

Deep Learning-based Propensity Scores for Confounding Control in Comparative Effectiveness Research

A Large-scale, Real-world Data Study

Janick Weberpals,^a Tim Becker,^b Jessica Davies,^c Fabian Schmich,^a Dominik Rüttinger,^d
Fabian J. Theis,^{e,f} and Anna Bauer-Mehren^a

Background: Due to the non-randomized nature of real-world data, prognostic factors need to be balanced, which is often done by propensity scores (PSs). This study aimed to investigate whether autoencoders, which are unsupervised deep learning architectures, might be leveraged to compute PS.

Methods: We selected patient-level data of 128,368 first-line treated cancer patients from the Flatiron Health EHR-derived de-identified database. We trained an autoencoder architecture to learn a lower-dimensional patient representation, which we used to compute PS. To compare the performance of an autoencoder-based PS with established methods, we performed a simulation study. We assessed the balancing and adjustment performance using standardized mean differences, root mean square errors (RMSE), percent bias, and confidence interval coverage. To illustrate the application of the autoencoder-based PS, we emulated the PRONOUNCE trial by applying the trial's protocol elements within an observational database setting, comparing two chemotherapy regimens.

Results: All methods but the manual variable selection approach led to well-balanced cohorts with average standardized mean differences <0.1. LASSO yielded on average the lowest deviation of resulting estimates (RMSE 0.0205) followed by the autoencoder approach (RMSE 0.0248). Altering the hyperparameter setup in sensitivity analysis, the autoencoder approach led to similar results as LASSO (RMSE 0.0203 and 0.0205, respectively). In the case study, all methods provided a similar conclusion with point estimates clustered around the null (e.g., $HR_{\text{autoencoder}} = 1.01$ [95% confidence interval = 0.80, 1.27] vs. $HR_{\text{PRONOUNCE}} = 1.07$ [0.83, 1.36]).

Conclusions: Autoencoder-based PS computation was a feasible approach to control for confounding but did not perform better than some established approaches like LASSO.

Keywords: Autoencoder; Causal inference; Comparative effectiveness research; Deep learning; Electronic health records; Machine learning; Propensity scores

(Epidemiology 2021;32: 378–388)

Submitted December 3, 2019; accepted January 27, 2021

From the ^aData Science, Pharmaceutical Research and Early Development Informatics (pREDi), Roche Innovation Center Munich (RICM), Penzberg, Germany; ^bxValue GmbH, Willich, Germany, on behalf of Data Science IV, Pharmaceutical Research and Early Development Informatics (pREDi), Roche Innovation Center Munich (RICM), Penzberg, Germany; ^cF. Hoffmann-La Roche Ltd, Welwyn Garden City, United Kingdom; ^dEarly Clinical Development Oncology, Pharmaceutical Research and Early Development (pRED), Roche Innovation Center Munich (RICM), Penzberg, Germany; ^eInstitute of Computational Biology, German Research Center for Environmental Health, Helmholtz Center Munich, Neuherberg, Germany; and ^fDepartment of Mathematics, Technical University of Munich, Garching, Germany.

Fabian J. Theis and Anna Bauer-Mehren contributed equally.

J.W. conceptualized the study, carried out the analysis and statistical programming, and drafted the article. A.B.-M. and F.J.T. supervised the project and gave significant advice at various stages of the project. T.B., F.J.T., and F.S. significantly contributed to the analysis of the data. F.S. and T.B. assisted with the machine learning setups used in this study and with the programming code. D.R. gave valuable insights to the clinical characteristics and clinical interpretation of the data. J.D. significantly contributed to the curation of the Flatiron Health database, ensured various data quality measures, contributed to the conceptualization of the case study, and helped with the interpretation of the data. All authors critically reviewed and edited the article draft. All authors agree to the submission of the article.

Copyright © 2021 Wolters Kluwer Health, Inc. All rights reserved.

ISSN: 1044-3983/21/323-378

DOI: 10.1097/EDE.0000000000001338

This work was funded by Roche as part of the main authors (J.W.) postdoctoral fellowship program (RPF-ID: 498).

J.W., F.S., D.R., J.D., and A.B.-M. are paid employees of Roche. J.W., F.S., D.R., J.D., and A.B.-M. report holding shares in Roche. T.B. is an employee of xValue and external consultant to Roche. F.J.T. reports receiving consulting fees from Roche Diagnostics GmbH and Cellarity Inc., and ownership interest in Cellarity, Inc.

The computing code used in this study is available as Python Jupyter Markdown scripts (.html) as supplementary material. All of the analyses described in the article were performed in R version 3.2.2. The PCA and autoencoder training was performed using scikit-learn and Keras with Tensorflow backend in Python version 3.6.0, respectively. The code that was used for the simulation is available as Rmarkdown. The data that support the findings of this study have been originated by Flatiron Health, Inc. These de-identified data may be made available upon request, and are subject to a license agreement with Flatiron Health; interested researchers should contact <DataAccess@flatiron.com> to determine licensing terms.

Ethical review: These research activities are covered in Flatiron's parent protocol which is reviewed and approved by a central IRB.

SDC Supplemental digital content is available through direct URL citations in the HTML and PDF versions of this article (www.epidem.com).

Correspondence: Anna Bauer-Mehren, Pharmaceutical Research and Early Development informatics (pREDi), Roche Diagnostics GmbH, PXID...2246, Nonnenwald 2, 82377 Penzberg, Germany. E-mail: anna.bauer-mehren@roche.com; Fabian J. Theis, Helmholtz Zentrum Munich, Institute of Computational Biology, Ingolstädter Landstr. 1, 85764 Neuherberg, Germany. E-mail: fabian.theis@helmholtz-muenchen.de.

Randomized controlled trials (RCTs) are the gold standard when evaluating the effects of interventions on health-related outcomes. However, the digitization of healthcare infrastructure, such as electronic health records (EHR), and a boost in computational power in the past years have led to an increase in evidence generated by routinely collected healthcare data, often termed real-world data.^{1–3}

Due to the heterogeneous and non-randomized nature of these data, such analyses inherit the chance to lead to misleading conclusions when biases, such as confounding bias, are not addressed appropriately.⁴ Therefore, propensity score (PS) techniques are popular analytical approaches to balance patient characteristics in observational research.⁵ In general, PSs are defined as an individual's (i) conditional probability to be assigned to a particular treatment (Z_i) given observed baseline covariates (X_i) with $Pr(Z_i = 1/X_i)$.⁶ By conditioning on the PS, researchers try to create positivity; that is, if a given combination of covariate values is observed in one cohort, it should also appear in the other cohort under comparison.⁷ Under the assumption of no unmeasured confounding and a correctly specified PS model, unbiased treatment effects may be estimated, for example, via matching or weighting on the PS.

There is ongoing debate about the ideal strategy to correctly specify the PS^{8,9} and in the majority of cases, logistic regression models are fitted using a set of a priori investigator-defined covariates.¹⁰ This approach is straightforward but may be error-prone when interaction terms or higher-order relationships are not appropriately modeled.¹¹ Moreover, as healthcare databases are getting increasingly complemented by more dimensions like genomics, selecting the correct set of covariates on a manual basis becomes infeasible and automatable data-adaptive methods are warranted.

With the ability to handle high-dimensional datasets in a nonlinear and automatable fashion, deep learning models are highly attractive approaches to solve these problems.¹² We aimed to investigate if autoencoders, which are unsupervised deep learning encoder-decoder architectures that learn a latent non-linear lower-dimensional covariate representation, might be leveraged as a data-adaptive alternative to compute PS for comparative effectiveness research.

The objective of this study is two-fold. First, we compare the performance of covariate balancing and confounding bias reduction with the autoencoder-based PS as compared with established adjustment strategies in a simulation framework among cancer patients with a first-line (1L) systemic anticancer treatment. In the second part of this study, we will emulate the 2015 published PRONOUNCE trial¹³ by applying the trial's major protocol elements to the observational database setting of this study to illustrate the application of the autoencoder-based PS to a real comparative effectiveness use case.

METHODS

Data Sources and Study Population

For this retrospective real-world data study, we used the nationwide Flatiron Health EHR-derived de-identified database which includes data from over 280 cancer clinics including more than 2.2 million US cancer patients available for analysis. The de-identified patient-level data in the EHRs include structured data (e.g., laboratory values and prescribed drugs) in addition to unstructured data collected via technology-enabled chart abstraction from physician's notes and other unstructured documents (e.g., biomarker reports). In this study, we selected patients out of tumor-specific databases and pooled them into a single cohort. Patients were eligible to be included if they were diagnosed with any primary tumor and received a 1L systemic anticancer treatment (CONSORT diagram, Figure 1).

Data Curation and Covariate Ascertainment

We considered covariates for modeling if they were applicable across all tumor types and for at least 20% of all patients (eTable 1; <http://links.lww.com/EDE/B777>). We imputed missing covariates or those with implausible values (as defined as being outside of $1.5 \times$ the interquartile range from the quartiles Q1 and Q3, respectively¹⁴) using median imputation for continuous covariates or assigning a missing indicator category to one-hot encoded categorical variables.^{15,16} In addition, we derived empirical covariates (EC) of lab and vital sign tests. As the Flatiron Health EHR-derived de-identified database does not contain records of claims, procedure codes, and outpatient diagnosis codes, these EC were derived from the frequency of clinical laboratory tests and vital sign tests (which corresponds to steps 1–3 of the high-dimensional PS algorithm¹⁷), which resulted in 123 additional covariates (eTable 1; <http://links.lww.com/EDE/B777>). All covariates were measured at or before the start of 1L therapy (=index date) with a maximum lookback window period of 90 days relative to the index date (eFigure 1; <http://links.lww.com/EDE/B777>).^{18,19}

Non-linear Latent Variables and Propensity Scores Computed by Autoencoder

The following section briefly illustrates the autoencoder-based PS computation (terminology used in this paragraph is defined in eAppendix1; <http://links.lww.com/EDE/B777> and in Bi et al²⁰).

Autoencoders are unsupervised neural network architectures that generally consist of an input layer, a lower-dimensional hidden “bottleneck” layer, and an output layer with the same dimensions as the input layer. Conceptually, the autoencoder-based PS computation can be described as follows (Figure 2). All available information about a patient may be defined as a high-dimensional covariate vector serving as the input layer. This input layer is sequentially compressed to arrive at a latent nonlinear lower-dimensional covariate representation in the

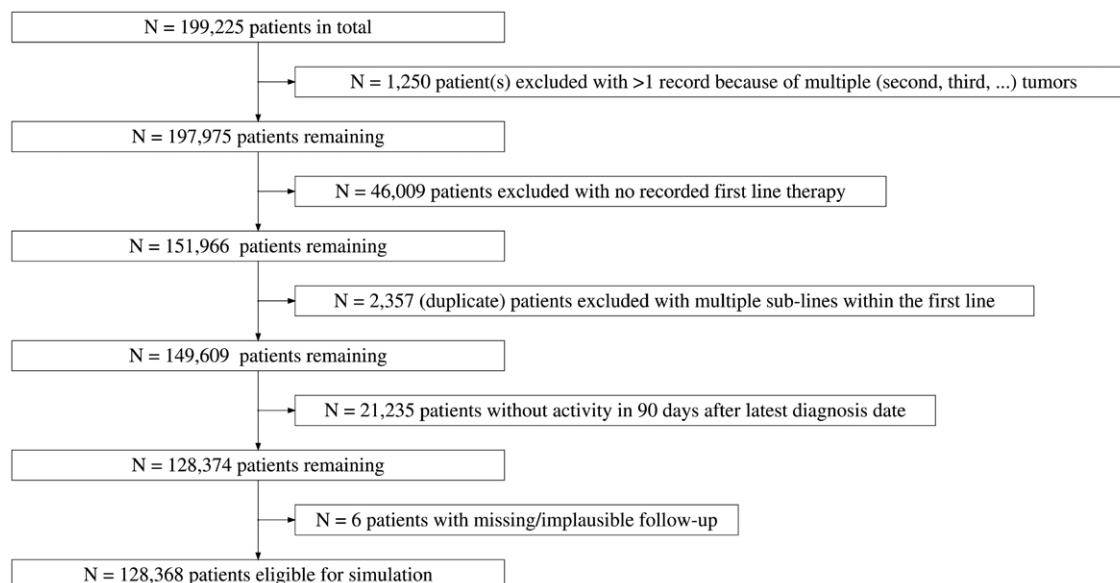


FIGURE 1. Consort diagram illustrating selection of eligible patients for simulation.

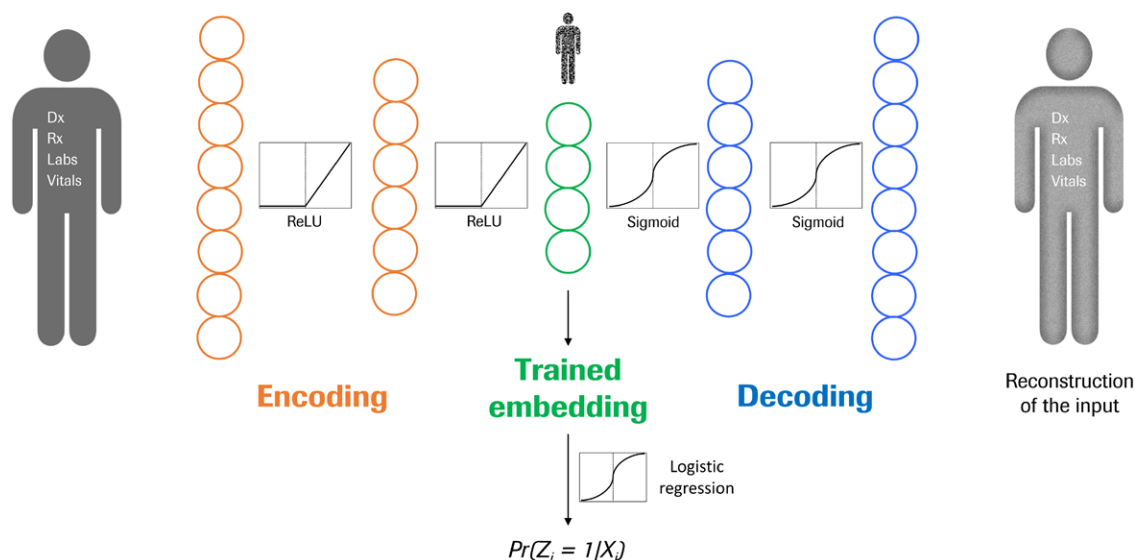


FIGURE 2. Conceptual architecture of patient representation learning and autoencoder-based propensity score computation. Figure is available in color online.

hidden bottleneck layer (encoding). Given the lower-dimensional information of the bottleneck layer, the actual input information can be reconstructed (decoding); the decoded information is leveraged in autoencoders to adjust the network parameters in each iteration by computing the loss between the actual data and the predicted reconstruction. Due to the compression and the optimization of parameters of the neural network in each encoding–decoding iteration step, the autoencoder learns essential features describing the highest variance of a dataset. This way the bottleneck layer captures the true data manifold in a much lower-dimensional representation (embedding) that can eventually be used to specify the PS.

After the above described general setup, we developed an autoencoder architecture (details on architecture, hyperparameters, and code can be found in eAppendix 1; <http://links.lww.com/EDE/B777>). To compute the PS based on the trained embedding, we used a logistic regression as the final output layer.

Propensity Score Estimation Methods for Comparison

To investigate the performance of an autoencoder-based PS, we chose established adjustment (multivariable regression) and PS estimation methods (manual variable selection, principal component analysis [PCA], and LASSO) for

comparison employing a simulation framework (Table 1). We additionally extended all machine learning models by the set of EC that were derived as described above (EC extended models 7–9). For more details see eAppendix 1; <http://links.lww.com/EDE/B777>.

Simulation Setup

The overall objective was to simulate different realistic scenarios of confounding bias between a fictional head-to-head drug comparison and to investigate the resulting balancing and adjustment after 1:1 PS matching with PS computed using the aforementioned PS estimation methods (Table 1). We defined the outcome of interest for this simulation study as overall survival, which we computed as the time from index date to death due to any reason or censoring.

The general simulation algorithm is illustrated in Figure 3A. In brief, all eligible patients were equally randomized to either a drug A or a drug B cohort to remove any prognostic association of the covariates to the assignment probability to one of the cohorts. This resulted in a hazard ratio (HR) for overall survival of 1.00 (95% confidence interval [CI] = 0.99, 1.01), which served as the true estimate in this simulation (eFigure 2; <http://links.lww.com/EDE/B777>). In a next step, we grouped patients into prognostic quartiles (Q1–Q4) according to their baseline hazards towards the outcome (overall survival) with patients in Q1 having a good prognosis (lowest hazard) to patients in Q4 having a poor prognosis (eFigure 3 and eTable 2; <http://links.lww.com/EDE/B777>). The prognostic quartiles are based on a published prognostic score for overall survival (eFigures 4 and 5; <http://links.lww.com/EDE/B777>) that was developed within a large pan-cancer cohort and is derived from a formula with strongly

prognostic demographic, clinical, routine hematology, and blood chemistry parameters (eTable 3; <http://links.lww.com/EDE/B777>) that were modeled within a Cox proportional hazard framework to derive a multivariable prognostic risk model for overall survival.²¹ The resulting prognostic score was validated in two independent phase I and III clinical studies. To simulate baseline imbalances, we exploited the correlation between prognostic score-based balance measures for PS models with bias in the treatment effect estimate using conditional resampling as described in the following.²² Out of the drug A cohort, we sampled 10,000 patients completely at random and independent of their assignment to the prognostic quartiles to arrive at a homogenous sample with a constant prognosis in each replication step. In contrast, we sampled 10,000 patients randomized to the drug B cohort with a conditional sampling probability based on their assignment to a prognostic quartile (e.g., scenario 1: patients in Q1 were sampled with a probability of 40%, in Q2 with 30%, in Q3 with 20%, and in Q4 with 10%). Because quartile membership is associated with overall survival, the conditional sampling of the drug B cohort (as compared to the random sampling of the drug A cohort) naturally induces a spurious association, which is solely driven by the variables defining the quartiles. We applied this sampling scheme in total 27 different sampling probabilities with 100 replications each to simulate various scenarios of confounding bias yielding biased estimates with different magnitudes and directions away from the true HR of 1.00 (Figure 3B).

We finally assessed the comparative performance of each PS computation method as to how much each method was able to adjust for the above described induced spurious association. For this purpose, we matched the resulting

TABLE 1. Models and Adjustment Strategies Compared in Simulation Framework

Model	Adjustment Strategy ^a	Data-adaptive Covariate Selection/ Transformation	Covariates Adjusted for or Potential Covariates to Choose from
1	Unadjusted	—	—
2	Multivariable regression (direct outcome model)	No	Age, cancer entity, gender, stage, histology, healthcare provider, race/ethnicity, time from initial cancer diagnosis to 1L initiation, calendar year of initial cancer diagnosis
3	Manual variable selection	No	Age, cancer entity, gender, stage, histology, healthcare provider, race/ethnicity, time from initial cancer diagnosis to 1L initiation, calendar year of initial cancer diagnosis
4	LASSO	Selection	All generally available covariates ^b . Algorithm picks covariates according to shrinkage/regularization
5	PCA	Transformation	All generally available covariates ^b . Algorithm computes linear transformation of all covariates in a dataset to principal components (PCs) of which the top <i>n</i> PCs, explaining 80% variance, were chosen
6	Autoencoder	Transformation	All generally available covariates ^b . Algorithm computes lower-dimensional representation of <i>j</i> dimensions based on non-linear data operations into latent-space variables
7	LASSO EC	Transformation	Model 4 + 123 empirical covariates ^c
8	PCA EC	Selection	Model 5 + 123 empirical covariates ^c
9	Autoencoder EC	Transformation	Model 6 + 123 empirical covariates ^c

^aIn model 2 the estimate is directly computed from a multivariable regression while models 3–9 are based on propensity score matching.

^bTotal of 318 demographic, clinical, cancer-/disease-specific covariates (eTable 1; <http://links.lww.com/EDE/B777>).

^cTotal of 123 empirical frequency covariates derived, corresponds to steps 1–3 of the high-dimensional propensity score algorithm (eTable 1; <http://links.lww.com/EDE/B777>). 1L indicates first-line systemic cancer treatment; LASSO, least absolute shrinkage and selection operator; PC(A), principal component (analysis).

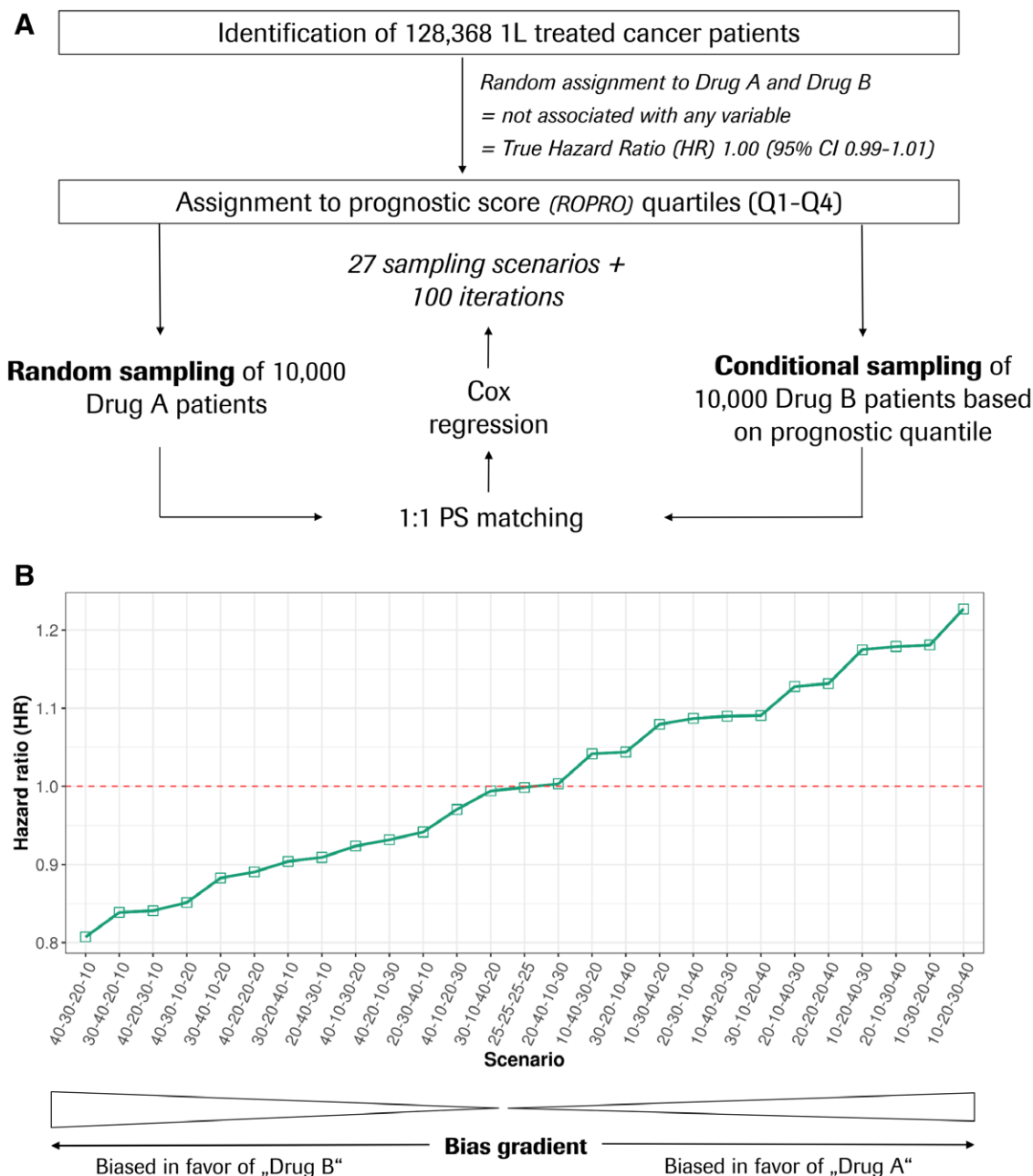


FIGURE 3. A, Sampling algorithm for simulation and (B) overview of magnitude of induced confounding bias by simulation scenario. Q indicates quartile; ROPRO, real-world prognostic score. Figure is available in color online.

cohorts without replacement in a 1:1 ratio with a caliper width of 0.2 SDs of the predicted PS logit,²³ and HRs were estimated using Cox proportional hazards regression models with a robust variance estimator.²⁴ Simulations of treatment effects other than a null treatment effect were not considered to avoid complications with the collapsibility²⁵ and proportional hazards assumption²⁶ of HRs.

We assessed the overall balance in the distribution of important baseline covariates after PS matching using standardized mean differences (SMD) with a cutoff of <0.1

indicating sufficient balance.²⁷ To assess the average deviation of the resulting HRs and the true HR of 1.00, we computed the root mean square error (RMSE) as performance metrics. To measure the uncertainty of the point estimates, we computed the coverage probability as the proportion of times the estimated 95% CI included the true HR of 1.00.^{28,29} Additionally, we estimated the absolute bias (in percent) as

$$\left| \frac{\text{HR}_{\text{pooled}} - \text{HR}_{\text{True}}}{\text{HR}_{\text{True}}} \times 100 \right| \text{ for each simulation scenario.}^{30}$$

Case Study

To illustrate the application of the autoencoder-based PS in comparative effectiveness research, we emulated the PRONOUNCE trial by applying the major protocol design elements of this trial within the observational Flatiron Health EHR-derived de-identified database.

In brief, the PRONOUNCE trial was a randomized, open-label, phase III trial aimed at evaluating the comparative efficacy of carboplatin/pemetrexed followed by pemetrexed maintenance versus bevacizumab/carboplatin/paclitaxel followed by bevacizumab maintenance as 1L treatment among advanced nonsquamous non-small-cell lung cancer patients.¹³ In terms of overall survival, the PRONOUNCE trial did not find a difference in treatment efficacy for either of the combinations, which served as our expected outcome for the emulation of this trial.

For the implementation of the major in-/exclusion criteria and study design elements, we followed the target trial emulation framework by Hernán and Robins³¹ and summarized the comparison to the original in eTable 4 and eFigure 6; <http://links.lww.com/EDE/B777>. Instead of a random assignment to either treatment strategy in a 1:1 ratio in the original trial, we applied PS matching (applying the different PS computation approaches) in a 1:1 ratio (nearest neighbor without replacement as main analysis)³² and standardized mortality ratio (SMR) weighting (sensitivity analysis).³³ We derived estimates for overall survival using Cox proportional hazards regression with the initiation of maintenance therapy as start of follow-up. The causal contrast of interest was analyzed as the counterfactual comparison of initiators of the two different treatment strategies as an observational equivalent of the RCT's intent-to-treat analysis. Further details are outlined in the supplementary methods (eAppendix 1; <http://links.lww.com/EDE/B777>).

RESULTS

The characteristics of the eligible simulation population are displayed in eTable 5; <http://links.lww.com/EDE/B777>. Results of the hyperparameter selection and evaluation are illustrated in the supplementary material (eFigures 7–11; <http://links.lww.com/EDE/B777>) and computation times for the autoencoder models and simulations are summarized in eFigures 12 and 13, and eTable 6; <http://links.lww.com/EDE/B777>, respectively.

Simulation—Balancing Properties

Figure 4 summarizes the average balancing performance of important baseline characteristics by simulation scenario. In general, most PS estimation methods led to sufficient balancing of important patient characteristics at baseline (SMD <0.1). In some scenarios, imbalances for some covariates were observed for PS computed using manual variable selection. Investigating SMDs by scenario indicated that those imbalances resulted from some of the more extreme confounded scenarios (eFigure 14; <http://links.lww.com/EDE/B777>).

Simulation—Root Mean Square Error, Percent Bias, and Coverage

The overall results across all simulated scenarios and iterations are illustrated in Table 2. Estimates without any adjustment resulted on average in high RMSEs (0.1205) and bias (10.4% bias) and low coverage (16.41%). When covariates were manually chosen (models 2 and 3), the PS method led on average to a lower RMSE (0.0670 vs. 0.0790), bias (5.73% vs. 6.75%), and a higher coverage (32.81% vs. 27.67%) as compared with choosing the same covariates for direct outcome regression, respectively. Point estimates were observed to scatter broadly around the null for both methods (eFigure 15; <http://links.lww.com/EDE/B777>). Comparisons between model standard errors and empirical standard errors indicated a less reliable variance estimation for models 1–3 (eTable 7; <http://links.lww.com/EDE/B777>). The PCA PS estimation method led to a noticeable improvement in adjustment performance as compared with selecting covariates manually with a RMSE of 0.0293 and 0.0329 for PCA and PCA EC, respectively. Employing an autoencoder-based estimation of the PS led to further improvements in RMSEs of 0.0248 and 0.0265, bias of 2.00% and 2.15%, and coverage of 87.70% and 85.19% for autoencoder and autoencoder EC, respectively. The best adjustment performance was observed with both LASSO approaches with around 1.7% bias and nearly 94% coverage.

We observed the same pattern when we compared the point estimates by simulated scenarios (Figure 5). As expected, unadjusted estimates ranged from approximately 0.8 to over 1.2. Both LASSO approaches followed by the autoencoder approaches demonstrated the best adjustment performance in most of the cases. In particular, we observed that the LASSO EC model had the best CI coverage to include the true HR in at least 95% of the times in 14 out of the 27 simulated scenarios (Figure 5 and eTables 8–10; <http://links.lww.com/EDE/B777>). When the percent bias was compared by simulated scenario, the results were consistent with <2% (LASSO and LASSO EC) and 3% (autoencoder and autoencoder EC) bias in almost all of the scenarios (Figure 6 and eTable 9; <http://links.lww.com/EDE/B777>).

Sensitivity Analyses

When we changed the autoencoder architecture from three hidden layers to one in sensitivity analysis I, the performance of the autoencoder-based models slightly improved (eTable 11; <http://links.lww.com/EDE/B777>). The overall performance remained nearly the same when the main architecture was altered to having a 128-dimensional bottleneck layer size in sensitivity analysis II (eTable 12; <http://links.lww.com/EDE/B777>). Combining the architecture alterations from sensitivity analyses I and II, results of the autoencoder approach were comparable to the ones of the LASSO approaches with an average RMSE of 0.0203 (eTable 13; <http://links.lww.com/EDE/B777>). When taking all possible PCs, instead of those

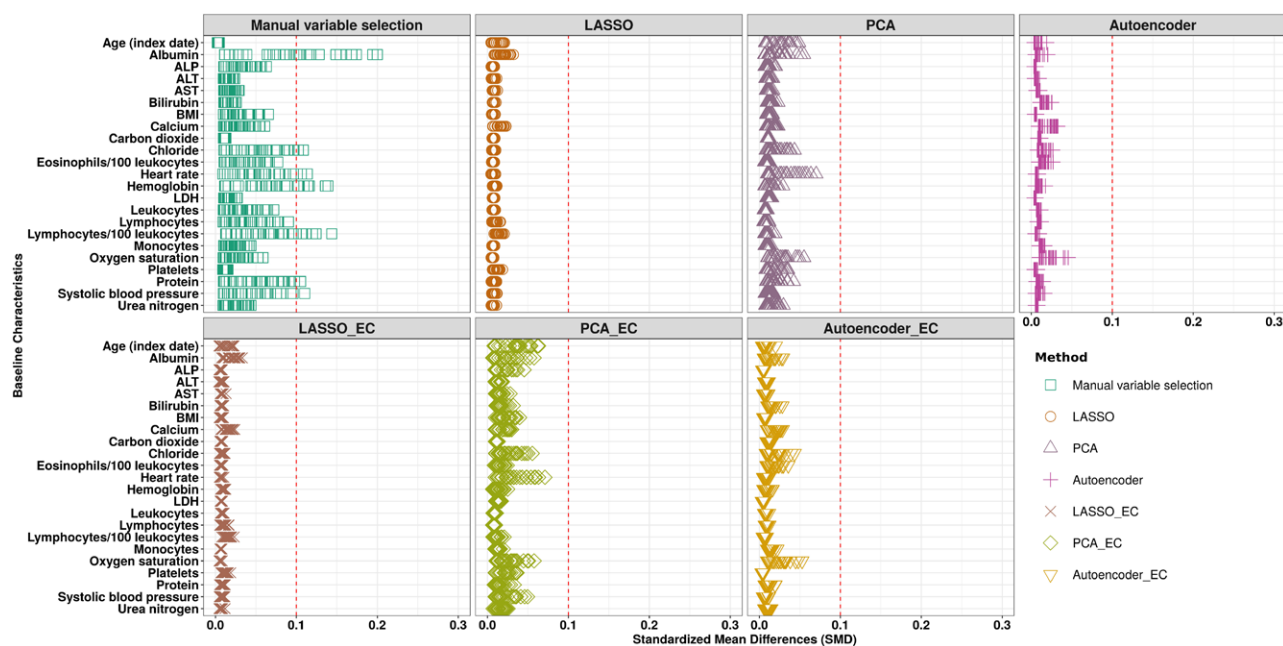


FIGURE 4. Baseline covariate balance by propensity score computation method and simulation scenario. Average SMDs are displayed for each of the 27 scenarios per baseline characteristic. ALP indicates alkaline phosphatase; ALT, alanine aminotransferase; AST, aspartate aminotransferase; BMI, body mass index; ECOG, Eastern Cooperative Oncology Group (ECOG) Performance Status; LASSO, least absolute shrinkage and selection operator; LDH, lactate dehydrogenase; NLR, neutrophil-to-lymphocyte ratio; PC(A), principal component (analysis). Figure is available in color online.

TABLE 2. Summary of Adjustment Performance Across All Scenarios

Method	RMSE	Bias (%)	CI Coverage (%)
Unadjusted	0.1205	10.4	16.41
Multivariable regression	0.0790	6.75	27.67
Manual variable selection	0.0670	5.73	32.81
LASSO	0.0205	1.65	93.74
PCA	0.0293	2.39	79.59
Autoencoder	0.0248	2.00	87.70
LASSO EC	0.0210	1.69	93.52
PCA EC	0.0329	2.71	74.00
Autoencoder EC	0.0265	2.15	85.19

LASSO indicates least absolute shrinkage and selection operator; PC(A), principal component (analysis).

describing 80% of the cumulative variance explained, the performance according to RMSE and bias decreased while the coverage improved (eTable 14; <http://links.lww.com/EDE/B777>). Increasing the number of replications to 500 did not noticeably change the results of the main analysis, indicating that 100 replications per scenario were sufficient (eTable 15; <http://links.lww.com/EDE/B777>).

Case Study

There were 781 patients eligible for the case study (eFigure 16; <http://links.lww.com/EDE/B777>). The results are summarized in Figure 7. All analyses suggested a null

association with the unadjusted point estimate being slightly below the null. All adjusted models ranged between point estimates of 1.00–1.09 with the autoencoder analysis being slightly closer to the null ($HR_{\text{autoencoder}} = 1.01$ [95% CI = 0.80, 1.27] vs. $HR_{\text{PRONOUNCE}} = 1.07$ [0.83, 1.36]) as compared with the autoencoder EC model ($HR = 1.09$ [95% CI = 0.87, 1.37]). SMR weighting led to very similar estimates with the exception of the LASSO approaches having much wider CIs (eFigure 17; <http://links.lww.com/EDE/B777>).

DISCUSSION

In this RWD study, we developed a novel automated autoencoder-based approach and compared it with established approaches. Using a comprehensive simulation framework, we observed that in terms of confounding control, the autoencoder-based approach led to reasonable results, but did not perform substantially better than some of the established approaches such as LASSO. In an empirical case study emulating the PRONOUNCE trial using observational data, the autoencoder-based results were consistent with the conclusion of the original trial.

PSs are frequently used analytical tools (eFigure 18; <http://links.lww.com/EDE/B777>) since they enable researchers to collapse many dimensions of confounding covariates into a single dimension while still maintaining sufficient precision. The advantage of deep learning-based PS is the ability to easily handle large amounts of data involving complex

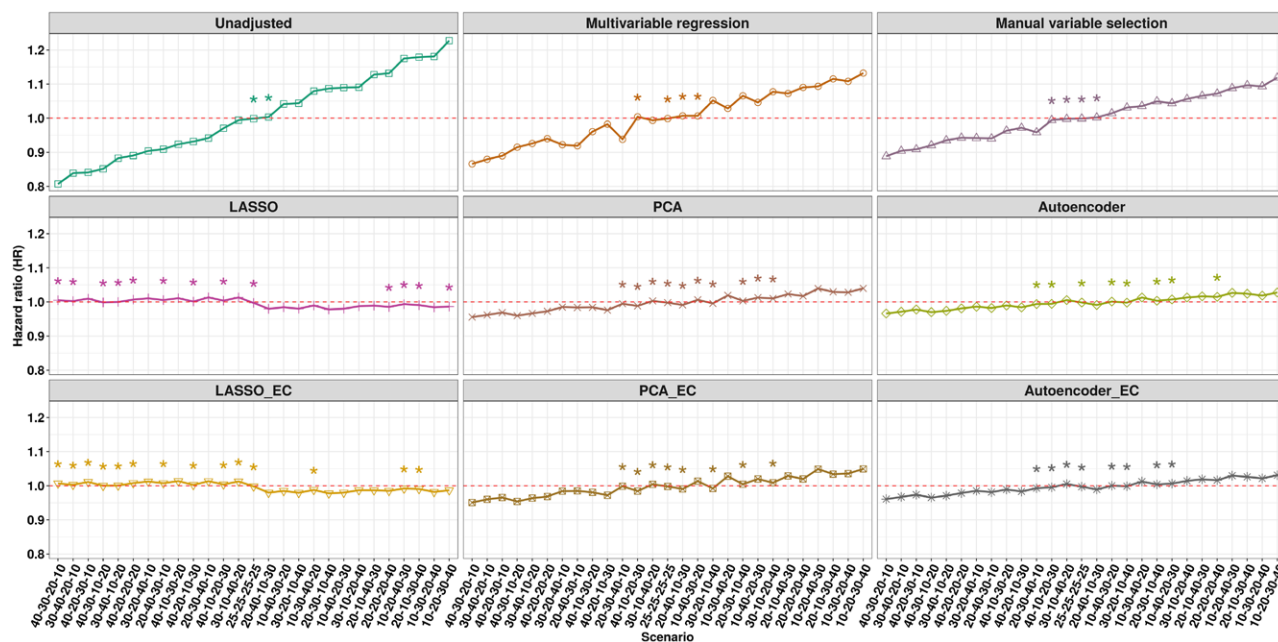


FIGURE 5. Average HRs for each of the 27 simulated scenarios and paneled PS estimation method. *Indicates that the CI coverage for the respective scenario included the true HR at least 95% of the times. The red dashed line indicates the true HR that is intended to be recovered by the propensity score adjustment. LASSO indicates least absolute shrinkage and selection operator. Figure is available in color online.

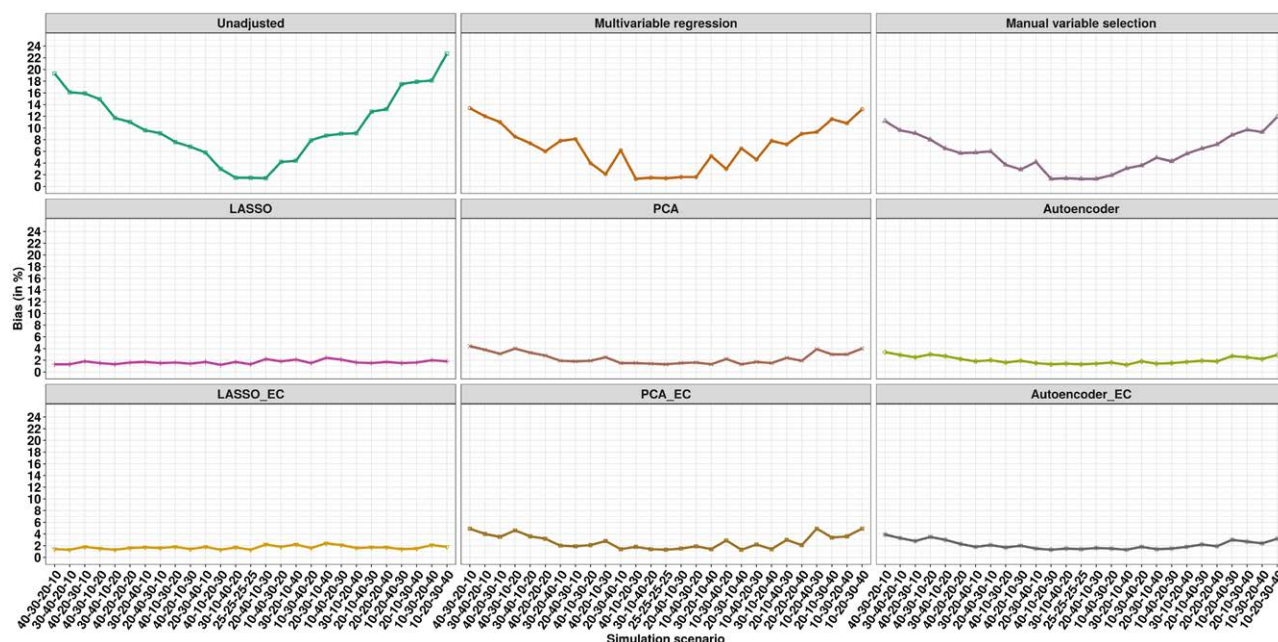


FIGURE 6. Comparison of average absolute % bias by simulation scenario for each PS estimation method. LASSO indicates least absolute shrinkage and selection operator. Figure is available in color online.

associations between covariates. An earlier study from 2008 investigated different techniques in PS estimation with various non-linear and non-additive associations on 10 binary/continuous covariates and concluded that even a rather simple neural network outperformed recursive partitioning algorithms in terms of providing the least numerically biased estimates.³⁴

This may suggest that the appropriate modeling of potentially non-linear covariate structures may be of relevant importance for confounding control. Especially analyses in EHR data may benefit from autoencoder-based PS as these usually capture routine care laboratory measurements and vital sign parameters which have been shown to be of paramount prognostic

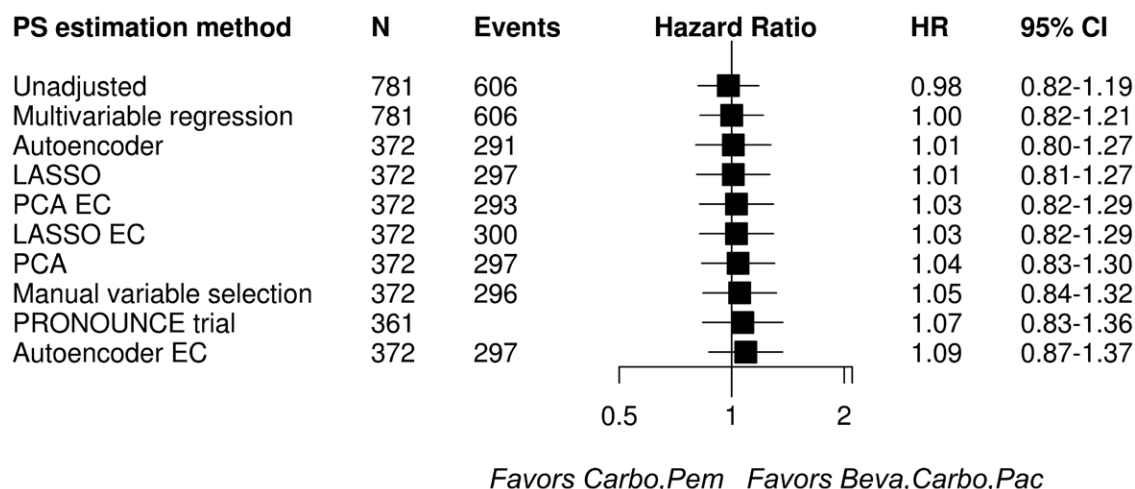


FIGURE 7. Forest plot illustrating HRs and 95% confidence intervals (CIs) for overall survival by PS estimation method. HR indicates hazard ratio; LASSO, least absolute shrinkage and selection operator; PCA, principal component analysis.

value.^{21,35} This may explain why in this study the autoencoder-based PS performed better than the PCA approaches since, in case of no nonlinearity, both methods should in principle lead to similar results.³⁶ However, given that continuous covariates are still usually rather rare in healthcare databases, we may have underestimated the abilities of the autoencoder-based approach in this study and further studies are warranted once multimodal data elements, such as medical images and sequencing data, complement contemporary databases.³⁷

Application and Use Cases

The autoencoder-based PS can be generally used in any type of comparative effectiveness study where sufficient confounder balancing between two cohorts is required. In the here presented comparative effectiveness case study, it was possible to derive the same qualitative conclusion as in the PRONOUNCE trial by applying autoencoder-based PS. Although the primary objective in the case study was to test the use of autoencoder-based PS in a real comparative effectiveness research setup, the equal results of all methods may be explained with the fact that confounding bias was not as strong in this particular research question as compared to some of the more extreme scenarios in the simulation. This seems plausible given that due to the variety of possible treatments and sometimes lacking evidence for the most effective combination, the selective channeling of patients with higher risk (as often observed with prescription drugs like COX-2 inhibitors vs. nonselective nonsteroidal anti-inflammatory drugs) may not be apparent. This may underline an attractive feature of the autoencoder-based PS, which could be used as an automated and data-adaptive sensitivity analysis in comparative effectiveness studies with unknown extent of confounding bias.

Especially in the era of precision medicine, in which treatment decisions for specific subpopulations of patients are based on distinct molecular characteristics, comparative

effectiveness research might play an increasingly important role in addressing challenges, for example, in the area of early clinical development of new therapeutics. Here, designs such as external control arms are interesting approaches that could benefit by advanced analytics like deep learning-based PS. A recent proof-of-concept study assessed how well external controls could have approximated the actual standard-of-care controls in nine lung cancer trials.³⁸ The authors reported that the comparison of estimates between RCTs and external controls resulted in a Pearson correlation coefficient of 0.86. This is an encouraging example suggesting that external control arms come with a sufficient validity and can play an important role in facilitating real-world data to support early clinical development and regulatory submissions.^{3,39}

Strengths and Limitations

Due to the nature of routinely collected health records, there is missing data. In this study, we employed median imputation and assigning a missing-indicator category to one-hot encoded categorical variables because this or similar approaches were suggested to have good performance in studies with large datasets where multiple imputation would be computationally very expensive and generally not operationalizable.^{16,40} This approach is also supported by various recent prediction models trained on EHR data which reported outstanding performance.^{15,41}

In addition, data-adaptive approaches always inherit the risk of including covariates that may be collider covariates (M-bias), instrumental covariates (Z-bias), or causal intermediates. Colliders are covariates that open a causal path from exposure to outcome.⁴² Including such covariates in the PS computation may induce a spurious association where in fact there is none. Besides directed acyclic graphs, there is no formal way to test for colliders, making it difficult to exclude such variables before PS computation. However, Schneeweiss⁴³ found that under realistic scenarios a collider-induced bias was negligible

and outweighed by the adjustment effect for other covariates. Instrumental variables (IV) are covariates that are only associated with the exposure but not with the outcome. IVs are frequently used to control for unmeasured confounding⁴⁴ but also introduce bias (Z-bias) when conditioning on them. Especially in oncology, calendar period effects are strong predictors for therapy decisions once new breakthrough treatments are approved.⁴⁵ Although there is a theoretical chance to have unintentionally included IVs, Myers et al⁴⁶ showed that only in the presence of strong unmeasured confounding does Z-bias have effects worth mentioning. While expert knowledge plays an important role in avoiding covariates that could mediate the association between exposure and outcome,⁴⁷ the risk of adjusting for causal intermediates can also be mitigated with appropriate study designs such as an active comparator, new user design as applied in this study.^{4,43}

A unique strength of this study is the novelty approach to learn patient representations for PS computation in a data-adaptive manner, which we found to have a reasonable performance and which may serve as a promising tool for the future once more data elements complement contemporary databases. Applying comprehensive sensitivity analyses, we found the methodology to be robust as all setups and scenarios resulted in a similar conclusion. The observation that the autoencoder architecture with less hidden layers and a larger bottleneck layer led to results closer to LASSO gave some concern that this may have been the consequence of overfitting of the main model. Nevertheless, differences were marginal and did not change the main conclusion while the hyperparameter setup of the main model was found to be a reasonable trade-off between compactness of the resulting embedding and sufficient reconstruction performance and generalizability.

It is further important to credit that the autoencoder approach is a pure unsupervised method, which means that the confounding control in this study has been solely achieved without optimizing the network towards the probability of patients receiving the treatment, which needs to be acknowledged when comparing to supervised approaches like LASSO. Hence, potential deep-learning architectural extensions would be of utmost interest, for example, by jointly modeling targets and inputs using end-to-end learning architectures.

A limitation of this simulation is that due to the non-collapseability of HRs, only a null treatment effect could be simulated which may in future research be addressed by estimating risk-differences and more sophisticated simulation techniques such as plasmode simulations.^{48,49} In addition, variance estimation seemed to be less reliable for models 1–3 (eTable 7; <http://links.lww.com/EDE/B777>), limiting the ability to make final conclusions about their true CI coverage.

For this study, it was possible to use a large underlying population to train and empirically examine the comparative performance of the proposed autoencoder approach. This real-world database provided comprehensive oncology-specific data, which underwent a rigorous data quality assurance process before release.

Finally, it is important to acknowledge that this study primarily focused on the analytical aspects to reduce confounding. Carefully chosen analysis always needs to go along with a causal study design to avoid serious biases such as reverse causality and immortal time bias, which are known as sources for much larger bias than conventional confounding.^{43,50}

CONCLUSIONS

In summary, we developed an autoencoder-based PS computation that our assessment found to be a feasible approach to reduce confounding bias, although not with a substantially stronger performance than some of the established approaches such as LASSO. As a promising tool for the future, it may be considered alongside with established approaches in non-randomized comparisons in comparative effectiveness research.

ACKNOWLEDGMENTS

We would like to thank all staff members of Flatiron Health for their dedicated data collection and curation.

REFERENCES

1. Basch E, Schrag D. The evolving uses of “Real-World” data. *JAMA*. 2019;321:1359–1360.
2. Corrigan-Curay J, Sacks L, Woodcock J. Real-world evidence and real-world data for evaluating drug safety and effectiveness. *JAMA*. 2018;320:867–868.
3. U.S. Food and Drug Administration. Framework for FDA’s Real-World Evidence Program. Published online December 2018. Available at: <https://www.fda.gov/media/120060/download>. Accessed 7 January 2020.
4. Lund JL, Richardson DB, Stürmer T. The active comparator, new user study design in pharmacoepidemiology: historical foundations and contemporary application. *Curr Epidemiol Rep*. 2015;2:221–228.
5. Rosenbaum PR, Rubin DB. The central role of the propensity score in observational studies for causal effects. *Biometrika*. 1983;70:41–55.
6. Austin PC. An introduction to propensity score methods for reducing the effects of confounding in observational studies. *Multivariate Behav Res*. 2011;46:399–424.
7. Westreich D, Cole SR. Invited commentary: positivity in practice. *Am J Epidemiol*. 2010;171:674–677; discussion 678.
8. Brookhart MA, Schneeweiss S, Rothman KJ, Glynn RJ, Avorn J, Stürmer T. Variable selection for propensity score models. *Am J Epidemiol*. 2006;163:1149–1156.
9. Schneeweiss S, Eddings W, Glynn RJ, Paterno E, Rassen J, Franklin JM. Variable selection for confounding adjustment in high-dimensional covariate spaces when analyzing healthcare databases. *Epidemiology*. 2017;28:237–248.
10. Stürmer T, Joshi M, Glynn RJ, Avorn J, Rothman KJ, Schneeweiss S. A review of the application of propensity score methods yielded increasing use, advantages in specific settings, but not substantially different estimates compared with conventional multivariable methods. *J Clin Epidemiol*. 2006;59:437–447.
11. Glynn RJ, Schneeweiss S, Stürmer T. Indications for propensity scores and review of their use in pharmacoepidemiology. *Basic Clin Pharmacol Toxicol*. 2006;98:253–259.
12. Goodfellow I, Bengio Y, Courville A, Bengio Y. *Deep Learning*. Vol 1. MIT press Cambridge; 2016.
13. Zinner RG, Obasaju CK, Spigel DR, et al. PRONOUNCE: randomized, open-label, phase III study of first-line pemetrexed + carboplatin followed by maintenance pemetrexed versus paclitaxel + carboplatin + bevacizumab followed by maintenance bevacizumab in patients with advanced nonsquamous non-small-cell lung cancer. *J Thorac Oncol*. 2015;10:134–142.
14. Tukey JW. *Exploratory Data Analysis*. Mass Addison-Wesley; 1977.

15. Tomašev N, Glorot X, Rae JW, et al. A clinically applicable approach to continuous prediction of future acute kidney injury. *Nature*. 2019;572:116–119.
16. Eraslan G, Simon LM, Mircea M, Mueller NS, Theis FJ. Single-cell RNA-seq denoising using a deep count autoencoder. *Nat Commun*. 2019;10:390.
17. Schneeweiss S, Rassen JA, Glynn RJ, Avorn J, Mogun H, Brookhart MA. High-dimensional propensity score adjustment in studies of treatment effects using health care claims data. *Epidemiology*. 2009;20:512–522.
18. Czwikla J, Jobski K, Schink T. The impact of the lookback period and definition of confirmatory events on the identification of incident cancer cases in administrative data. *BMC Med Res Methodol*. 2017;17:122.
19. Schneeweiss S, Rassen JA, Brown JS, et al. Graphical depiction of longitudinal study designs in health care databases. *Ann Intern Med*. 2019;170:398–406.
20. Bi Q, Goodman KE, Kaminsky J, Lessler J. What is machine learning? A primer for the epidemiologist. *Am J Epidemiol*. 2019;188:2222–2239.
21. Becker T, Weberpals J, Jegg AM, et al. An enhanced prognostic score for overall survival of patients with cancer derived from a large real world cohort. *Ann Oncol*. 2020;31:1561–1568.
22. Stuart EA, Lee BK, Leacy FP. Prognostic score-based balance measures can be a useful diagnostic for propensity score methods in comparative effectiveness research. *J Clin Epidemiol*. 2013;66(8 suppl):S84–S90.e1.
23. Austin PC. A comparison of 12 algorithms for matching on the propensity score. *Stat Med*. 2014;33:1057–1069.
24. Lin DY, Wei LJ. The robust inference for the Cox proportional hazards model. *J Am Stat Assoc*. 1989;84:1074–1078.
25. Greenland S, Robins JM, Pearl J. Confounding and collapsibility in causal inference. *Stat Sci*. 1999;14:29–46.
26. Stensrud MJ, Hernán MA. Why test for proportional hazards? *JAMA*. 2020;323:1401–1402.
27. Austin PC. Balance diagnostics for comparing the distribution of baseline covariates between treatment groups in propensity-score matched samples. *Stat Med*. 2009;28:3083–3107.
28. Burton A, Altman DG, Royston P, Holder RL. The design of simulation studies in medical statistics. *Stat Med*. 2006;25:4279–4292.
29. Morris TP, White IR, Crowther MJ. Using simulation studies to evaluate statistical methods. *Stat Med*. 2019;38:2074–2102.
30. Desai RJ, Wyss R, Abdia Y, et al. Evaluating the use of bootstrapping in cohort studies conducted with 1:1 propensity score matching-A plasmode simulation study. *Pharmacoepidemiol Drug Saf*. 2019;28:879–886.
31. Hernán MA, Robins JM. Using big data to emulate a target trial when a randomized trial is not available. *Am J Epidemiol*. 2016;183:758–764.
32. Ho D, Imai K, King G, Stuart EA. MatchIt: nonparametric preprocessing for parametric causal inference. *J Stat Softw*. 2011;42:1–28.
33. Greifer N. *WeightIt: Weighting for Covariate Balance in Observational Studies (R Package Version 0.5.1)*. 2019. Available at: <https://cran.r-project.org/web/packages/WeightIt/WeightIt.pdf>. Accessed 16 January 2019.
34. Setoguchi S, Schneeweiss S, Brookhart MA, Glynn RJ, Cook EF. Evaluating uses of data mining techniques in propensity score estimation: a simulation study. *Pharmacoepidemiol Drug Saf*. 2008;17:546–555.
35. Garrido-Laguna I, Janku F, Vaklavas C, et al. Validation of the Royal Marsden Hospital prognostic score in patients treated in the phase I clinical trials program at the MD Anderson Cancer Center. *Cancer*. 2012;118:1422–1428.
36. Miotto R, Li L, Kidd BA, Dudley JT. Deep patient: an unsupervised representation to predict the future of patients from the electronic health records. *Sci Rep*. 2016;6:26094.
37. Binder H. Big data und deep learning in der onkologie. *Onkol*. 2018;24:361–367.
38. Carrigan G, Whipple S, Capra WB, et al. Using electronic health records to derive control arms for early phase single-arm lung cancer trials: proof-of-concept in randomized controlled trials. *Clin Pharmacol Ther*. 2020;107:369–377.
39. Burcu M, Dreyer NA, Franklin JM, et al. Real-world evidence to support regulatory decision-making for medicines: considerations for external control arms. *Pharmacoepidemiol Drug Saf*. 2020;29:1228–1235.
40. Rodenburg FJ, Sawada Y, Hayashi N. Improving RNN performance by modelling informative missingness with combined indicators. *Appl Sci*. 2019;9:1623.
41. Rajkomar A, Oren E, Chen K, et al. Scalable and accurate deep learning with electronic health records. *Npj Digit Med*. 2018;1:1–10.
42. Liu W, Brookhart MA, Schneeweiss S, Mi X, Setoguchi S. Implications of M bias in epidemiologic studies: a simulation study. *Am J Epidemiol*. 2012;176:938–948.
43. Schneeweiss S. Automated data-adaptive analytics for electronic healthcare data to study causal treatment effects. *Clin Epidemiol*. 2018;10:771–788.
44. Brookhart MA, Rassen JA, Schneeweiss S. Instrumental variable methods in comparative safety and effectiveness research. *Pharmacoepidemiol Drug Saf*. 2010;19:537–554.
45. Mack CD, Brookhart MA, Glynn RJ, Stürmer T. Calendar time as an instrumental variable in nonexperimental comparative effectiveness research of emerging therapies. *Value Health*. 2013;16:A129–A130.
46. Myers JA, Rassen JA, Gagne JJ, et al. Effects of adjusting for instrumental variables on bias and precision of effect estimates. *Am J Epidemiol*. 2011;174:1213–1222.
47. Corraíni P, Olsen M, Pedersen L, Dekkers OM, Vandenbroucke JP. Effect modification, interaction and mediation: an overview of the theoretical insights for clinical investigators. *Clin Epidemiol*. 2017;9:331–338.
48. Franklin JM, Schneeweiss S, Polinski JM, Rassen JA. Plasmode simulation for the evaluation of pharmacoepidemiologic methods in complex healthcare databases. *Comput Stat Data Anal*. 2014;72:219–226.
49. Huitfeldt A, Stensrud MJ, Suzuki E. On the collapsibility of measures of effect in the counterfactual causal framework. *Emerg Themes Epidemiol*. 2019;16:1.
50. Weberpals J, Jansen L, van Herk-Sukel MPP, et al. Immortal time bias in pharmacoepidemiological studies on cancer patient survival: empirical illustration for beta-blocker use in four cancers with different prognosis. *Eur J Epidemiol*. 2017;32:1019–1031.