

When deep learning meets causal inference: a computational framework for drug repurposing from real-world data

Ruoqi Liu¹, Lai Wei², and Ping Zhang^{1,2,*}

¹Department of Computer Science and Engineering. The Ohio State University, 2015 Neil Ave, Columbus OH 43210, USA.

²Department of Biomedical Informatics. The Ohio State University, 1800 Cannon Drive, Columbus OH 43210, USA.

*e-mail: zhang.10631@osu.edu

ABSTRACT

Drug repurposing is an effective strategy to identify new uses for existing drugs, providing the quickest possible transition from bench to bedside. Existing methods for drug repurposing that mainly focus on pre-clinical information may exist translational issues when applied to human beings. Real world data (RWD), such as electronic health records and insurance claims, provide information on large cohorts of users for many drugs. Here we present an efficient and easily-customized framework for generating and testing multiple candidates for drug repurposing using a retrospective analysis of RWDs. Building upon well-established causal inference and deep learning methods, our framework emulates randomized clinical trials for drugs present in a large-scale medical claims database. We demonstrate our framework in a case study of coronary artery disease (CAD) by evaluating the effect of 55 repurposing drug candidates on various disease outcomes. We achieve 6 drug candidates that significantly improve the CAD outcomes but not have been indicated for treating CAD, paving the way for drug repurposing.

Introduction

Drug repurposing (a.k.a., drug repositioning) is a strategy to accelerate the drug discovery process by identifying novel uses for existing approved drugs [1]. The primary advantage of drug repurposing over traditional drug development is that it starts from compounds with well-characterized pharmacology and safety profiles and can significantly reduce the risk of adverse effects and attrition in clinical phases [2].

While many successful repurposed drugs (e.g., Viagra for erectile dysfunction) have been discovered serendipitously [3], computation-based repurposing methods are developed recently by leveraging structural features of compounds or proteins [4, 5], genome-wide association study (GWAS) [6], transcriptional responses [7], and gene expression [8]. These methods focus primarily on using pre-clinical information. Unfortunately, the clinical therapeutic effects in humans are not always consistent with pre-clinical outcomes [9].

In healthcare, real world data (RWD) [10] refers to longitudinal observational data derived from sources that are associated with outcomes in a heterogeneous patient population in real-world settings, such as patient surveys, electronic health records (EHRs), and claims and billing activities. Since RWDs are direct observations from human bodies, they become a promising source for drug repurposing. Few researchers have already validated a small number of repurposing drug candidates on RWD [11, 12]. However, there are some limitations with these approaches. First, most studies are complementary (i.e., the original hypotheses usually come from other studies). Second, their studied number of repurposing candidates is limited and unable to proactively generate *de novo* repurposing drug candidates.

In this study, we follow protocols of randomized clinical trial (RCT) design [13], and computationally screen repurposing candidates for beneficial effect by explicitly emulating the corresponding clinical trials using RWDs. Considering the inherent characteristics of RWD (i.e., temporal sequence data and existing confounding variables [14]), we apply deep learning and causal inference methodologies to control the confounders in RWD, and systematically estimate the drug effects on various disease outcomes. Specifically, the estimated drug effects are obtained by long short-term memory (LSTM) [15] and inverse probability of treatment weighting (IPTW) [16], on MarketScan claims data [17].

As a test case, we apply the proposed drug repurposing framework to coronary artery disease (CAD) cohorts of millions of patients and emulate RCTs for multiple drug candidates, estimating their effects on CAD progression outcomes.

In general, our contributions are three folds:

- We develop a framework for high-throughput screening of on-marked drugs by emulating, for each drug, an RCT that evaluates its beneficial effect. The repurposed drug candidates can be proactively generated on existing large-scale RWDs.
- We present an innovative study design for the estimation of the drug's effect from longitudinal observational data. The study CAD cohorts are automatically derived under our framework, which accelerates the process of computational drug repurposing.
- We propose a deep learning based propensity score estimation model to correct for confounding and selection

biases. Experimental comparisons to the logistic regression based propensity score estimation model show that our proposed deep learning model effectively estimate drug effects from RWDs, paving the way for drug repurposing.

Overall framework

We develop a high throughput computational drug repurposing pipeline (Fig. 1) that, given a disease cohort (i.e., CAD patients) extracts a list of potential repurposing drug ingredients and, for each, identifies the corresponding user and non-user sub-cohorts. It then computes, for all patients in both sub-cohorts, a large number of features (confounding factors), as well as the disease progression outcomes. The treatment effects are estimated after correcting for confounding and selection biases using the deep learning framework (Fig. 2). Here, the proposed framework is equipped with attention mechanism that provides the interpretability of the model. The drug ingredients with beneficial effect and statistical significance will be considered as repurposed drug candidates and suggested to be used for treating CAD. Algorithm 1 overviews the steps of estimating the effect of assigned treatment on the outcome from observational data.

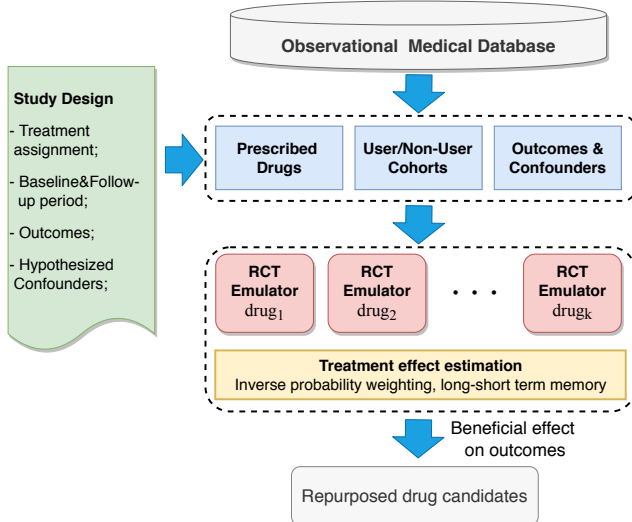


Figure 1. Flowchart of overall drug repurposing framework.

Results

In this section, we first introduce the dataset we use for this study. Then we demonstrate the performance of our model in CAD drug repurposing experiments. We identify more than 90 million patients in MarketScan [17] data from 2012 to 2017, which contain individual-level, de-identified healthcare claims information from employers, health plans, hospitals, Medicare, and Medicaid programs. MarketScan claims data is primarily used for evaluating health utilization and services. The overall patients' distribution of the recording period is

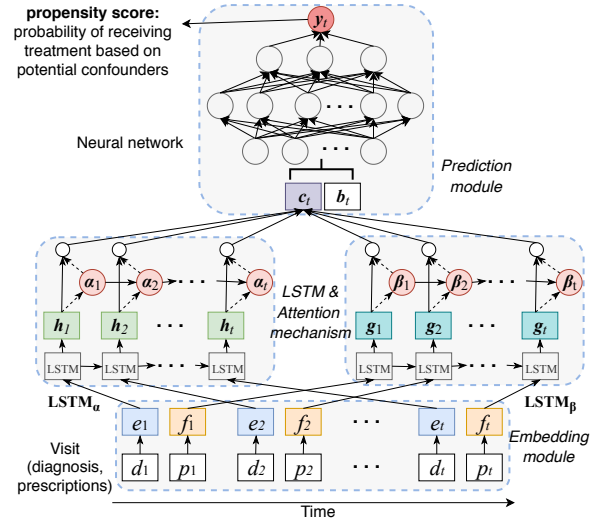


Figure 2. Illustration of the deep learning model for predicting treatment probability (a.k.a., propensity score) that used for correcting confounding from temporal time sequence data (including diagnoses d_t , prescriptions p_t and demographics b_t). It consists of three main components: embedding module, recurrent neural network and prediction module.

shown in Fig. 3(a). We take both inpatient and outpatient claims into consideration. CAD cohort criteria are defined using International Classification of Diseases (ICD) codes [18] (Supplementary Table 1 for definitions). In total, there are 1,178,997 CAD patients. We refer to the first date when patients were determined to have CAD as *CAD initiation date*. Figure 3(b) shows the patients distribution of time before/after *CAD initiation date*.

Dataset

We identify three categories of study variables: demographic characteristics, diagnosis codes and prescription medication. Demographic characteristics in MarketScan CAD data include information on age and gender for each patient. Figure 3(d) shows the age and gender statistics and distribution of our dataset. Because a majority of data come from commercial claims, race and ethnicity information is incomplete and is not included in the analysis. Diagnosis codes in MarketScan CAD data are defined using the ICD codes for billing purposes. There are 57,089 ICD-9/10 codes considered in the dataset. Prescription medications in MarketScan CAD data also contain all prescription drug claims which contain prescription drug name (generic and brand), national drug code (NDC), and the number of days of supply approved. By matching NDC to observational medical outcomes partnership (OMOP) ingredient concept ID [19], we get 1,353 unique drugs in the dataset for drug repositioning screening. For drugs with multiple ingredients, we consider each active ingredient separately in the mapping processes.

Algorithm 1 The algorithmic framework to estimate effect of assigned treatment on the outcomes

Input: patient data: assigned treatment, outcomes, values for potential confounders

Output: repurposed drug candidates, and their estimated effect, unbalanced feature ratio and significance

- 1: Generate user and non-user sub-cohorts for the treatment
- 2: Compute balancing weights for all patients in both sub-cohorts via LSTM based IPTW
- 3: Estimate the effect over multiple outcomes after correcting for the biases in the confounders (Eq. (1))
- 4: Compute the unbalanced feature ratio for the treatment after re-weighting using standardized difference (Eq. (2))
- 5: Estimate the significance of effect and compute adjusted p-value using bootstrapping
- 6: **if** estimated effect < 0 and adjusted p-value < 0.05 and unbalanced feature ratio $< 2\%$ **then**
- 7: **return** the estimated effect, unbalanced feature ratio and computed p-value
- 8: **end if**

To evaluate the drug effect, we define a set of clinically-relevant events linked the CAD as the disease outcomes (e.g., heart failure onset and stroke onset) after consulting domain experts. The definition is based on the ICD codes and can be found in Table 2 and Table 3 of Supplementary Materials. Since CAD is the major risk factor for both heart failure [20] and stroke [21], we hypothesize that an effective drug will lower the risks of CAD patients develop those diseases. Figure 3(c) demonstrates the time to develop outcomes from the CAD *initiation date*. The confounding variables affect both treatment assignment of patients and an outcome used in the trial. We consult domain experts to compile a list of hypothesized confounders for the CAD case study with respect to the study variables illustrate above: demographics, co-morbidities (diagnosis codes) and co-prescribed drugs.

Model performance

Evaluation metrics

Treatment effect estimation In this study, we leverage *average treatment effect* (ATE) to examine the treatment effect at the population level, which is defined as

$$ATE = \mathbb{E}(Y_1) - \mathbb{E}(Y_0) \quad (1)$$

where $\mathbb{E}(Y_1)$ and $\mathbb{E}(Y_0)$ are the expected potential treated and control outcome of the whole population respectively. The values of ATE are used to determine whether the given treatment can improve disease outcomes or not.

Testing feature balance We evaluate the performance of models by measuring features' balance between the weighted user and non-user sub-cohorts generated by the IPTW. Given patient weights from IPTW, we quantify the balance for each feature using its standardized mean difference (SMD), which

is the difference in the variable means between the two treatment groups, divided by the combined standard deviation. To be exact, we use the following definition of the standardized difference,

$$SMD = \frac{|\mu_{user} - \mu_{nonuser}|}{\sqrt{(s_{user}^2 + s_{nonuser}^2)/2}} \quad (2)$$

where μ_{user} and $\mu_{nonuser}$ are the mean in user cohort and nonuser cohort; s_{user}^2 and $s_{nonuser}^2$ are sample variance of variables in two sub-cohorts. For binary variables, the variance s^2 is calculated by $\mu(1 - \mu)$. We consider a standardized difference greater than 0.1 as unbalanced [22] and compute the unbalanced feature ratio (i.e., $\frac{\#unbalanced\ feature}{\#all\ features}$) before/after weighting to evaluate the performance of balancing. The user and non-user sub-cohorts are considered as balanced if their unbalanced feature ratio is below 2% after weighting.

Confidence intervals and significance of effect We use the bootstrapping method [23] to calculate the confidence intervals of estimators of $\mathbb{E}(Y_1)$ and $\mathbb{E}(Y_0)$, and statistical significance of ATE. For each candidate ingredient, we repeatedly generate multiple different control drugs via random sampling with replacement, and the analysis is repeated in each bootstrap sample. The 95% confidence interval is then computed by using the standard normal approximation: ± 1.96 times the estimate of the standard error. The p-value of the effect estimator can be computed by the normal cumulative distribution function of estimators. We further use adjusted p-value [24] as a statistically significant measurement. We consider a repurposing drug candidate as significant if its adjusted p-value is below 0.05.

Performance over repurposing drug candidates

We identify 55 ingredients as drug repurposing candidates following the study design (see Methods). Then we estimate the treatment effect on various disease outcomes (i.e., heart failure and stroke) and demonstrate the distribution of ATE in Fig. 4. Here, we only show the drug candidates with the balanced user and non-user sub-cohorts after re-weighting and statistically significant estimates (adjusted p-value). All the drugs are ranked from the left side to the right side according to the increasing order of estimated ATE values. Based on the definition of ATE (i.e., the weighted average of observed outcomes from the user and non-user sub-cohorts), the drug ingredients with ATE values that smaller than 0 are identified to improve the disease outcomes, while the drug ingredients with ATE values larger than 0 are identified to worsen the disease outcomes. For the drugs with beneficial effects, we color those with known CAD indication in red and those without known CAD indication in blue (The drug label information is collected from SIDER [25] database and DrugBank [26]).

From the results, we observe that 9 drugs yield a beneficial effect on disease outcomes among 16 selected significant drug candidates. Specifically, only 3 drugs have been used as CAD indication according to their drug labels information. The remaining 6 drugs which have not been indicated

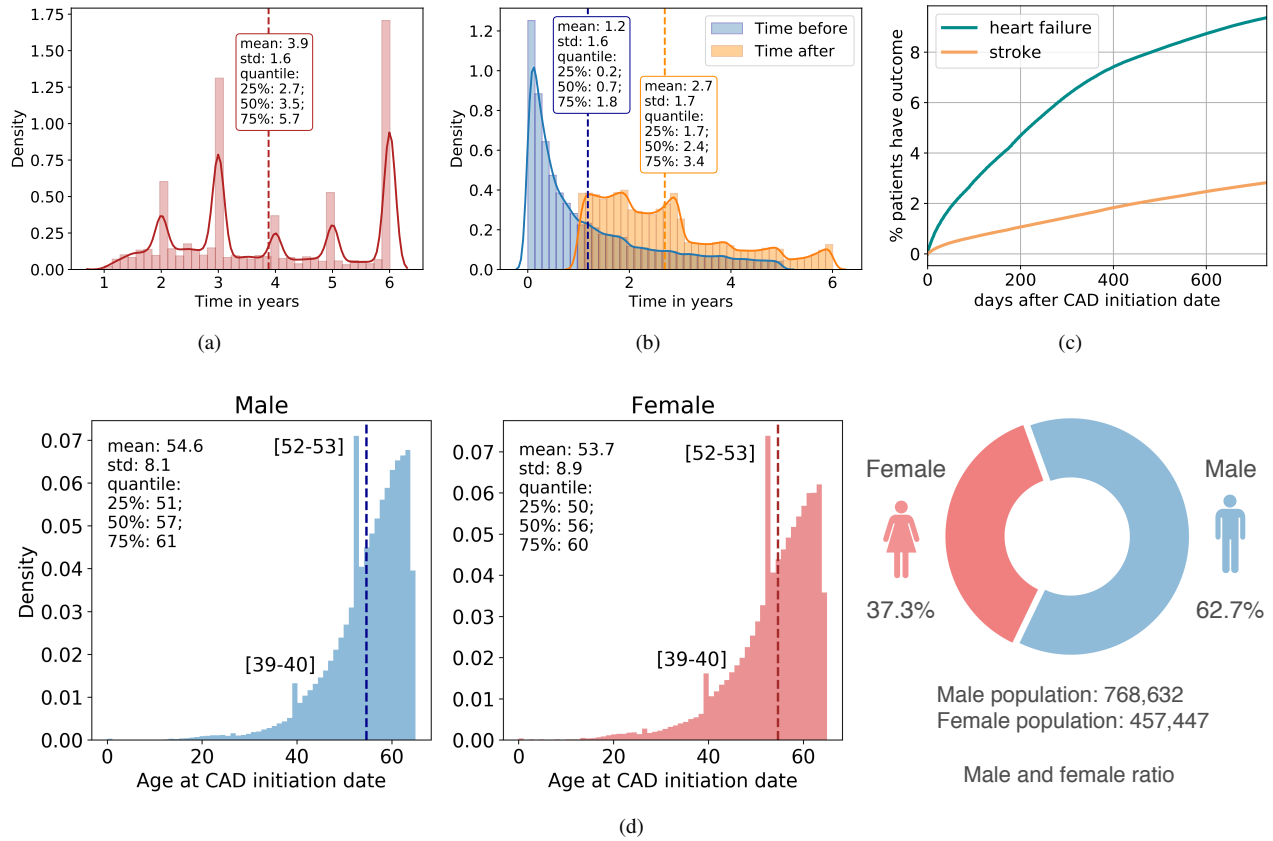


Figure 3. CAD cohorts characteristics. Figure 3(a) shows the patients’ distribution of total time in the database. Figure 3(b) shows the patients distribution of time before/after CAD initiation date. Figure 3(c) shows the growth of the number of patients develop outcomes after CAD initiation date. Figure 3(d) shows the gender distribution with age at CAD initiation date.

for treating CAD but can improve the disease outcomes are considered as repurposed drug candidates. We find evidence support for these 6 drug candidates from related literature and web resources as follows: (1) Metoprolol is one of the most commonly used beta-blockers for treating high blood pressure and chest pain. It shows beneficial effects in patients with heart failure associated with CAD [27]; (2) Fenofibrate is mainly used to treat abnormal blood lipid levels and also appears to decrease the risk of CAD in patients with diabetes mellitus [28]; (3) Hydrochlorothiazide which is often used to treat high blood pressure and diabetes insipidus [29], has already completed phase 4 trials for CAD treatment [30]; (4) Pravastatin has also been studied to have a beneficial effect on CAD [31]; (5) For simvastatin, results from randomized clinical trials show that it can reduce the occurrence of heart failure in patients with CAD [32]; (6) Valsartan, a kind of angiotensin receptor blocker, results in improved coronary micro-vascular flow reserve, suggesting direct beneficial in hypertensive patients with stable CAD [33].

We further list the sub-cohort size, feature balancing and estimated ATE values for each drug candidate in Table. 1. The results of all 55 drugs can be found in Table 4 of Supplementary Materials. The first column lists the drug names

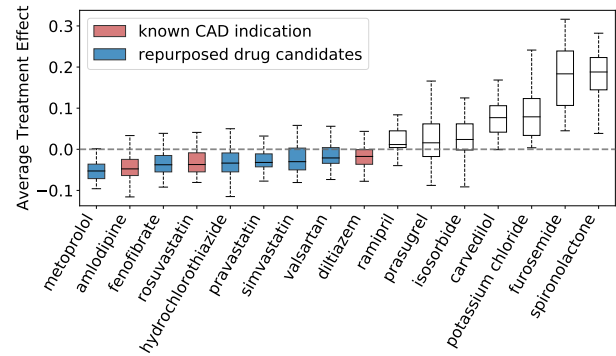


Figure 4. Distribution of estimated ATE of drugs on defined outcomes across the 50 bootstrap samples. All these showing drugs satisfy two conditions: adjusted p value ≤ 0.05 and post unbalanced ratio $\leq 2\%$. Within the boxplot, the central line denotes the median, and the bottom and the top edges denote the 25th(Q1) and 75th(Q3) and percentiles respectively. The whiskers extend to 1.5 times the interquartile range.

corresponding to drugs in Fig. 4. The second and third columns denote the number of patients in user and non-user

Table 1. Sub-cohorts size, feature balancing and estimated effects for CAD over balanced and statistically significant drug ingredients. Bold in the table denotes the ingredients *without* known CAD indication (repurposed drug candidates).

# drug_name	# user	# non-user	pre.unbalanced covariates	post.unbalanced covariates	# covariates	post.unbalanced.ratio %	pre.ATE	post.ATE
metoprolol	9730	29190	38.308	23.231	1270	1.8	-0.023	-0.043
fenofibrate	1352	4056	39.340	13.200	1038	1.3	-0.051	-0.038
rosuvastatin	2420	7260	24.020	9.620	1097	0.9	-0.063	-0.030
hydrochlorothiazide	2001	6003	32.500	15.320	1076	1.4	-0.055	-0.029
amlodipine	4613	13839	21.340	8.300	1180	0.7	-0.050	-0.026
pravastatin	2007	6021	11.260	9.640	1085	0.9	-0.016	-0.022
simvastatin	1605	4815	10.060	13.240	1044	1.3	-0.032	-0.020
valsartan	1316	3948	24.940	13.740	1026	1.3	0.010	-0.015
diltiazem	1044	3132	28.360	13.080	1007	1.3	-0.010	-0.013
isosorbide	1482	4446	33.320	9.560	1039	0.9	0.045	0.034
prasugrel	1316	3948	41.500	18.340	1019	1.8	-0.043	0.036
ramipril	887	2661	25.340	14.840	973	1.5	0.020	0.043
potassium chloride	1110	3330	43.460	20.240	1016	2.0	0.169	0.090
carvedilol	3959	11877	38.280	8.140	1154	0.7	0.198	0.124
furosemide	1545	4635	50.880	17.080	1064	1.6	0.301	0.179
spironolactone	1292	3876	70.620	12.920	1034	1.3	0.393	0.190

sub-cohorts, respectively. The next two columns denote the average number of unbalanced covariates before and after re-weighting. The "post unbalanced ratio" column represents the percentage of unbalanced covariates after re-weighting (i.e., the number of unbalanced covariates divided by the total number of covariates). And the last two columns are the estimated ATE before and after re-weighting. We rank the drugs by increasing of re-weighted ATE values. We see that our proposed method successfully corrects for most biases in the original data which results in a decrease in the number of unbalanced covariates.

Case studies: attention visualization

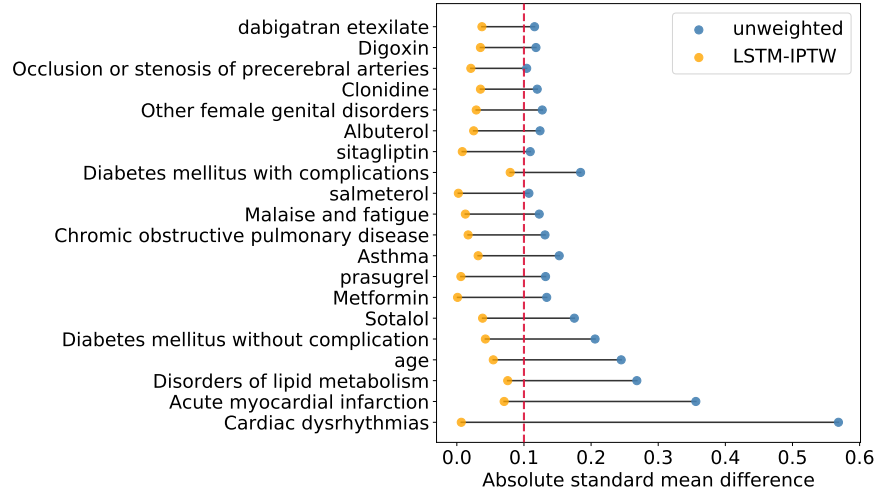
Having presented that our model successfully identified repurposed drug candidates for CAD treatment, we further demonstrate the interpretability of our framework achieves via attention mechanism. To exemplify this, we select two case drug candidates: diltiazem and fenofibrate. According to Table. 1, diltiazem and fenofibrate both have beneficial effect on CAD disease outcomes. Diltiazem has already been used for treating CAD [34], while fenofibrate does not have CAD indication on its drug label.

We want to identify the covariates that significantly biased between the user and non-user cohorts in original data but balanced after re-weighting. The learned attention weights enable visualization of each covariate and its SMD values before/after balancing between user and non-user cohorts. We select the top 20 well-balanced (i.e., large deviations of SMD during balancing) covariates and plot the distribution of SMD values for two case drugs in Fig. 5. The original unweighted data are denoted in blue dots and LSTM weighted data are in orange dots. The covariates are ordered from bottom to top

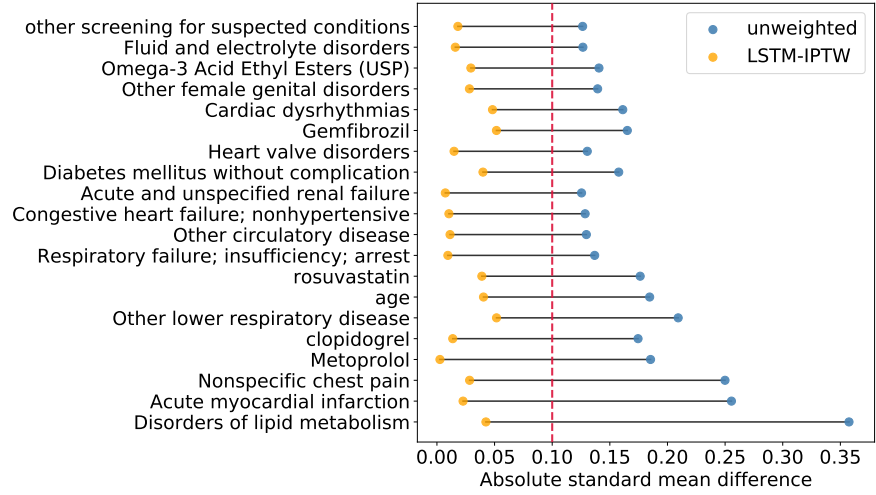
according to the increase of differences between SMD values of unweighted data and LSTM weighted data. According to the figure, we see that for both two drugs, the SMD values in original data are greater than 0.1 (i.e., the threshold of balancing), which indicates that the original observational data is highly biased and exists much confounding variables. The maximum SMD value is about 0.6 for diltiazem and 0.35 for fenofibrate. While the SMD values estimated in the LSTM weighted data are smaller than 0.1, which means that no major biases between user and non-user cohort in terms of selected covariates. The selected covariates include demographics (e.g., age), co-prescribed drugs (metformin, metoprolol, etc.) and co-morbidities (e.g., acute myocardial infarction, cardiac dysrhythmias, etc.). After correcting for these confounding variables, we can have a more accurate estimation of the treatment effect on the diseases.

Discussion

In this study, we present a computational drug repurposing framework for high-throughput screening of on-marked drugs by emulating a corresponding RCT for each drug and evaluating its treatment effect on various disease outcomes. We propose a deep learning based propensity score model for correcting selection biases and confounding in longitudinal observational data. We demonstrate our framework in a case study of CAD and evaluate 55 different repurposing drug candidates on two disease outcomes. According to the results, we obtain 6 drug candidates (i.e., metoprolol, fenofibrate, hydrochlorothiazide, pravastatin, simvastatin, valsartan) that improve the CAD outcomes with statistical significance but have not been indicated for treating CAD.



(a)



(b)

Figure 5. The SMD values of top 20 well balanced covariates. Fig. 5(a) shows results of diltiazem. Fig. 5(b) shows results of fenofibrate.

We also develop a base version of our model that uses logistic regression (LR) for computing propensity score and treatment effect estimation. We conduct comparison experiments on the base model (LR-IPTW) and our model (LSTM-IPTW) on the above two case drugs and show the results for diltiazem in Fig. 6 (The results for fenofibrate can be found in Supplementary Materials).

As the feature balancing is one of the most important evaluation metrics, we first plot the distribution of absolute SMD values computed by LSTM-IPTW and LR-IPTW in Fig. 6(a) and Fig. 6(d). In both LSTM weighted data and LR weighted data, many features exhibit large absolute SMD values (greater than 0.1) in the original data, while most features exhibit low absolute SMD (below 0.1) after re-weighting. Specifically, less features exhibit absolute SMD values above 0.1 threshold after weighted by LSTM model than weighted by LR model.

This indicates that the data is well-balanced by LSTM-IPTW and the estimated ATE from LSTM-IPTW should be more accurate than LR-IPTW. Figure 6(b) and Figure 6(e) show the propensity distribution plot over user and non-user cohorts using LSTM-IPTW and LR-IPTW models. We observe that the propensity distribution of LSTM-IPTW is more smooth (i.e., the propensities are normally distributed) than the distribution of LR-IPTW. Under LR-IPTW model, many of the patients in non-user cohorts are predicted to have a propensity of 0. We also evaluated our models using conventional metrics. The ROC curve is a standard metric widely used to estimate the performance of prediction models. The area under ROC curve (AUC) characterize the accuracy of the prediction results. Figure 6(c) and Figure 6(f) show the ROC curves for LSTM-IPTW model and LR-IPTW model. The "Propensity" curves in the figures are the standard ROC curves of LSTM

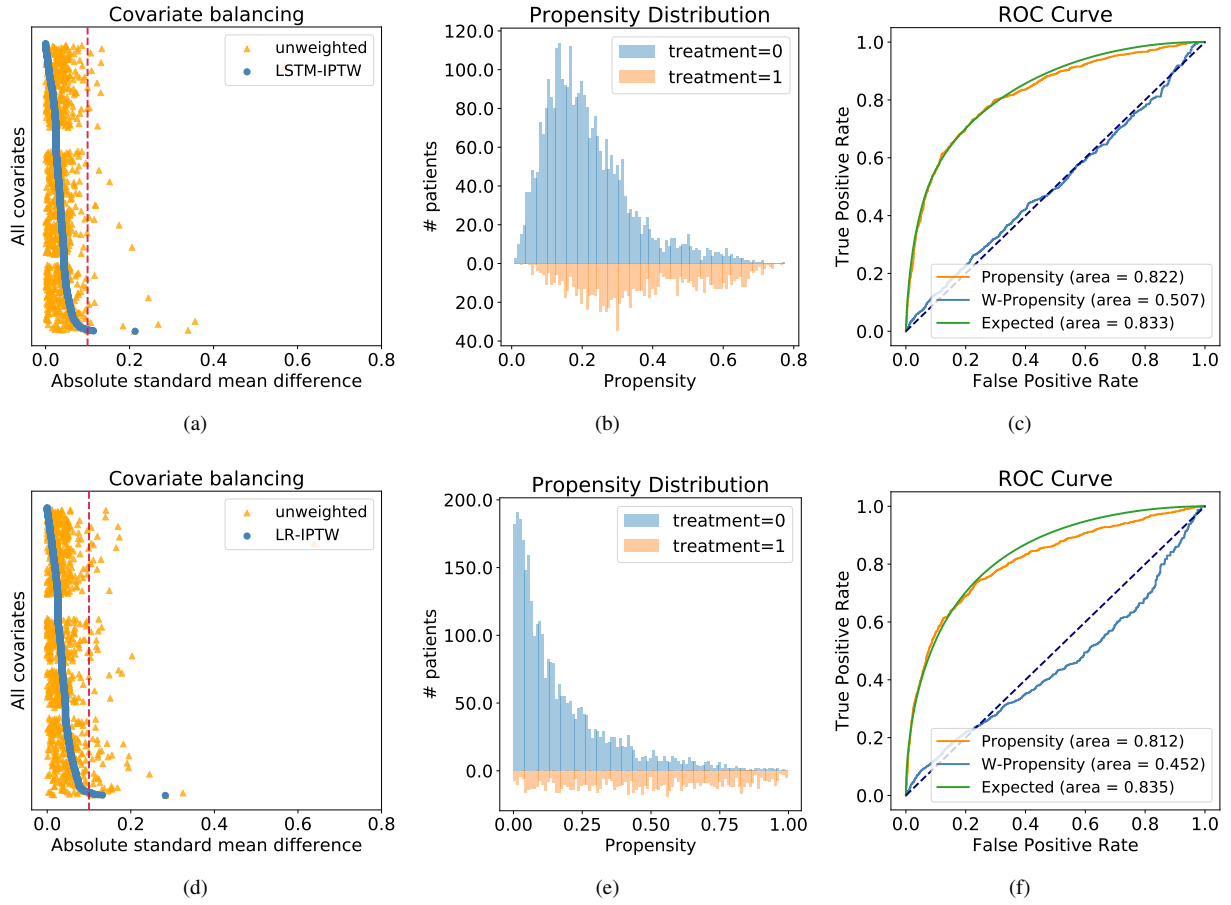


Figure 6. Performance comparison of LSTM-IPTW and LR-IPTW using drug candidate: diltiazem (*with* known CAD indication). The three figures on the top are results obtained from LSTM-IPTW, while the figures on the bottom are from LR-IPTW. Figure 6(a) and Figure 6(d) show the absolute SMD of each covariate in the original data (orange triangles) and in the weighted data (blue circles). Figure 6(b) and Figure 6(e) show the distribution of estimated propensity scores over user (orange area) and non-user (blue area) cohorts. Figure 6(c) and Figure 6(f) show the ROC curves for the propensity model (orange), expected value (green) and weighted propensity (blue).

model and LR model. By comparing the AUC values of two models, we see that LSTM model yields more accurate prediction results than the LR model. With the accurate treatment predictions, the model would generate better weights for balancing and treatment effect estimates in the following tasks. Besides the standard ROC curve, we also show another two curves: weighted propensity curve and expected curve. The weighted propensity curve is obtained by re-weighting the standard ROC curve using the weights drawn from the propensity model (the same weights applied in covariates balancing and effect estimates). This curve should be very close to the curve that would arise by a random assignment (i.e., with an AUC close to 0.5) which indicates our assumption that the weighting can emulate an RCT. From the plots, we can find that LSTM-IPTW performs better than LR-IPTW in terms of more close value to 0.5. Compared to the standard propensity ROC curve, "Expected" ROC curve duplicates the population and assign weights to each individual based on the propensity.

In this setting, each patient contributes their propensity to the true positives and $(1 - \text{propensity})$ to the false positives. The standard propensity ROC curve should be close to the expected propensity ROC. We observe that LSTM-IPTW's "Propensity" curve is much close to its "Expected" curve when compared to LR-IPTW.

This study can be extended in multiple directions in the future. For the study, we use the hypothesized confounders including demographics, co-morbidities and co-prescribed drugs. Some other potential confounders such as time elapsed from the first disease diagnosis to index date and outcome value calculated over the baseline period could be considered to build the model in the future work. Also, we can consider drug combinations and estimate the effect of entire combinations on the disease outcomes.

In summary, we demonstrate that the proposed computational drug repurposing framework can successfully identify drug candidates that have a beneficial effect on disease out-

comes but not have been indicated for CAD patients yet. The new LSTM-IPTW model shows better performance for correcting biases and estimating treatment effect than LR-IPTW, and remaining the interpretability for recognizing significant confounding.

Methods

In this section, we first introduce the study design, which includes definitions of cohorts and study variables. Then we illustrate our deep learning model with three main components in detail.

Study design

Our framework identifies drug repurposing candidates using MarketScan CAD data to emulate a bulk of corresponding RCTs. Below we describe the design of the emulated trials and the key components of our framework for CAD drug repurposing.

User and non-user cohorts

Given the drug tested in trial, a patient is assigned to the *user cohort* if the following inclusion criteria are satisfied: (1) the patient has persistently prescribed the drug (e.g., the interval between two prescriptions is less than 30 days); (2) the patient is eligible for trial at the time of the first prescription for the drug. In the CAD study, this condition is that the first prescription is after CAD initiation data; (3) the patient had a history in the database of at least one year (365 days) prior to the first prescription of the drug.

Estimating the effect of a drug requires comparing the user cohort to a control group assigned with alternative drugs. Once the alternative drugs are determined, the *non-user cohort* is defined by the same inclusion criteria described above – but with respect to the alternative drugs. To avoid overlap between the user and non-user cohorts, the framework further excludes from the non-user cohort any patient prescribed with the trial's drug. In our study design, alternative drugs are selected randomly from the prescribed ingredients, excluding the trial drug itself. Such a control group directly compares the trial's drug to drugs of the same therapeutic indication, reducing confounding by indication. We use the term *index-date* to refer to the date of the first prescription of the assigned drug, that is, the first time the trial's drug (respectively, the alternative drug) was prescribed for patients in the user (respectively, non-user) cohort.

Baseline and follow-up periods

We refer to the time period prior to the *index-date* for which we have information on the patient as the baseline period. We use the baseline period for characterizing the patients prior to the beginning of the treatment with the assigned drug. The follow-up period starts at the *index-date*, that is, at the beginning of the treatment with the trial's drug in the user cohort, and the control-drug in the non-user cohort. The effect of the drug is evaluated during the follow-up period. In the

CAD study, the baseline period is at least 365 days, and the follow-up period is 2 years (730 days). Figure 7 demonstrates the definition of user and non-user cohorts.

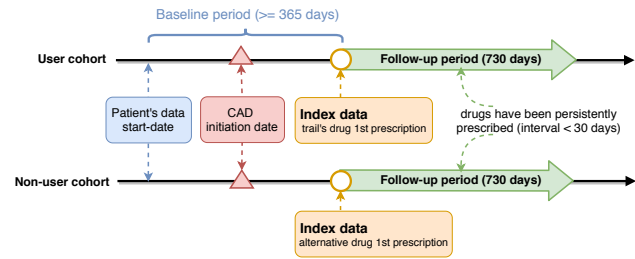


Figure 7. The definition of user and non-user cohorts

Outcomes and hypothesized confounders

The effect of the drug during the follow-up period is defined with respect to various disease outcomes. In this CAD drug repurposing case study, we consult domain experts to define a set of clinically-relevant events linked with CAD as the outcome, e.g., heart failure onset (Supplemental Table 2) and stroke onset (Supplemental Table 3). The treatment effect is estimated on these outcomes during the follow-up period (i.e., 730 days after *index-date*). The patient is considered to have the disease outcome if either of them happens in follow-up period.

Confounders are variables affecting both treatment assignment of patients and an outcome used in the trial, thus creating a "backdoor path" that may hinder the true effect of the drug on the outcome. We consult domain experts to compile a list of hypothesized confounders for the CAD study, including demographics (e.g., age at the index date, sex), co-morbidities (e.g., indicator per each ICD-9/10 diagnosis class) and co-prescribed drugs. Since confounders affect treatment assignment, they are computed on the baseline period.

Repurposing drug ingredients

We regard a drug as a repurposing candidate if it satisfies the following condition: (1) an active ingredient (i.e., the ingredient directly involved in achieving the mediation objectives). (2) persistently prescribed to a large enough number of patients in the disease cohort. Specifically, an ingredient is considered being used by a patient only if it was prescribed in two or more distinct dates, as least one month apart. And a minimum of 500 patients prescribed a certain ingredient was required. For each repurposing candidate, we can compute the user and non-user cohorts according to the above definition of cohorts. After obtaining the corresponding user and non-user cohorts, we can extract outcomes and hypothesized confounders for each individual patient from the database. Every patient in their sub-cohort is represented by a sequence of events, with each event providing the patient information (i.e., co-morbidities, co-prescribed drugs, etc.) that corresponds to each visit. The available data within these visits during the baseline period, combined with demographic characteristics

(i.e., age and gender collected at CAD initiation date) are used as input to the model.

Model

Estimation of average treatment effect

Our proposed framework evaluates the effect of a certain drug (i.e., trial's drug) on a clinical outcome with respect to alternative treatments. Let $\alpha = 1$ denote the treatment corresponding to the trial's drug, and $\alpha = 0$ denote the alternative treatments. We define average treatment effect (ATE) of a drug on the potential outcome Y as $ATE = \mathbb{E}(Y_1) - \mathbb{E}(Y_0)$, with $\mathbb{E}(Y_\alpha)$ denoting the potential expected prevalence of patients that would have experienced an outcome event during a complete follow-up period if all patients in the trial had been assigned with treatment α . The potential outcomes are referred as counterfactual as only one of these is observed for any given individual. By running RCT, we can measure the outcomes within user and non-user groups into which individuals are randomly assigned: $\mathbb{E}(Y_1)$ can be directly estimated as $\mathbb{E}(Y|\alpha = 1)$ and $\mathbb{E}(Y_0)$ as $\mathbb{E}(Y|\alpha = 0)$. However, in observational data (e.g., our MarketScan CAD data), treatment assignment is usually far from being random, which may depend on confounders (affecting both treatment assignment and outcome). We need to assign weights to the individuals in each group to avoid the influence of confounders.

In order to control the influence of confounders, we apply inverse probability of treatment weighting (IPTW) to create a pseudo-population from the original one by assigning a weight w_i^α to an individual i with treatment α . The weight is defined as the inverse of conditional probability (aka propensity score) that an individual is treated with α given the confounding values. One common issue with IPTW is that individuals with a propensity score very close to 0 will end up with an extremely high weight, potentially making the weighted estimator unstable. We address this problem by adopting an alternative weighting function called standardized IPTW [22], which uses the marginal probability of treatment instead of 1 in the weight numerator.

For estimating the propensity score, logistic regression is the most popular method in statistics [35, 36]. In longitudinal observational data, those observational covariates are not a set of static feature vectors (one for each patient), but irregularly sampled time series (recording diagnoses, medications, etc. at each timestamp). Thus, logistic regression is not ideal for effectively modeling longitudinal observational data.

Model for propensity score weighting

The schematic view of our model is shown in Figure 2, which consists of three main components: embedding module, recurrent neural network and prediction module. Briefly, the model estimates the propensity score by first transforming the input features using an embedding layer. These embedded features are then fed into LSTM, the output of which at every time point is aggregated through an attention layer for automatically focusing on important time points. The aggregated

features are fed into a prediction module that provides the probability of receiving treatment. Each of these is discussed below in detail.

Embedding module

The embedding module is to convert the initial high-dimensional and sparse input features into a lower-dimensional and continuous data representations, which is beneficial to the following prediction task. As shown in Fig. 2, the input features consist of three components: diagnosis, prescription and demographic information (age and gender). The diagnosis codes for each patient at each timestamp can be denoted as $\{\mathbf{d}_1, \mathbf{d}_2, \dots, \mathbf{d}_t\}$, and prescription can be denoted as $\{\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_t\}$. Here, \mathbf{d}_t and \mathbf{p}_t are both one dimensional binary vector with the size of diagnosis code dictionary (r) and prescription code dictionary (s), respectively. For each element in the vector, the value one in the j -th column indicates that code j is documented in t -th visit. We use two linear embedding modules to represent diagnosis and prescription respectively. That is, we define $\mathbf{e}_t = \mathbf{W}_{emb}^d \mathbf{d}_t$, $\mathbf{f}_t = \mathbf{W}_{emb}^p \mathbf{p}_t$, where $\mathbf{e}_t \in \mathbb{R}^m$ denotes the embedding of the input vector $\mathbf{d}_t \in \mathbb{R}^r$, m the size of the diagnosis embedding dimension, $\mathbf{W}_{emb}^d \in \mathbb{R}^{m \times r}$ the embedding matrix. $\mathbf{f}_t \in \mathbb{R}^n$ denotes the embedding of the input vector $\mathbf{p}_t \in \mathbb{R}^s$, n the size of the diagnosis embedding dimension, $\mathbf{W}_{emb}^p \in \mathbb{R}^{n \times s}$ the embedding matrix. The age is normalized into range of $[0, 1]$ using min-max normalization and the gender is represented as a binary vector. Having the embedded vectors of patients, we input them to LSTM.

Recurrent neural network and Attention mechanism

Long short-term memory (LSTM) [15], which is a kind of recurrent neural network (RNN) equipped with memory cells, can better model temporality of observational data. LSTM and its variations are widely adopted in the scenario that contains sequential and temporal data, such as in language translation [37], speech recognition [38] and image captioning [39]. A common LSTM unit contains a cell, an input gate, an output gate and a forget gate. The cell can remember values over irregular time intervals and the three gates moderate the flow of information into and out of the cell. The inputs to the LSTM are embedded confounding vectors from the embedding module and the output of which is patient's latent health status at the time of visit. We use two LSTMs, $LSTM_\alpha$ and $LSTM_\beta$ to separately model diagnosis and prescription codes of patients.

$$\begin{aligned} \mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_t &= LSTM_\alpha(\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_t) \\ \mathbf{g}_1, \mathbf{g}_2, \dots, \mathbf{g}_t &= LSTM_\beta(\mathbf{f}_1, \mathbf{f}_2, \dots, \mathbf{f}_t) \end{aligned} \quad (3)$$

where $\mathbf{h}_t \in \mathbb{R}^u$, $\mathbf{g}_t \in \mathbb{R}^v$ are hidden state vectors at t -th visit, and u, v denote the size of hidden layer of $LSTM_\alpha$ and $LSTM_\beta$. Then those patient hidden states are aggregated through two separate attention layers for automatically focusing on impor-

tant visits.

$$\begin{aligned}
\alpha_i &= \text{Softmax}(\mathbf{W}_\alpha^\top \mathbf{h}_i + \mathbf{b}_\alpha), \quad \text{for } i = 1, 2, \dots, t \\
\mathbf{c}_\alpha &= \sum_{i=1}^t \alpha_i \mathbf{h}_i \\
\beta_i &= \text{Softmax}(\mathbf{W}_\beta^\top \mathbf{g}_i + \mathbf{b}_\beta), \quad \text{for } i = 1, 2, \dots, t \\
\mathbf{c}_\beta &= \sum_{i=1}^t \beta_i \mathbf{g}_i
\end{aligned} \tag{4}$$

where $\mathbf{W}_\alpha \in \mathbb{R}^u$, $\mathbf{b}_\alpha \in \mathbb{R}^u$, $\mathbf{W}_\beta \in \mathbb{R}^v$ and $\mathbf{b}_\beta \in \mathbb{R}^v$ are the parameters to learn. Using the generated attention weights for diagnosis and prescription, we obtain the aggregated vectors $\mathbf{c}_\alpha \in \mathbb{R}^u$ and $\mathbf{c}_\beta \in \mathbb{R}^v$ as defined in Eq. 4. Then we combine \mathbf{c}_α , \mathbf{c}_β with vectorized age and gender to predict the probability of receiving a treatment (propensity score).

Prediction module

The aggregated patient states from attention layer $\mathbf{c}_\alpha, \mathbf{c}_\beta$, combined with the demographic features c_{demo} , are passed through a fully-connected neural network predict the probability of receiving a treatment as follows,

$$\hat{y} = \text{Sigmoid}(\mathbf{W}^\top \mathbf{c}_t + b) \tag{5}$$

where $\mathbf{c}_t = \text{ReLU}(\mathbf{W}_c[\mathbf{c}_\alpha, \mathbf{c}_\beta, c_{demo}] + \mathbf{b}_c)$, $\mathbf{W}_c \in \mathbb{R}^{k \times (u+v+2)}$, $\mathbf{b}_c \in \mathbb{R}^k$, $\mathbf{W} \in \mathbb{R}^k$, $b \in \mathbb{R}$ are the model parameters. We use cross-entropy to calculate the prediction loss as follows,

$$\mathcal{L} = -\frac{1}{N} \sum_{i=1}^N (y_i \log \hat{y}_i + (1 - y_i) \log (1 - \hat{y}_i)) \tag{6}$$

where y_i is the ground truth of observed treatment for patient i .

Experiment settings

The model is implemented and trained with Python 3.6 and PyTorch 1.4¹, on a high-performance computing cluster with four NVIDIA TITAN RTX 6000 GPUs. For each drug candidate, we train a model using the adaptive moment estimation (Adam) algorithm with a batch size of 50 subjects and the learning rate is 0.001. We run each model for 50 iterations for computing p-value and confidence interval. We randomly split the input data into training, validation and test sets with a ratio of 70%, 10%, 20%. The information from a given patient is only present in one set. The training set is to train the proposed models. The validation set is used to improve the models and select the best model hyperparameters.

Data availability

The access of the MarketScan data analyzed in this manuscript is provided by The Ohio State University. The dataset may be available from IBM at <https://www.ibm.com/products/marketscan-research-databases>.

Code availability

The code for this paper is available at <https://github.com/ruoqi-liu/DeepIPW>

References

1. Langedijk, J., Mantel-Teeuwisse, A. K., Slijkerman, D. S. & Schutjens, M.-H. D. Drug repositioning and repurposing: terminology and definitions in literature. *Drug discovery today* **20**, 1027–1034 (2015).
2. Ashburn, T. T. & Thor, K. B. Drug repositioning: identifying and developing new uses for existing drugs. *Nat. reviews Drug discovery* **3**, 673 (2004).
3. Pushpakom, S. *et al.* Drug repurposing: progress, challenges and recommendations. *Nat. Rev. Drug Discov.* **18**, 41 (2019).
4. Luo, H. *et al.* Dpdr-cpi, a server that predicts drug positioning and drug repositioning via chemical-protein interactome. *Sci. reports* **6**, 35996 (2016).
5. Dakshanamurthy, S. *et al.* Predicting new indications for approved drugs using a proteochemometric method. *J. medicinal chemistry* **55**, 6832–6848 (2012).
6. Sanseau, P. *et al.* Use of genome-wide association studies for drug repositioning. *Nat. biotechnology* **30**, 317 (2012).
7. Iorio, F. *et al.* Discovery of drug mode of action and drug repositioning from transcriptional responses. *Proc. Natl. Acad. Sci.* **107**, 14621–14626 (2010).
8. Sirota, M. *et al.* Discovery and preclinical validation of drug indications using compendia of public gene expression data. *Sci. translational medicine* **3**, 96ra77–96ra77 (2011).
9. Buchan, N. S. *et al.* The role of translational bioinformatics in drug discovery. *Drug discovery today* **16**, 426–434 (2011).
10. Sherman, R. E. *et al.* Real-world evidence—what is it and what can it tell us. *N Engl J Med* **375**, 2293–2297 (2016).
11. Cheng, F. *et al.* Network-based approach to prediction and population-based validation of in silico drug repurposing. *Nat. communications* **9**, 2691 (2018).
12. Xu, H. *et al.* Validating drug repurposing signals using electronic health records: a case study of metformin associated with reduced cancer mortality. *J. Am. Med. Informatics Assoc.* **22**, 179–191 (2014).
13. Hernán, M. A. & Robins, J. M. Using big data to emulate a target trial when a randomized trial is not available. *Am. journal epidemiology* **183**, 758–764 (2016).
14. D’Agostino, R. B. Estimating treatment effects using observational data. *Jama* **297**, 314–316 (2007).
15. Hochreiter, S. & Schmidhuber, J. Long short-term memory. *Neural computation* **9**, 1735–1780 (1997).
16. Hirano, K., Imbens, G. W. & Ridder, G. Efficient estimation of average treatment effects using the estimated propensity score. *Econometrica* **71**, 1161–1189 (2003).
17. Truven health marketscan. Accessed Nov. 2019.

¹<https://pytorch.org/>

18. for Health Statistics, N. C. Classification of diseases, functioning, and disability (2019). Accessed Nov. 2019.
19. The observational health data sciences and informatics (ohdsi) (2019). Accessed Nov. 2019.
20. Causes of heart failure (2017). Accessed Nov. 2019.
21. Conditions that increase risk for stroke (2018). Accessed Nov. 2019.
22. Austin, P. C. An introduction to propensity score methods for reducing the effects of confounding in observational studies. *Multivar. behavioral research* **46**, 399–424 (2011).
23. Efron, B. & Tibshirani, R. Bootstrap methods for standard errors, confidence intervals, and other measures of statistical accuracy. *Stat. science* 54–75 (1986).
24. Benjamini, Y. & Hochberg, Y. Controlling the false discovery rate: a practical and powerful approach to multiple testing. *J. Royal statistical society: series B (Methodological)* **57**, 289–300 (1995).
25. Kuhn, M., Campillos, M., Letunic, I., Jensen, L. J. & Bork, P. A side effect resource to capture phenotypic effects of drugs. *Mol. systems biology* **6** (2010).
26. Wishart, D. S. *et al.* Drugbank 5.0: a major update to the drugbank database for 2018. *Nucleic acids research* **46**, D1074–D1082 (2018).
27. Fisher, M. L. *et al.* Beneficial effects of metoprolol in heart failure associated with coronary artery disease: a randomized trial. *J. Am. Coll. Cardiol.* **23**, 943–950 (1994).
28. Wong, T. Y., Simó, R. & Mitchell, P. Fenofibrate—a potential systemic treatment for diabetic retinopathy? *Am. journal ophthalmology* **154**, 6–12 (2012).
29. Hydrochlorothiazide-the american society of health-system pharmacists (2019). Accessed Nov. 2019.
30. Hydrochlorothiazide completed phase 4 trials for coronary artery disease / high blood pressure (hypertension) treatment (2018). Accessed Nov. 2019.
31. Jukema, J. W. *et al.* Effects of lipid lowering by pravastatin on progression and regression of coronary artery disease in symptomatic men with normal to moderately elevated serum cholesterol levels: the regression growth evaluation statin study (regress). *Circulation* **91**, 2528–2540 (1995).
32. Kjekshus, J., Pedersen, T. R., Olsson, A. G., Færgeman, O. & Pyörälä, K. The effects of simvastatin on the incidence of heart failure in patients with coronary heart disease. *J. cardiac failure* **3**, 249–254 (1997).
33. Higuchi, T., Abletshauser, C., Nekolla, S. G., Schwaiger, M. & Bengel, F. M. Effect of the angiotensin receptor blocker valsartan on coronary microvascular flow reserve in moderately hypertensive patients with stable coronary artery disease. *Microcirculation* **14**, 805–812 (2007).
34. Resource, S. . . S. E. Indications of diltiazem in sider (2019). Accessed Nov. 2019.
35. D’Agostino Jr, R. B. Propensity score methods for bias reduction in the comparison of a treatment to a non-randomized control group. *Stat. medicine* **17**, 2265–2281 (1998).
36. Weitzen, S., Lapane, K. L., Toledano, A. Y., Hume, A. L. & Mor, V. Principles for modeling propensity scores in medical research: a systematic literature review. *Pharmacoepidemiol. drug safety* **13**, 841–853 (2004).
37. Sutskever, I., Vinyals, O. & Le, Q. Sequence to sequence learning with neural networks. *Adv. NIPS* (2014).
38. Graves, A., Mohamed, A.-r. & Hinton, G. Speech recognition with deep recurrent neural networks. In *2013 IEEE international conference on acoustics, speech and signal processing*, 6645–6649 (IEEE, 2013).
39. Mao, J. *et al.* Deep captioning with multimodal recurrent neural networks (m-rnn). *arXiv preprint arXiv:1412.6632* (2014).

Acknowledgements

This work was funded in part by the National Center for Advancing Translational Research of the National Institutes of Health under award number CTSA Grant UL1TR002733. The content is solely the responsibility of the authors and does not necessarily represent the official views of the National Institutes of Health.

Author contributions

PZ conceived the project. RL and PZ developed the method. RL conducted the experiments. RL, LW and PZ analyzed the results. RL, LW and PZ wrote the manuscript. All authors read and approved the final manuscript.

Competing interests

The authors declare no competing interests.

Additional information

Supplemental Table 1. The definition of coronary artery disease (CAD) from observational health data. Supplemental Table 2-3. The definition of heart failure and stroke from observational health data. Supplemental Table 4. Main results for all 55 repurposing drug candidates. Supplemental Figure 1. Performance comparison of LSTM-IPTW and LR-IPTW on case drug: fenofibrate.