *Front. Genet.* **11**, 424 (2020).
- 1578 68. Watanabe, K., Taskesen, E., van Bochoven, A. & Posthuma, D. Functional mapping
1579 and annotation of genetic associations with FUMA. *Nat. Commun.* **8**, 1826 (2017).
- 1580 69. Functional Mapping and Annotation of Genome-wide association studies.

- 1581 <https://fuma.ctglab.nl/>.
- 1582 70. GEO Accession viewer.
- 1583 <https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE87112>.
- 1584 71. GTEx Portal. <https://www.gtexportal.org/home/>.
- 1585 72. Euesden, J., Lewis, C. M. & O'Reilly, P. F. PRSice: Polygenic Risk Score software.
- 1586 *Bioinformatics* **31**, 1466–1468 (2015).