1	
2	
3	
4	Using Topic Modeling via Non-negative Matrix Factorization to Identify
5	Relationships between Genetic Variants and Disease Phenotypes: A Case
6	Study of Lipoprotein(a) (LPA)
7	
8 9	Juan Zhao ¹ , QiPing Feng ² , Patrick Wu ¹ , Jeremy L. Warner ^{1,3} , Joshua C. Denny ^{1,3} , Wei-Qi Wei ^{1*}
10	¹ Department of Biomedical Informatics, Vanderbilt University Medical Center, Nashville, TN, USA
11	² Division of Clinical Pharmacology, Vanderbilt University Medical Center, Nashville, TN, USA
12	³ Department of Medicine, Vanderbilt University Medical Center, Nashville, TN, USA
13	
14	* Corresponding author
15	Email: wei-qi.wei@vanderbilt.edu
16	Department of Biomedical Informatics
17	2525 West End Ave., Suite 1500
18	Nashville, TN 37203
19	Tel: (615)343-1956
20 21	

22 Abstract

23 Genome-wide and phenome-wide association studies are commonly used to identify 24 important relationships between genetic variants and phenotypes. Most of these studies have 25 treated diseases as independent variables and suffered from heavy multiple adjustment burdens 26 due to the large number of genetic variants and disease phenotypes. In this study, we propose 27 using topic modeling via non-negative matrix factorization (NMF) for identifying associations between disease phenotypes and genetic variants. Topic modeling is an unsupervised machine 28 learning approach that can be used to learn the semantic patterns from electronic health record 29 30 data. We chose rs10455872 in LPA as the predictor since it has been shown to be associated with 31 increased risk of hyperlipidemia and cardiovascular diseases (CVD). Using data of 12,759 32 individuals from the biobank at Vanderbilt University Medical Center, we trained a topic model 33 using NMF from 1,853 distinct phecodes extracted from the cohort's electronic health records 34 and generated six topics. We quantified their associations with rs10455872 in LPA. Topics 35 indicating CVD had positive correlations with rs10455872 (P < 0.001), replicating a previous 36 finding. We also identified a negative correlation between LPA and a topic representing lung 37 cancer (P < 0.001). Our results demonstrate the applicability of topic modeling in exploring the 38 relationship between the genome and clinical diseases.

40 Author summary

41	Identifying the clinical associations of genetic variants remains crucial in understanding
42	how the human genome modulates disease risk. Traditional phenome-wide association studies
43	consider each disease phenotype as an independent variable, however, diseases often present as
44	complex clusters of comorbid conditions. In this study, we propose using topic modeling to
45	model electronic health record data as a mixture of topics (e.g., disease clusters or relevant
46	comorbidities) and testing associations between topics and genetic variants. Our results
47	demonstrated the feasibility of using topic modeling to replicate and discover novel associations
48	between the human genome and clinical diseases.
49	

50

52 Introduction

53	Elucidating associations between genetic variants and human diseases creates new
54	avenues for disease prevention and enables more precise treatment of diseases [1,2]. During the
55	past two decades, genetic studies have uncovered thousands of genetic variants that influence
56	risk for disease phenotypes [3], e.g., the discovery of a variant in proprotein convertase
57	subtilisin/kexin type 9 (PCSK9[4]) associated with low plasma low-density lipoprotein, which
58	led to a new therapeutic drug class that was approved by the US Food and Drug Administration
59	in 2015. Many of these discoveries come from large-scale association analyses. The two most
60	notable approaches are genome-wide (GWAS) and phenome-wide association studies (PheWAS)
61	[2, 5]. For a given phenotype, GWAS scans hundreds of thousands to millions of single
62	nucleotide polymorphisms (SNPs) across the genome in a hypothesis-free approach. PheWAS,
63	on the contrary, analyzes thousands of disease phenotypes compared to a single SNP. In a
64	GWAS, the outcome variable is a disease phenotype and predictor variables are SNPs. While in
65	a PheWAS, the outcome variable is a SNP and predictor variables are disease phenotypes.
66	Association analyses test a large number of predictor variables at one time and assume
67	that each variable has an independent effect. However, diseases often occur together as a group
68	of comorbidities, e.g. hyperlipidemia (HLD) and cardiovascular diseases (CVDs). Conventional
69	association analyses may not capture the inter-connections among variables such as phenotypes
70	and thus may not be sensitive enough to identify important genotype-phenotype relationships.
71	Moreover, association analyses also face the challenge of scaling to an increasing number of
72	phenotypes. Previously, we have described a "networked PheWAS" approach which can address
73	interconnectivity but still requires a degree of supervised interpretation [6].

74 This study aimed to test the feasibility of topic modeling for identifying relationships 75 between genetic variants and disease phenotypes. Topic modeling is an unsupervised machine 76 learning method that was initially introduced as a text mining technique [7]. It has been 77 demonstrated to extract latent topics or themes from documents, aiding in the understanding of 78 large amounts of data [8]. Compared to traditional clustering approaches such as K-means 79 clustering that partitions a collection of documents into several disjoint clusters (i.e., topics) 80 based on a similarity measure, topic modeling assigns a document to multiple clusters with 81 different scores. Therefore, each document is characterized by one or more topics. In addition to 82 its wide adoption in the text mining field, topic modeling has achieved many successes in 83 computer vision and biomedical science. Recently, a few groups have used this approach to 84 analyze electronic health records (EHRs) [9,10] and genetic data to capture the characteristic of 85 data [11,12].

86 We hypothesized that topic modeling would be useful in replicating known findings and 87 uncovering previously unidentified relationships between genetic variants and disease 88 phenotypes. To test this hypothesis, we used topic modeling via non-negative matrix 89 factorization (NMF) [13,14] to identify latent topics (e.g. disease clusters or relevant 90 comorbidities) from EHR data. We then tested associations between the EHR-derived topics and 91 a LPA SNP (rs10455872). We chose the SNP because previous studies have shown that high-92 levels of the LPA product (Lp(a)) is associated with increased risks of developing HLD and 93 CVD [15]. Specifically, rs10455872, as a single variant, explains 20-30% of the variation in 94 circulating Lp(a) levels, which makes it an ideal candidate for this study [16].

95 **Results**

96	We applied a topic modeling algorithm using NMF on the dataset of 12,759 individuals
97	and obtained six potentially meaningful topics from the EHRs (Fig 1). The learned topics (i.e.,
98	clusters of disease phenotypes) were consistent with the comorbidities associated with the
99	phenotypes most prevalent in the cohort. For example, topic #0 represented diseases of
100	respiratory failure, topic #2 defined diseases related to CVD (e.g., HLD, hypertension, and
101	chronic ischemic heart disease), topic #3 represented phenotypes relevant to lung cancer and its
102	treatment, topic #4 was related to diabetes and its comorbidities; and topic #5 was related to liver
103	disease and its sequelae.
104	
105	Fig 1. Word clouds for six topics. The size of the words (phecode) in each cloud
106	indicate the weights of the phenotypes on the topic. Phenotypes with larger-sized words had
107	greater influence on the topic compared to phenotypes with smaller-sized words. For each word
108	cloud, we listed the top 60 words to provide a better visual presentation of what each topic
109	represents.
110	Fig 2 shows the distribution of the numbers of topics in the cohort. Topic #2 was the most
111	prevalent (33%) topic in the cohort. Topics #1 and #3 were the second and third most prevalent
112	topics in the cohort.
113	
114	Fig 2. Topic distribution in the cohort. To visualize the prevalence of each topic in the
115	cohort, we assigned an individual to the topic with the maximum score.
116	We also used t-Distributed Stochastic Neighbor Embedding (t-SNE) [17] to transform the
117	individual-phenotypes matrix (W) into a 2-dimensional (2D) space to visualize the quality of

topic modeling (Fig 3). Each data point in the figure corresponds to one individual. We labeledeach individual with the assigned topic.

120

121 Fig 3. t-SNE plot of visualizing the patient clusters in a projected 2D metric map.

122 The perplexity was set to 30. We used PCA initialization as it is more globally stable. Each point

123 represents an individual. Topic #2 contains the most individuals in the cohort.

124 We then applied the Pearson correlation coefficient (PCC) to examine the association

between each topic and rs10455872. Statistical test results suggest that topic #2 and #3 were

significantly associated with rs10455872 (Table 1). Topic #2, a group of lipid and cardiovascular

127 diseases, had a positive correlation with rs10455872 (r=0.072, p=5.8e-16). We also found that

128 topic #3, a group of phenotypes relevant to lung cancer, had a negative correlation with

129 rs10455872 (r=-0.039, p=8.5e-6). Although the *r* coefficient is weaker than the topic#2, these

- 130 correlations are highly statistically significant.
- 131

Table 1. Pearson correlation between LPA variant for each topic

Topic	Top phenotypes in this topic	r	P value
#0	Respiratory failure, Pneumonia, Pleurisy, Pulmonary	0.011	0.199
	collapse; interstitial/compensatory emphysema, Hypotension		
	NOS, Tachycardia NOS, Other dyspnea, Hypopotassemia,		
	Sepsis, Septicemia		
#1	Pain in joint, Other tests, Back pain, Pain in limb, Malaise	-0.008	0.358
	and fatigue, Cough, Nonspecific chest pain, Essential		
	hypertension, Osteoarthrosis NOS, Abdominal pain		

#2	Coronary atherosclerosis, Essential hypertension,	0.072	5.8e-16
	Hyperlipidemia, Congestive heart failure NOS, Nonspecific		
	chest pain, Atrial fibrillation, Chronic ischemic heart disease,		
	Shortness of breath, Nonrheumatic mitral valve disorders,		
	Cardiomegaly		
#3	Chemotherapy, Tobacco use disorder, Lung cancer, Other	-0.039	8.5e-6
	diseases of lung, Malaise and fatigue, Secondary malignancy		
	of lymph nodes, Secondary malignancy of lung, Nausea and		
	vomiting, Nonspecific chest pain, Shortness of breath		
#4	Type 2 diabetes, Hypertensive chronic kidney disease,	0.002	0.783
	Chronic renal failure, Insulin pump user, Type 2 diabetic		
	neuropathy, Chronic Kidney Disease, Stage III, Type 2		
	diabetic nephropathy, Type 1 diabetes, Polyneuropathy in		
	diabetes, Acute renal failure		
#5	Ascites (nonmalignant), Abdominal pain, Cirrhosis of liver	-0.02	0.021
	without mention of alcohol, Thrombocytopenia, Liver		
	abscess and sequelae of chronic liver disease, Portal		
	hypertension, Chronic nonalcoholic liver disease, Disorders		
	of liver, Esophageal bleeding, Nausea and vomiting		

132

133 **Discussion**

Topic modeling has been widely used in the field of text mining. In this paper, we applied
this technique to explore associations between disease phenotypes and genetic variants. We

136	assumed that some disease phenotypes found simultaneously in a large EHR have correlated
137	semantic meanings and thus can be learned as topics. We examined the associations between a
138	LPA variant (rs10455872) and the six topics derived from EHRs. We observed the expected
139	association between rs10455872 and a topic representing CVD/HLD. We also found a novel
140	association, as of this writing [18], between the LPA variant and a lung cancer topic.
141	The LPA gene encodes lipoprotein (a), a major component of the Lp(a) particle.
142	Individuals with elevated Lp(a) levels are more likely to develop CVD compared to those with
143	normal or low Lp(a) level [16,19]. Approximately 70% of Lp(a) variation can be attributed to
144	variants at the LPA locus [20–22], and rs10455872 alone explains ~25% variation in circulating
145	Lp(a) levels [16]. Further, a previous genetic study suggested that LPA variants were strong
146	predictors for CVD risk [16]. In a more recent study of >10,000 patients taking statins, our group
147	found that rs10455872 predicted residual CVD risk while on lipid-lowering treatment [23]. This
148	study's finding of a significant association between rs10455872 and the CVD/HLD topic
149	demonstrates the feasibility of topic modeling as a critical tool for uncovering genotype-
150	phenotype relationships.
151	We also observed a negative correlation between the LPA variant and the cancer/lung

151 We also observed a negative correlation between the LPA variant and the cancer/lung 152 cancer topic, i.e., possessing this variant is protective. Previous epidemiological studies have 153 reported that individuals with low Lp(a) levels had increased risk of all-cause and cancer-related 154 mortality [24]. Mieno et al. found that hypolipoproteinemia(a) is a risk factor for cancer except 155 for lung cancer. Nevertheless, there are few reports on a relationship between cancer and LPA 156 polymorphism or expression levels. Our previous PheWAS analysis of a separate cohort 157 identified an association between rs10455872 and cancer diagnosis code with borderline 158 significance [23]. To further explore this association between rs10455872 and the cancer/lung

159 cancer topic, we queried gene2pheno (https://imlab.shinyapps.io/gene2pheno ukb neale/), 160 which is a publicly available database for testing associations between predicted gene expression 161 levels and phenotypes using data from the UK Biobank. Genetically predicted LPA expression 162 levels were associated with death from T cell lymphomas (p=6.9 e-5, Underlying (primary) 163 cause of death: ICD10: C84.5 Other and unspecified T-cell lymphomas). Given that lung cancer 164 is strongly mediated by environmental exposure and that tobacco use disorder was also part of 165 topic #4, it is possible that the SNP is a marker for propensity to smoking, e.g., similar to what 166 was shown for rs16969968 [25]. Further genetic and epidemiological studies are needed to 167 elucidate the relationship between Lp(a) levels and cancer incidence. 168 Topic modeling approaches require pre-specification of the number of topics. In this 169 study, we set k=6, because we aimed to capture the most prevalent diseases such as CVD and to 170 quantify the association. Increasing k allows the quantification of associations between genetic 171 variants and rare diseases but risks fracturing common phenotype clusters. It can be seen that 172 (Fig 3), except for topic #4 (diabetes), the learned topics formed distinct clusters, indicating a 173 good quality of topic modeling. Some of points in topic #4 (diabetes) were close with topic #2 174 (CVD), which was expected, because type II diabetes is an important risk factor that increases 175 the risk of developing CVD. Compared to the other topics, #1 (Pain), #2 (CVD), and #3 (Lung 176 Cancer) have more concentrated clusters. 177 For optimal selection of k, common approaches have used different values of k to look at

the error in optimization and selected the best value by having domain experts review the topics to identify which set of topics are most meaningful, and have estimated k using singular value decomposition (SVD) to look at the decay of singular values [26–28]. To provide evidence for 181 the stabilities of our results, we also set different numbers of topics k=10, 20, 30 (Supplementary 182 Table 1) and examined the PCC. Results were consistent with topics at k=6.

In summary, unlike traditional PheWAS that have treated each disease phenotype as a distinct variable, topic modeling via NMF generates more abstract latent factors from disease phenotypes and significantly reduces the number of multiple tests. Our results demonstrate the power of topic modeling in the detection of comorbidities and previously unexplored genotypephenotype relationships among a large cohort.

188 **Limitations**

189 There are several limitations in this study. First, we tested only one genetic variant in one 190 gene. Rs10455872 explains approximate 25% change in circulating Lp(a) levels according to 191 previous studies; however, it would be interesting to generate a genetic risk score for Lp(a) levels 192 and test its association with disease phenotypes in the future. Second, we used a binary value to 193 indicate if an individual had a diagnosis code. A method accounting for disease severity (e.g., 194 counts of diagnosis codes) could be used in future studies. Finally, the current study was limited 195 to using billing codes to phenotype individuals. We did not include other information, e.g. 196 laboratory test and medications, to assign more accurate phenotypes. This problem can be solved 197 in the future using more sophisticated "deep" phenotyping methods that include more features 198 from EHRs.

200 Materials and methods

201 Study cohort

We used data from BioVU, the de-identified DNA biobank at Vanderbilt University Medical Center (VUMC), to conduct this study. BioVU contains DNA samples from >250,000 individuals that are linked with their de-identified EHRs, including diagnostic and procedure codes, clinical notes, laboratory values and medications. We identified 12,759 adult individuals of European ancestry (F/M: 6,018/6,741; age: 70.3 \pm 12.3) who had both EHRs and genotyped data of rs10455872 available.

208 rs10455872 Genotyping

209 We extracted each individual's rs10455872 information from existing genotyped data.

210 All genotyping was previously conducted using commercially available genome-wide SNP

211 arrays with quality control criteria for variants followed by a standard imputation process using

212 1000 Genomes Project allele frequency estimates.

Among the cohort of 12,759 individuals, we observed 85.2% AA, 14.2% AG, 6.1% GG. The minor allele frequency (MAF) of the rs10455872 G allele is 7.7% in our cohort, consistent with the 7% MAF in the European population [29]. We used 0, 1, 2 to represent the number of *LPA* rs10455872 G alleles that an individual carry.

217 **Disease Phenotypes**

Following established protocols used in past studies [30], we grouped each individual's
ICD-9-CM (International Classification of Disease, 9th edition) codes into disease phecodes.

220	There were 1853 phecodes present in the 12,759 individuals. For each phecode, we labeled
221	individuals without the phecode with a '0' and those with the phecode with a '1'.
222	We applied a topic modeling algorithm using NMF on the dataset of 12,759 individuals
223	to learn potentially meaningful topics from the EHRs. Then, we quantified the association
224	between each learned topic with rs10455872 using PCC. The workflow of this experiment is
225	demonstrated in Fig 4.
226	
227	Fig 4. Illustration of topic modeling on EHRs using NMF
228	Topic modeling via Non-negative Matrix Factorization (NMF)
229	We used NMF as our topic modeling approach. NMF is a low-rank matrix approximation
230	algorithm that has been widely used for feature reduction for high-dimensional data. The
231	assumption is that given a large and sparse matrix X of size $n \times m$ representing a collection of n
232	high dimensional data points in R^m . X is low rank which means that most data points can be
233	approximately represented by a linear combination of a small set of k basis vectors $H \in \mathbb{R}^{k \times n}$.
234	The linear combination is a coefficients matrix $W \in \mathbb{R}^{n \times k}$ providing a lower-dimensional
235	encoding for X, which result in a feature reduction for a high-dimensional data. Since NMF
236	restricts the X non-negative and enforces the H and W to be also non-negative, NMF has good
237	interpretability and has been commonly used as a topic modeling approach in text mining.
238	We considered each individual's EHR as one document, and each document was
239	described by disease phenotypes represented by the phecodes (Fig 1). Since we had 12,759 (n)
240	individuals' EHRs and 1,853 (<i>m</i>) unique phecodes, we used matrix $X \in \mathbb{R}^{n \times m}$ to represent the
241	input data, where each row of X represented an individual, and each column of X was a phecode.
242	The entry of the matrix $X_{ij} \in X$ was a binary value (0 or 1) indicating whether <i>i</i> th individual had

the *j*th phecode. This representation is similar with the bag-of-word model, where each documentis associated with a set of words, and word ordering in the documents does not matter.

Given that non-negative input matrix $\in \mathbb{R}^{n \times m}$, and an expected number of topics $k \leq m$, NMF generates two matrices, $W \in \mathbb{R}^{n \times k}$, and $H \in \mathbb{R}^{k \times m}$. Both W and H are non-negative entries, such that

248

$$\min_{W \ge 0, H \ge 0} \|X - WH\|_F^2 + \lambda R(W, H)$$
(1)

249 H(k,:) is a latent topic – phenotype matrix. Specifically, each row of *H* corresponds to a 250 disease topic, and each topic is represented by a set of relevant phenotypes that co-occurred in 251 several individuals' EHRs, with specific cores indicating their relevance to this topic. Through 252 H(k,:), we extracted a semantic meaning of each topic, e.g. what kind of diseases or 253 comorbidities are represented by the topic.

W(i,k) is an individual-topic matrix. Each row of *W* corresponds to an individual's score on each topic that indicates the diseases and comorbidities carried by the individual. An individual that has a large score for a disease topic indicates that there a higher probability for an association between individual and the topic. W(i,k) is then used for the association calculation between the topics and rs10455872, which is described further below.

259 R(W, H) is the regularization term that combines L1 and L2 norms, which is defined as:

260
$$R(W,H) = \gamma (||W||_F + ||H||_F) + \frac{1}{2}(1-\gamma) (||W||_F^2 + ||H||_F^2), \qquad (3)$$

261 where γ is the ratio for L1 penalty. Adding the regularization term is necessary for 262 balancing the sparsity of the topics, meaning that an individual may have several topics at the 263 same time. Moreover, addition of the regularization term minimizes the effect of outliers on the 264 model.

265 Statistical analysis

266	We applied the PCC to quantify the association between each individuals' scores on

- specific topic and each individual's rs10455872 status, for each learned topic. PCC measures the
- strength of a linear association between two variables. PCC also can generate a correlation

269 coefficient denoted by $r \in [-1,1]$, which shows the direction of the correlation.

- 270 We used the individual-topic matrix, $W \in \mathbb{R}^{n \times k}$, generated by NMF to calculate the PCC
- 271 with the genetic variants. Each column vector of W(:,j) of the matrix W represented a topic
- vector with scores on all the individuals. We used each column vector as the predictor variable *x*
- and the number of minor alleles (0, 1, or 2) at rs10455872 of each patient as variable y.

274 **References**

- Denny JC, Van Driest SL, Wei W-Q, Roden DM. The Influence of Big (Clinical) Data and Genomics on Precision Medicine and Drug Development. Clinical Pharmacology & Therapeutics. 2018;103: 409–418. doi:10.1002/cpt.951
- Manolio TA. Genomewide Association Studies and Assessment of the Risk of Disease. Feero
 WG, Guttmacher AE, editors. New England Journal of Medicine. 2010;363: 166–176.
 doi:10.1056/NEJMra0905980
- Welter D, MacArthur J, Morales J, Burdett T, Hall P, Junkins H, et al. The NHGRI GWAS
 Catalog, a curated resource of SNP-trait associations. Nucleic Acids Res. 2014;42: D1001 1006. doi:10.1093/nar/gkt1229
- Cohen JC, Boerwinkle E, Mosley TH, Hobbs HH. Sequence variations in PCSK9, low LDL,
 and protection against coronary heart disease. N Engl J Med. 2006;354: 1264–1272.
 doi:10.1056/NEJMoa054013
- 5. Denny JC, Ritchie MD, Basford MA, Pulley JM, Bastarache L, Brown-Gentry K, et al.
 PheWAS: demonstrating the feasibility of a phenome-wide scan to discover gene-disease
 associations. Bioinformatics. 2010;26: 1205–1210. doi:10.1093/bioinformatics/btq126
- Warner JL, Denny JC, Kreda DA, Alterovitz G. Seeing the forest through the trees:
 uncovering phenomic complexity through interactive network visualization. J Am Med
 Inform Assoc. 2015;22: 324–329. doi:10.1136/amiajnl-2014-002965

- Arora S, Ge R, Moitra A. Learning Topic Models Going Beyond SVD. Proceedings of the
 2012 IEEE 53rd Annual Symposium on Foundations of Computer Science. Washington, DC,
 USA: IEEE Computer Society; 2012. pp. 1–10. doi:10.1109/FOCS.2012.49
- 8. MacMillan K, Wilson JD. Topic supervised non-negative matrix factorization.
 arXiv:170605084 [cs, stat]. 2017; Available: http://arxiv.org/abs/1706.05084
- Pinoli P, Chicco D, Masseroli M. Enhanced probabilistic latent semantic analysis with
 weighting schemes to predict genomic annotations. 13th IEEE International Conference on
 BioInformatics and BioEngineering. 2013. pp. 1–4. doi:10.1109/BIBE.2013.6701702
- Wahabzada M, Mahlein A-K, Bauckhage C, Steiner U, Oerke E-C, Kersting K. Plant
 Phenotyping using Probabilistic Topic Models: Uncovering the Hyperspectral Language of
 Plants. Scientific Reports. 2016;6: 22482. doi:10.1038/srep22482
- McCoy TH, Castro VM, Snapper LA, Hart KL, Perlis RH. Efficient Genome-wide
 Association in Biobanks Using Topic Modeling Identifies Multiple Novel Disease Loci. Mol
 Med. 2017;23: 285–294. doi:10.2119/molmed.2017.00100
- 307 12. Backenroth D, He Z, Kiryluk K, Boeva V, Pethukova L, Khurana E, et al. FUN-LDA: A
 308 Latent Dirichlet Allocation Model for Predicting Tissue-Specific Functional Effects of
 309 Noncoding Variation: Methods and Applications. The American Journal of Human Genetics.
 310 2018;102: 920–942. doi:10.1016/j.ajhg.2018.03.026
- 311 13. Sra S, Dhillon IS. Generalized Nonnegative Matrix Approximations with Bregman 312 Divergences. In: Weiss Y, Schölkopf B, Platt JC, editors. Advances in Neural Information 313 Processing **Systems** 18. MIT Press: 2006. pp. 283 - 290.Available: 314 http://papers.nips.cc/paper/2757-generalized-nonnegative-matrix-approximations-with-315 bregman-divergences.pdf
- Kim H, Park H. Sparse non-negative matrix factorizations via alternating non-negativity constrained least squares for microarray data analysis. Bioinformatics. 2007;23: 1495–1502.
 doi:10.1093/bioinformatics/btm134
- Nordestgaard BG, Chapman MJ, Ray K, Borén J, Andreotti F, Watts GF, et al. Lipoprotein(a)
 as a cardiovascular risk factor: current status. Eur Heart J. 2010;31: 2844–2853.
 doi:10.1093/eurheartj/ehq386
- 16. Clarke R, Peden JF, Hopewell JC, Kyriakou T, Goel A, Heath SC, et al. Genetic Variants
 Associated with Lp(a) Lipoprotein Level and Coronary Disease. New England Journal of
 Medicine. 2009;361: 2518–2528. doi:10.1056/NEJMoa0902604
- Maaten L van der, Hinton G. Visualizing Data using t-SNE. Journal of Machine Learning
 Research. 2008;9: 2579–2605.
- 32718. rs10455872 SNPedia [Internet]. [cited 23 May 2018]. Available:328https://www.snpedia.com/index.php/Rs10455872

- 19. Khera AV, Emdin CA, Drake I, Natarajan P, Bick AG, Cook NR, et al. Genetic Risk,
 Adherence to a Healthy Lifestyle, and Coronary Disease. New England Journal of Medicine.
 2016;375: 2349–2358. doi:10.1056/NEJMoa1605086
- Barlera S, Specchia C, Farrall M, Chiodini BD, Franzosi MG, Rust S, et al. Multiple QTL
 influence the serum Lp(a) concentration: a genome-wide linkage screen in the PROCARDIS
 study. Eur J Hum Genet. 2007;15: 221–227. doi:10.1038/sj.ejhg.5201732
- 335 21. Berglund L, Ramakrishnan R. Lipoprotein(a): an elusive cardiovascular risk factor.
 336 Arterioscler Thromb Vasc Biol. 2004;24: 2219–2226.
 337 doi:10.1161/01.ATV.0000144010.55563.63
- 338 22. Sandholzer C, Hallman DM, Saha N, Sigurdsson G, Lackner C, Császár A, et al. Effects of
 339 the apolipoprotein(a) size polymorphism on the lipoprotein(a) concentration in 7 ethnic
 340 groups. Hum Genet. 1991;86: 607–614.
- Wei W-Q, Li X, Feng Q, Kubo M, Kullo IJ, Peissig PL, et al. LPA Variants are Associated
 with Residual Cardiovascular Risk in Patients Receiving Statins. Circulation. 2018;
 CIRCULATIONAHA.117.031356. doi:10.1161/CIRCULATIONAHA.117.031356
- 24. Low Lipoprotein(a) Concentration Is Associated with Cancer and All-Cause Deaths: A
 Population-Based Cohort Study (The JMS Cohort Study) [Internet]. [cited 14 May 2018].
 Available: http://journals.plos.org/plosone/article?id=10.1371/journal.pone.0031954
- Lips EH, Gaborieau V, McKay JD, Chabrier A, Hung RJ, Boffetta P, et al. Association
 between a 15q25 gene variant, smoking quantity and tobacco-related cancers among 17 000
 individuals. Int J Epidemiol. 2010;39: 563–577. doi:10.1093/ije/dyp288
- Bioucas-Dias JM, Nascimento JMP. Estimation of signal subspace on hyperspectral data.
 Image and Signal Processing for Remote Sensing XI. International Society for Optics and
 Photonics; 2005. p. 59820L. doi:10.1117/12.620061
- Tan VYF, Févotte C. Automatic Relevance Determination in Nonnegative Matrix
 Factorization with the /spl beta/-Divergence. IEEE Transactions on Pattern Analysis and
 Machine Intelligence. 2013;35: 1592–1605. doi:10.1109/TPAMI.2012.240
- Kanagal B, Sindhwani V. Rank Selection in Low-rank Matrix Approximations: A Study of
 Cross-Validation for NMFs. : 5.
- 29. 1000 Genomes Project Consortium, Auton A, Brooks LD, Durbin RM, Garrison EP, Kang
 HM, et al. A global reference for human genetic variation. Nature. 2015;526: 68–74.
 doi:10.1038/nature15393
- 30. Martin PA, Thorburn MJ, Smith-Read EH. Chromosomal rearrangements in three
 generations of a Jamaican family. A possible further example of recombinational imbalance.
 Cytogenetics. 1970;9: 360–368.

365 Supporting Information

366 S1 Table. Results with *topic k=10, 20,30*

- 368
- 369



Topic #2



Topic #3

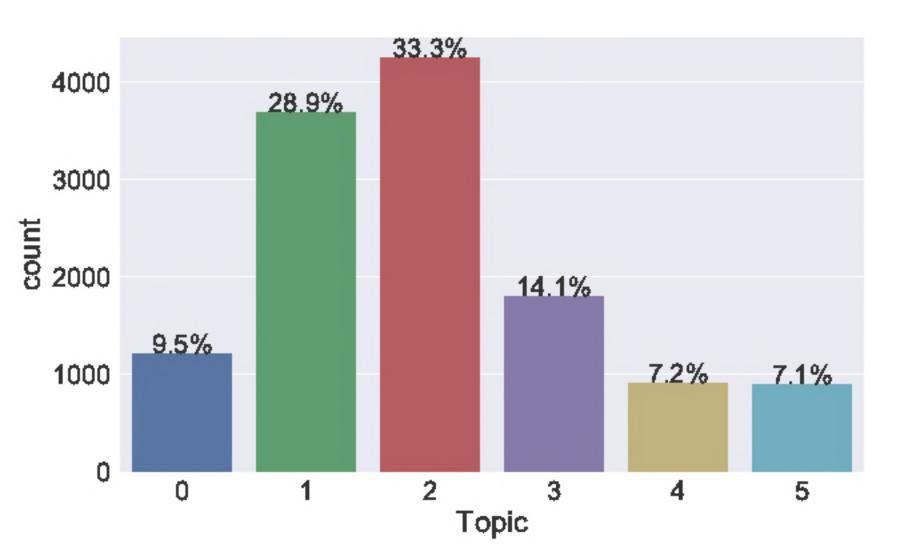


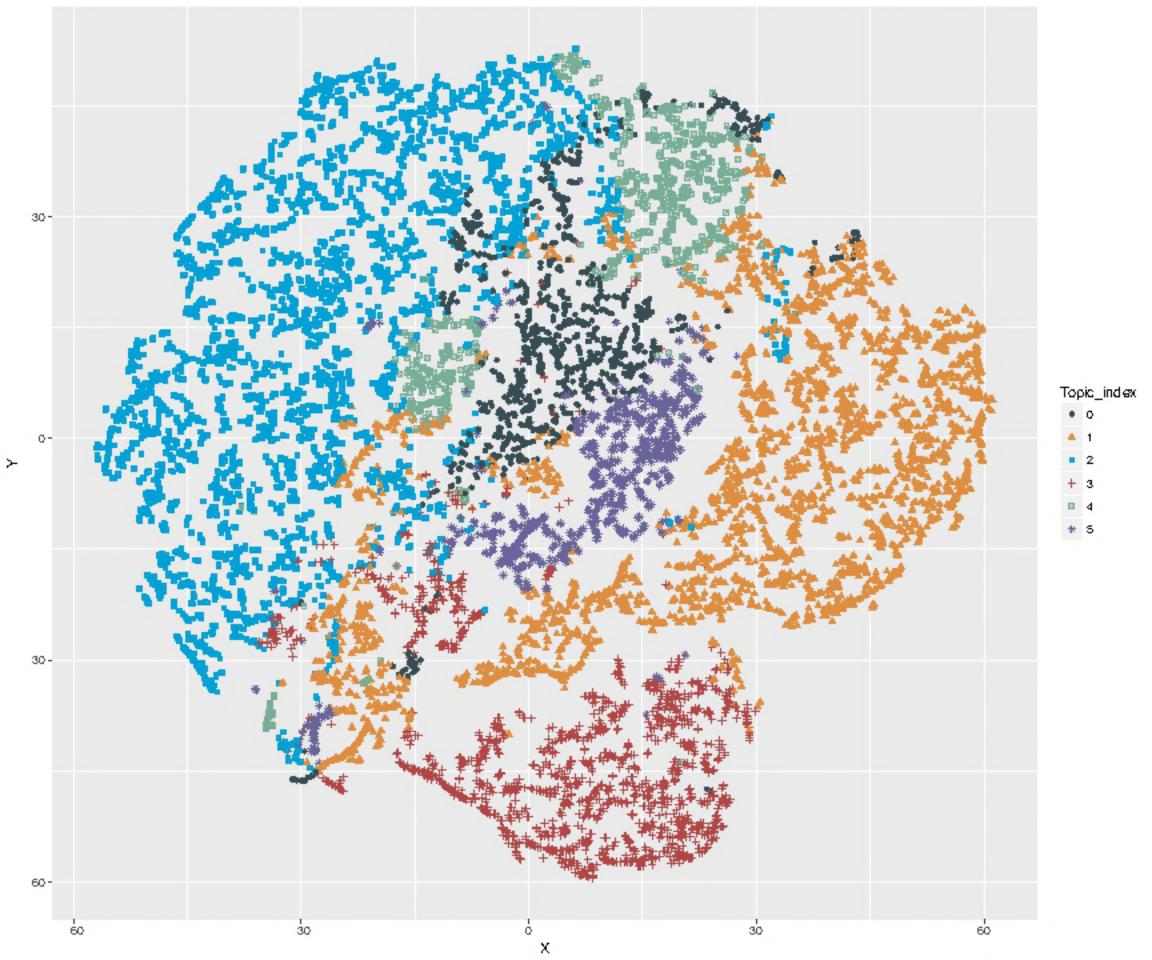
Topic #4

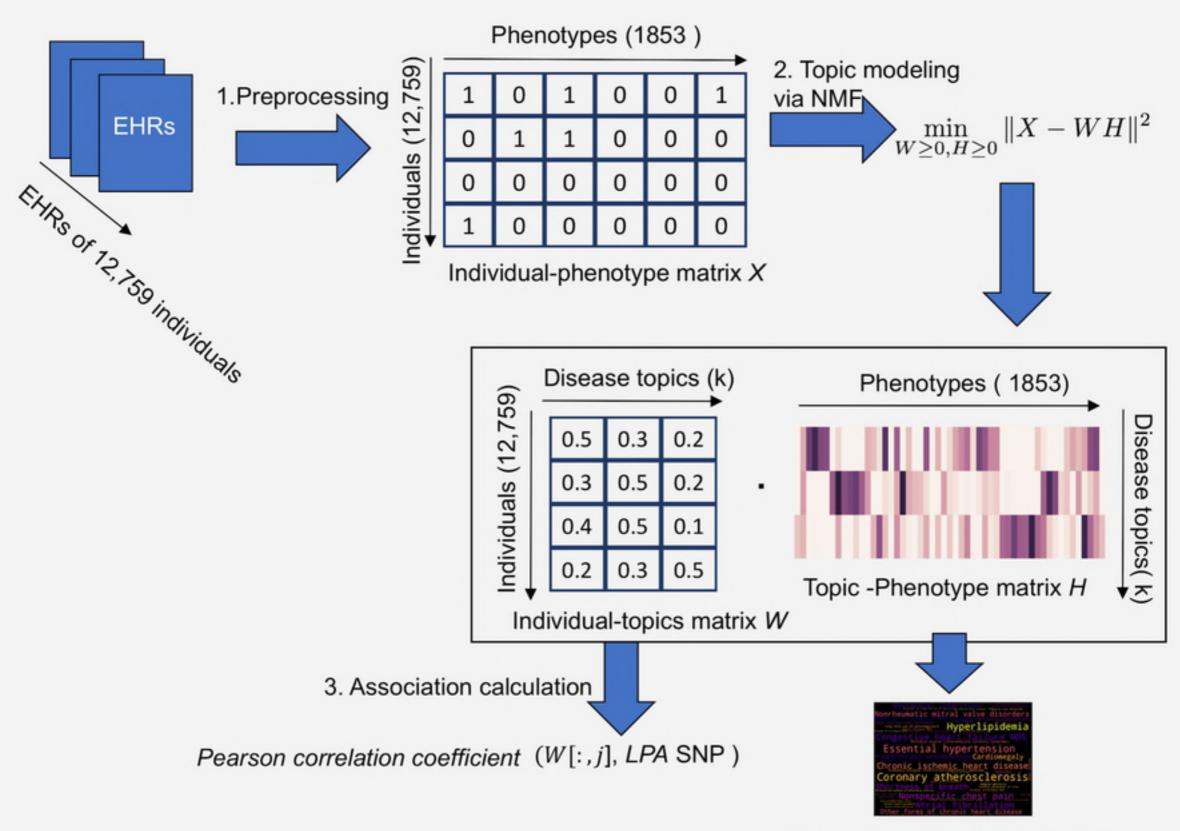
Anemia in chronic kidney disease, Stage III Type 2 diabetic nephropathy Stud overload Chronic Kidney Disease, Stage IV Stud overload Chronic Kidney Disease, Stage IV Stud overload Sypertensive chronic kidney & Stud overload Disorders of the kidney & Stude overload Stude overload Disorders of the kidney & Stude overload Disorders of the kidney & Stude overload Disorders of the kidney & Stude overload Stude overload Disorders of the kidney & Stude overload Stude overload Disorders of the kidney & Stude overload Stude overload Stude overload Disorders of the kidney & Stage Stude overload Stude overl

Topic #5









^{4.} Topic - Phenotype visualization