bioRxiv preprint doi: https://doi.org/10.1101/737049. The copyright holder for this preprint (which was not peer-reviewed) is the author/funder. It is made available under a CC-BY 4.0 International license.

Using indication embeddings to represent patient health for drug safety

studies

Rachel D. Melamed, PhD^{1,2}

¹Biomedical Data Science, University of Chicago, 900 E 57 St, Chicago, IL, USA

²melamed@uchicago.edu

Key words: representation, causal inference, cohort studies, embeddings, dimensionality

reduction

Abstract

Objective

The electronic health record is a rising resource for quantifying medical practice and discovering adverse effects of drugs. One of the challenges of applying these methods to health care data is the high dimensionality of the health record. Methods to discover effects of drugs in health data must account for tens of thousands of potentially relevant confounders. Our goal in this work is to reduce the dimensionality of the health data with the aim of accelerating the application of retrospective cohort studies to this data.

Materials and Methods

In this work, we develop indication embeddings, a way to reduce the dimensionality of health data while capturing the information relevant to treatment decisions. We evaluate these embeddings using external data on drug indications. Then, we use the embeddings as a substitute for medical history to match patients, and develop evaluation metrics for these matches.

Results

We demonstrate that these embeddings recover therapeutic uses of drugs. We use embeddings as an informative representation of relationships between drugs, between health history events and drug prescriptions, and between patients at a particular time in their health history. We show that using embeddings to match cohorts improves the balance of the cohorts in terms of poorly measured risk factors like smoking.

Discussion and Conclusion

Unlike other embeddings inspired by word2vec, indication embeddings are specifically designed to capture the medical history leading to prescription of a new drug. For retrospective cohort studies, our low-dimensional representation helps in finding comparator drugs and constructing comparator cohorts.

Introduction

The EHR is an increasingly complete record of human health and medical practice. Soon EHR will be combined with genomic information, wearable sensors, and other health information to enable scientific discovery from human health histories. Recent studies have used the EHR to evaluate variation in drug treatment decisions.[1,2] Researchers aiming to create a complete knowledge base of health and medicine have harnessed EHR data to identify symptoms related

to diseases,[3,4] and to annotate the medical conditions for which drugs are prescribed, known as the indications for a drug.[5,6] Another area of research uncovers adverse effects among people taking a drug, as a way to complement randomized trials.[7–9]

All of these efforts must confront the high dimensional nature of health data: there are over 10,000 ICD-9 codes specifying particular diagnoses, and thousands of common drugs, alongside other medical events such as procedures and tests. The most straightforward way to code such discrete data for use in an algorithm results in a high-dimensional vector that is sparse, such as a one-hot (dummy variable) vector for each ICD code or drug (Figure 1A). Furthermore, these vectors lack information: the sparse, high-dimensional vector that represents a diagnosis for diabetes will be no nearer to the vector for insulin, than to the vector for acne. To address this issue, a number of studies have explored alternative lower-dimensional representations of health data.

In particular, Miotto[10] proposed autoencoders to create dense vector summaries of patient health at a time point, and they showed that these vectors could be used to predict disease incidence. Choi[11] adapted the Word2Vec method to medical histories. This method, originally popularized for representing natural language[12] transforms the high-dimensional sparse representation of a word to a low-dimensional vector such that words that are used together have closer vectors. These representations have shown success in the task of phenotyping patient cohorts.[13,14] Thus, embeddings have become a popular way to represent disease, but so far, have not been used to represent medical decision making, or to relate health state to treatments.

In this work, we describe a new strategy to create embeddings that represent the relationship between medical history and the first prescription of a drug. This could include diseases that are indications for prescription of a drug, or drugs that act to ameliorate the side effects of other drugs (such as potassium supplements for diuretic prescriptions), or procedures that require prescription of a drug (such as colonoscopy preparations). By creating these embeddings, we aim to capture not just patient state, but the variation in patient state that leads to drug prescription. We describe a method for constructing such embeddings and the characteristics of these embeddings.

We assess the performance of our embeddings for identification of de facto drug indications. While databases such as UMLS[15] and MEDI[16] codify the most well-established uses for drugs, these may not reflect off-label uses and changes in medical practice. Our method provides a data-driven way to identify such prescription choices. But rather

than simply classifying a medication as used or not-used for a particular indication, we identify a richer context for medication use.

We explore the application to the important task of conducting drug safety studies using observational data. Comparative cohort studies are a major tool for assessing drug safety, but these studies require appropriate adjustment for confounding variables. Propensity score methods adjust for confounders by estimating the effect of these confounders on treatment assignment. Confronting the high dimensionality of the health record, with tens of thousands of potential confounders, recent studies have proposed improvements to the propensity score, including large-scale regularized estimation, and automatic determination of relevant confounders.[8,17–19] Fewer studies have focused on the selection of patient cohorts, a necessary step preceding the application of the propensity score. This step is usually designed by experts, limiting the throughput of comparative cohort studies. Including non-comparable patients can add bias and noise to effect estimates.[20,21] For example, Weinstein[20] showed that people with a history of gastrointestinal bleeding are unlikely to be prescribed paracetamol rather than ibuprofen, complicating efforts to estimate the effect of ibuprofen on bleeding. Coarsened Exact Matching (CEM) [22] is one way to create more comparable cohorts, but it requires experts to select a few important variables for matching, and it is not well suited to the context of sparse, high-dimensional data

To address the deficit in methods for constructing cohorts, we examine the utility of our embeddings for this purpose. The first step in constructing such cohorts typically involves picking a comparable drug to the drug of interest. We assess the use of drug embeddings to identify the most comparable drugs: that is, the drugs given in the most similar contexts. The second step in constructing cohorts is to select a control set of people taking the comparator drug who are roughly similar to the treated group. Since the embeddings represent medical events leading up to drug prescription, the patients with the most similar events should not only have similar likelihood of being prescribed the drug of interest, but they may also be prescribed that drug for the same reason. We create a matching method that combines CEM with matching on the embedding vector. We compare this form of matching to other forms of constructing matched cohorts.

Our embeddings create a representation of medical history centered on provider treatment choices. We expect that our approach will enable researchers to more rapidly conduct cohort studies and discover characteristics of medical practice.



Figure 1 A. Illustration of transformation of embeddings. B. Outline of skipgram creation and neural network to create indication embeddings

Methods

All methods described below were applied to the Marketscan IBM claims data set, which contains prescriptions, coded diagnoses and procedures, each time-stamped. We match NDC codes to drug generic names using the MarketScan RED BOOK[™] Supplement (includes variables related to drug prescription).

Generation of indication embeddings

To create the indication embeddings, we adapt the method discussed in Mikolov.[12] In that approach, each word in a sentence can be a *label*, and a randomly chosen *context* window size (number of words before and after) around the label word is chosen. Each word in the window forms a *context* for the *label*. Skip-grams, which are pairs of *(context, label)* are the input to the learning procedure, which then creates embeddings that maximize the likelihood of this data. For our application, we create a number of changes to this procedure (Figure 1B). Only new drug prescriptions are *labels*, and events (ICD-9 codes, drug prescriptions, or procedures) preceding a new drug prescription are *context* for the new drug. Therefore, unlike in Word2Vec, there is an asymmetry between context and label. In another difference, medical data, unlike text data, has the element of time: we are more interested in events that happen soon before a prescription. Mikolov. et al., implemented this idea by selecting words in randomly chosen windows up to 10 words around the training word. Our context window is based on selecting a random number of weeks before prescription, by drawing from half-normal distribution with a standard deviation of 40 weeks. As well, some patients have more data than other patients, which can result in very sick and densely observed patients dominating the distribution of skip-grams. So, when patients have multiple events in two month period, we randomly select one of these events rather

than creating a skip-gram pair for each event. As in many word2vec implementations, this sampling is weighted to downsample the most frequent codes.

Given the set of skip-grams generated above, we train a simple neural network to create embeddings that best predict the new-drug labels. We use the usual method of cross-validated hyperparameter selection to select embedding size, regularization, and learning rate, resulting in creating 50-dimensional embedding vectors, which are finally vectors L2 normalized.

Evaluation of embedding vectors for prediction of therapeutic use

We use the dot product (cosine distance) between a drug's embedding and an ICD-9 code's embedding vector as a measure of distance: close drugs and diagnoses imply similar health context. We use this to evaluate the quality of our embedding vectors. We evaluate whether the closest drugs to an ICD-9 code are the therapies associated with that ICD-9 code (and vice versa), as recorded in MEDI.[16] We use MEDI as our gold-standard for true therapeutic uses to create an ROC curve. The 10,912 reported therapeutic relationships in the MEDI High Precision Set that overlap with the drugs and ICD-9 codes in our data form our gold-standard positive examples, and 752,812 relationships between the same drugs and ICD-9 codes, where the relationship does not appear in any MEDI database, comprise our negative examples. We evaluate the performance as compared to a baseline method: co-occurrence of drug with the ICD-9 code preceding the drug. Our baseline co-occurrence method uses two-by-two contingency tables for the co-occurrence of drug and diagnosis yields the relative reporting ratio; pairs of drug and diagnosis with the highest relative reporting ratio are expected to have therapeutic relationships. We also report the performance if the ratios are adjusted using the Gamma-Poisson Shrinker method,[23] an empirical Bayes approach which has been used to mitigate the influence of low counts in these two-by-two tables.

Creating propensity scores to assess comparability of drugs

Given a set of patients treated with either a drug of interest or a comparator drug, we aim to create a propensity score that models the probability of receiving treatment, as opposed to comparator. This amounts to a large-scale logistic regression to model p(treated = 1 | medical history, demographics). We used elastic-net regression to fit the model and estimate the propensity score, with hyperparameter tuning with cross validation. In the logistic regression model, we include the following as predictors:

- Gender
- Age, year, number of prescriptions, diagnoses, and procedures preceding treatment. These are modeled using B-splines with knots placed based on the quantiles of the distribution of these values in the treated population.
- Indicator variables for presence or absence of each drug, prescription or procedure in a patient's medical history. Since more recent events would be expected to be more relevant for predicting prescription, we experimented with various features to represent these events. By assessing performance of the resulting classifiers, we settled on representing each event with an indicator variable for if the event is recorded in the previous month before prescription; in the previous year; or ever in patient history

Implementation and analysis of patient matching

We summarize patients using a weighted average of the codes that appear in their history in the time before treatment. Each code is weighted using an exponential decay so that more recent codes count more, and the weights are normalized to sum to 1. This creates a 50-dimensional health summary vector for each patient. This method was chosen because it maximized the accuracy of predicting treatment, but other methods of encoding patient history had similar results. Then we employ a two-step matching scheme that uses Coarsened Exact Matching (CEM) to bin patients, then uses Mahalanobis matching to match patients within bins based on their health summary vectors. For the CEM step we create coarsened patient bins, defined as: gender; calendar year of prescription; age in bins of: 0-5,6-10,11-15,16-25,26-40,41-55,55-70; and number of unique drugs prescribed to that person (divided into percentiles). The CEM has a dual purpose: it ensures we are not creating entirely inappropriate matches (ie, matching children to seniors); and it reduces the number of possible matches for each treated person, making the Mahalanobis matching against propensity score matching, we use the propensity scores calculated as in the previous section, and perform nearest neighbor matching on the propensity score values.

Then, we analyze how similar matched pairs are in terms of smoking. There is an ICD-9 code for smoking, but we expect that it only represents some of the information related to smoking status. Therefore, we estimate the probability of smoking by fitting a simple logistic regression model trained to predict whether a person has an ICD-9 code for smoking (code 305.1), given the person's entire health history, as described in the previous section (removing ICD-9)

code 305.1). The top predictors, as expected, include codes for lung diseases and smoking cessation therapies. We use each person's predicted probability of smoking to evaluate the matching. Note that the same information (including the code for smoking) is provided for training the propensity score function; however this does not guarantee that propensity score matching will match people for probability of smoking. Then, we calculate the correlation between paired patients from the two matching schemes.

Results

Overview of indication embeddings

We apply the method outlined in Figure 1 to create a skip-gram set and to learn embeddings. This method yields two types of embeddings: indication embeddings, which signify the health context for a treatment choice; and new-drug (incident drug) embeddings, which represent the treatment choices given in those contexts. The indication embeddings create a unique vector for each possible event in the health history, comprising 3,014 common drugs, 13,090 ICD-9 codes, and 14,187 CPT codes. Figure 2A shows the UMAP[24] placement of all such events, which puts codes with the most similar indication embeddings near each other. Some selected medical events are highlighted. For example, myocardial infarction (MI), medications related to MI, and procedure codes for coronary bypass graft, are all located near each other. Similarly, in a visualization of the new-drug embeddings, drugs from the same REDBOOK therapeutic group appear near each other. In Figure 2B, we group codes into Clinical Classification Software category, for diagnoses, and REDBOOK therapeutic group, for drugs, and then calculate the average indication embedding distance between codes in a pair of groups. For example, at the top, gastrointestinal drugs are near the gastrointestinal

bioRxiv preprint doi: https://doi.org/10.1101/737049. The copyright holder for this preprint (which was not peer-reviewed) is the author/funder. It is made available under a CC-BY 4.0 International license.



Figure 2: Indication embedding overview. A. Each point is one event (prescription, diagnosis code, or procedure code), visualized with UMAP. Selected related sets of events are highlighted. Circles are ICD-9 codes, + symbols are medications, and triangles are procedures. B. For each pair of categories of disease codes (black labels on rows) and drug codes (red labels), we show the average distance between codes (color scale).

Using indication embeddings to predict drug therapeutic uses

Our embeddings were created to summarize the relationship between events in the health history, and treatment choices. A slightly different, but related, task, is prediction of the FDA approved uses of drugs. Similar to Li and Xiao,[6] we use MEDI as our gold standard, as it is based on human curated databases. We examine whether embedding distance between a drug and an ICD-9 code is a good predictor of whether that drug and ICD-9 code have an indication relationship. For comparison, we use the relative reporting ratio: how much more frequently an ICD-9 code is observed before drug prescription, versus the rate expected, and a stabilized version known as the Gamma-Poisson Shrinker. The AUC is .82 for the embeddings, and .80 for both of the disproportionality methods; this is a similar improvement in ROC AUC over disproportionality methods to what Li and Xiao were able to achieve with



their method. Figure 3 shows the variation in ability of embeddings to predict MEDI relationships per drug therapeutic group, and per CCS class.

Figure 3 We calculate an area under the ROC for each drug, and for each ICD-9 code. Left: Distribution of these ROC values per drug category. Right: per ICD-9 group

Although the results are promising, more interesting is the ability to detect subtleties of the context in which drugs are prescribed. Figure 4 shows the variation in the embedding distance between

antidepressants and their most associated diagnoses. Antidepressants include drugs with a range of mechanisms of action, and they are prescribed for diverse reasons. Most have side effects that can influence prescription choice. For instance, tricyclic antidepressants (TCAs) are no longer the standard of care as a first-line treatment for depression (while selective serotonin reuptake inhibitors (SSRI), serotonin norepinephrine reuptake inhibitors (SNRI), mirtazapine, and bupropion are preferred[25]). TCAs can be a second-line depression treatment, but another common use is for chronic pain.[26] This is reflected by the closer embeddings distance of these antidepressants to neuropathy and spine diseases, as well as Irritable Bowel Syndrome (IBS). Pediatric neuropsychiatric diseases are closest to fluvoxamine, while mirtazapine is closest among antidepressants to neurodegenerative diseases of the elderly, which are often associated with depression. Indeed, mirtazapine is a widely used in populations with dementia.[27] This



Figure 4 The dot-products between drug vectors (row) and groups of diagnosis codes (columns) for the diagnoses most related to the antidepressants. Antidepressants are colored by class: TCA are blue; SSRI red; SNRI green; other black.

New-drug embeddings point to appropriate comparator drugs

We now turn to the application of the embeddings for drug safety studies, in particular, cohort studies. These study designs typically compare outcomes in the pool of new users of a drug of interest, against the people who took a comparator drug. Selecting an appropriate comparator drug is the first step toward removing differences between the exposed and comparator cohorts that can confound the results of cohort studies. Typically, a medical expert selects the best comparator, but here we assess the performance of embeddings for this task. Drugs that are most comparable should be given in the most similar medical contexts, which is exactly the information that we have shown above is captured in our embeddings. Therefore, we evaluate whether drugs with closer new-drug embeddings are more comparable. Comparability of the treated cohort and the comparator cohort can be assessed using the performance of a classifier to predict whether each person in the cohort study is on the treatment, versus the comparator drug.[28] The easier it is to distinguish these two populations, the less comparable they are. Thus we fit the propensity score, *p(treated | medical history)*, which summarizes the probability of each person getting the treatment of interest. Propensity score models for comparable populations have a low area under the ROC curve (AUC), while less comparable populations are easily separated and have a high AUC. Figure 5 shows that similar embeddings indeed point to more appropriate comparator drugs for a given treatment, even in the complex case of neuropsychiatric drugs. Our embeddings, then, can be used to suggest the most appropriate comparator drug among a number of choices, and



Figure 5 Anticonvulsants (red) and antipsychotics (black) have similar drug embedding vectors, indicating similar medical contexts. We compare the embedding similarity (x-axes) against the AUC of the propensity scores (y-axes) for finding comparator drugs for olanzapine, left, an antipsychotic, and carbamazepine, right, an anticonvulsant.

they even have the potential to allow automatic selection of comparators.

Using embeddings to match patients in a cohort study

Choosing appropriate comparator drugs is only the first step in removing confounding differences between cohorts. This is typically followed by further adjustments, such as propensity score matching or propensity score weighting. Weighting can result in extreme values if very incomparable patients are included. Matching patients again forces researchers to confront the high dimensionality of the health care record. Exact matching of health histories is impossible. One alternative is coarsened exact matching,[22] but that method requires that researchers choose only a few important variables to match on, and it is not easily extended to the high-dimensional setting. In contrast, propensity score matching reduces the high-dimensional health history to a single dimension: the propensity for treatment, which is the probability of treatment given health history. This reduction can result in loss of valuable information regarding patient state, resulting in matching dissimilar people.[29] Mahalanobis matching, like coarsened exact matching, does not extend to the setting of sparse uninformatively coded data. Therefore, we experiment with a new matching method. We summarize each patient's health status as a weighted average of their embedding vectors. Then, we perform a two-step matching strategy that first, matches patients on coarsened versions of age, year, and number of prescriptions, as well as gender, then, within these matches, performs Mahalanobis matching on the lower dimensional *health summary vectors*.

An example of the difficulty of patient matching in observational cohort studies can be found in a comparison of two common antidepressants: bupropion and trazadone. In a study estimating the effect of bupropion on some outcomes, each person taking bupropion could be matched to a person taking trazodone using the propensity score. Each drug has a side effect risks and contraindications that can influence treatment choice and thus propensity for treatment assignment. In addition, bupropion is approved not only for depressive disorders, but also for smoking cessation. Since smoking can cause a number of downstream effects on health, it is desirable to match bupropion-users who are smokers to trazodone users who are smokers. In Figure 6 we show the UMAP visualization of the health summary vectors for a set of people on bupropion versus trazodone, within one coarsened bin (age. We also include people taking varenicline, a smoking cessation drug, for comparison. Since varenicline users presumably smoke, people nearer to the varenicline population should be more likely to be smokers. Naturally, the Mahalanobis matching on

health summary vectors (Figure 6B) results in points from closer parts of the plot being matched. To see whether this results in smokers being matched other smokers, we create a simple score to summarize smoking status (see Methods) and calculate how similar this score is between matched patient pairs, using Spearman correlation. The correlation is 0.57 for the Mahalanobis matched pairs, and 0.34 for the propensity matched pairs. This shows that Mahalanobis matching on health summary vectors improves matching on specific health risk factors. Such a matching can create more comparable cohorts, which can then be used with traditional methods in cohort studies to estimate the effects of drugs.[30]



Figure 6: Each point is one person on bupropion, trazodone, or varenicline, using a UMAP visualization of people's health summary vectors. Varenicline-takers cluster in the upper left, so people who are on the bottom/right probably are not in need of smoking cessation. The same points are shown in A and B. A. Lines connect 57 pairs matched by propensity score matching. B. Lines connect 61 pairs using Mahalanobis matching on health summary vectors.

Discussion and Conclusion

We have shown that indication embeddings form representations of medical events that reflect the complexity of health care. While the indication embeddings work well for predicting known drug indications, their true utility is in their ability to summarize the health status implied by presence of a given code in a medical history. The decision to prescribe a medication is a complex consideration of personal history, tolerability of side effects, and variation in severity of disease, making such a representation desirable. An important application of these embeddings is to accelerate discovery in drug safety cohort studies. Cohort study design currently heavily relies on clinician knowledge, but since these embeddings are trained on data that results from similar knowledge, they are a way to describe these

considerations in a computable form. We show how the embeddings can identify comparable drugs and suggest more comparable matched populations.

Our approach has a number of limitations related to the design choices. In our embedding construction, each medical event occurring in a patient's history is treated as unrelated to other medical events; an alternative design would instead combine patient history in order to predict the drug, given a combination of medical events. We chose our approach specifically to untangle the relationship between each event and treatment state. Other limitations include the drawbacks of claims data, which is not ideal for research purposes: some diseases have better coverage in terms of codes available than others. As well, we did not create embeddings for ICD-10 codes. The same method can be easily implemented for ICD-10 codes, but as they are only used for data in since 2015, they are not yet as valuable for use in retrospective cohort studies.

These embeddings available for public are reuse at https://figshare.com/projects/Using indication embeddings to represent patient health for drug safety studies/6 7532 and have applications beyond the ones described in this paper. Previous embeddings have focused on modeling disease onset but our focus on modeling health state leading to drug prescription makes these results unique. While we have focused our analysis of the embeddings on the relationships between drugs and diseases, these embeddings also relate to symptoms induced by diseases and procedures associated with drugs. Therefore, our embeddings could be used to discover these associations. As Marketscan aggregates data from many health systems, it should represent the standard of care across the USA. Then, the indication embeddings can also be incorporated into analyses on other USA health data sets which have the same standard codes, but smaller patient sizes, as a way to share information learned on a large national data set. In another area of future work, comparing the embeddings trained on medical data from different health systems will allow comparison of medical decisions. Future health data will contain genetic information, text, and image data, and these can be incorporated into the same framework. Other types of embeddings may be of interest for future work, such as embeddings trained on the time after drug prescription, which may indicate side effects of drugs.

Acknowledgements

This work was funded by NIH 5K01ES028055.

bioRxiv preprint doi: https://doi.org/10.1101/737049. The copyright holder for this preprint (which was not peer-reviewed) is the author/funder. It is made available under a CC-BY 4.0 International license.

References

1. Hripcsak G, Ryan PB, Duke JD *et al.* Characterizing treatment pathways at scale using the OHDSI network. *Proceedings of the National Academy of Sciences* 2016;**113**:7329–36.

2. Melamed RD, Rzhetsky A. Patchwork of contrasting medication cultures across the USA. *Nature Communications* 2018;9, DOI: 10.1038/s41467-018-06205-1.

3. Rotmensch M, Halpern Y, Tlimat A *et al.* Learning a Health Knowledge Graph from Electronic Medical Records. *Sci Rep* 2017;7, DOI: 10.1038/s41598-017-05778-z.

4. Li L, Ruau DJ, Patel CJ *et al.* Disease risk factors identified through shared genetic architecture and electronic medical records. *Science translational medicine* 2014;**6**:234ra57-234ra57.

5. Jung K, LePendu P, Chen WS et al. Automated Detection of Off-Label Drug Use. PLOS ONE 2014;9:e89324.

6. Li Y, Xiao C. Developing a Data-driven Medication Indication Knowledge Base using a Large Scale Medical Claims Database. *AMIA Summits on Translational Science Proceedings*.

7. Hernán MA, Alonso A, Logan R *et al.* Observational Studies Analyzed Like Randomized Experiments: An Application to Postmenopausal Hormone Therapy and Coronary Heart Disease. *Epidemiology* 2008;**19**:766–79.

8. Schneeweiss S, Rassen JA, Glynn RJ *et al.* High-dimensional Propensity Score Adjustment in Studies of Treatment Effects Using Health Care Claims Data. *Epidemiology* 2009;**20**:512–22.

9. Ryan PB, Madigan D, Stang PE et al. Medication-Wide Association Studies. CPT: Pharmacometrics & Systems Pharmacology 2013;2:e76–e76.

10. Miotto R, Li L, Kidd BA *et al.* Deep Patient: An Unsupervised Representation to Predict the Future of Patients from the Electronic Health Records. *Scientific Reports* 2016;6:26094.

11. Choi Y, Chiu CY-I, Sontag D. Learning Low-Dimensional Representations of Medical Concepts. *AMIA Jt Summits Transl Sci Proc* 2016;**2016**:41–50.

12. Mikolov T, Chen K, Corrado G et al. Efficient Estimation of Word Representations in Vector Space. arXiv:13013781 [cs] 2013.

13. Bai T, Chanda AK, Egleston BL *et al.* EHR phenotyping via jointly embedding medical concepts and words into a unified vector space. *BMC Medical Informatics and Decision Making* 2018;**18**:123.

14. Glicksberg BS, Miotto R, Johnson KW *et al.* Automated disease cohort selection using word embeddings from Electronic Health Records. *Pac Symp Biocomput* 2018;**23**:145–56.

15. Bodenreider O. The Unified Medical Language System (UMLS): integrating biomedical terminology. *Nucleic Acids Res* 2004;**32**:D267-270.

16. Wei W-Q, Cronin RM, Xu H *et al.* Development and evaluation of an ensemble resource linking medications to their indications. *Journal of the American Medical Informatics Association : JAMIA* 2013;**20**:954–61.

17. Brookhart MA, Schneeweiss S, Rothman KJ *et al.* Variable selection for propensity score models. *Am J Epidemiol* 2006;**163**:1149–56.

18. Tian Y, Schuemie MJ, Suchard MA. Evaluating large-scale propensity score performance through real-world and synthetic data experiments. *International Journal of Epidemiology* 2018, DOI: 10.1093/ije/dyy120.

19. Mittal S, Madigan D, Burd RS *et al*. High-dimensional, massive sample-size Cox proportional hazards regression for survival analysis. *Biostatistics* 2014;**15**:207–21.

20. Weinstein RB, Ryan P, Berlin JA *et al.* Channeling in the Use of Nonprescription Paracetamol and Ibuprofen in an Electronic Medical Records Database: Evidence and Implications. *Drug Saf* 2017:1–14.

21. Pearl J. On a Class of Bias-Amplifying Variables that Endanger Effect Estimates. *arXiv:12033503 [cs, stat]* 2012.

22. Iacus SM, King G, Porro G. Causal inference without balance checking: Coarsened exact matching. *Political Analysis* 2012;**20**:1–24.

23. DuMouchel W. Bayesian Data Mining in Large Frequency Tables, with an Application to the FDA Spontaneous Reporting System. *The American Statistician* 1999;**53**:177–177.

24. McInnes L, Healy J, Melville J. UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction. *arXiv:180203426 [cs, stat]* 2018.

25. Kennedy SH. A Review of Antidepressant Therapy in Primary Care: Current Practices and Future Directions. *Prim Care Companion CNS Disord* 2013;**15**, DOI: 10.4088/PCC.12r01420.

26. Antidepressant Pharmacotherapy: Considerations for the Pain Clinician - Jackson - 2003 - Pain Practice - Wiley Online Library.

27. Bannerjee S, Hellier J, Romeo R *et al.* Study of the use of antidepressants for depression in dementia: the HTA - SADD trial - a multicentre, randomised, double-blind, placebo-controlled trial of the clinical effectiveness and cost-effectiveness of sertraline and mirtazapine. *Health Technology Assessment* 2013;**17**, DOI: 10.3310/hta17070.

28. Walker AM, Patrick AR, Lauer MS *et al.* A tool for assessing the feasibility of comparative effectiveness research. *Comparative Effectiveness Research* 2013, DOI: 10.2147/CER.S40357.

29. King G, Nielsen R. Why propensity scores should not be used for matching. *Copy at http://j mp/1sexgVw Download Citation BibTex Tagged XML Download Paper* 2016;**378**.

30. Ho DE, Imai K, King G *et al.* Matching as Nonparametric Preprocessing for Reducing Model Dependence in Parametric Causal Inference. *Political Analysis* 2007;**15**:199–236.

Figure legends

Fig. 1: A. Illustration of transformation of embeddings. B. Outline of skipgram creation and neural network to create

indication embeddings.

Fig 2: Indication embedding overview. A. Each point is one event (prescription, diagnosis code, or procedure code),

visualized with UMAP. Selected related sets of events are highlighted. Circles are ICD-9 codes, + symbols are

medications, and triangles are procedures. B. For each pair of categories of disease codes (black labels on rows) and drug codes (red labels), we show the average distance between codes (color scale).

Fig 3: We calculate an area under the ROC for each drug, and for each ICD-9 code. Left: Distribution of these ROC values per drug category. Right: per ICD-9 group.

Fig 4: The dot-products between drug vectors (row) and groups of diagnosis codes (columns) for the diagnoses most related to the antidepressants. Antidepressants are colored by class: TCA are blue; SSRI red; SNRI green; other black.

Fig 5: Anticonvulsants (red) and antipsychotics (black) have similar drug embedding vectors, indicating similar medical contexts. We compare the embedding similarity (x-axes) against the AUC of the propensity scores (y-axes) for finding comparator drugs for olanzapine, left, an antipsychotic, and carbamazepine, right, an anticonvulsant.

Fig 6: Each point is one person on bupropion, trazodone, or varenicline, using a UMAP visualization of people's health summary vectors. Varenicline-takers cluster in the upper left, so people who are on the bottom/right probably are not in need of smoking cessation. The same points are shown in A and B. A. Lines connect 57 pairs matched by propensity score matching. B. Lines connect 61 pairs using Mahalanobis matching on health summary vectors.