

---

# Covariate Shift Corrected Conditional Randomization Test

---

**Bowen Xu\***

Harvard University  
bowenxu@g.harvard.edu

**Yiwen Huang\***

Department of Statistics  
Peking University  
2000010773@stu.pku.edu.cn

**Chuan Hong**

Department of Biostatistics and Bioinformatics  
Duke University  
chuan.hong@duke.edu

**Shuangning Li**

Booth School of Business  
University of Chicago  
shuangning.li@chicagobooth.edu

**Molei Liu<sup>†</sup>**

Department of Biostatistics  
Columbia Mailman School of Public Health  
m14890@cumc.columbia.edu

## Abstract

1 Conditional independence tests are crucial across various disciplines in determining  
2 the independence of an outcome variable  $Y$  from a treatment variable  $X$ , condition-  
3 ing on a set of confounders  $Z$ . The Conditional Randomization Test (CRT) offers  
4 a powerful framework for such testing by assuming known distributions of  $X | Z$ ;  
5 it controls the Type-I error exactly, allowing for the use of flexible, black-box  
6 test statistics. In practice, testing for conditional independence often involves  
7 using data from a source population to draw conclusions about a target population.  
8 This can be challenging due to covariate shift—differences in the distribution of  
9  $X$ ,  $Z$ , and surrogate variables, which can affect the conditional distribution of  
10  $Y | X, Z$ —rendering traditional CRT approaches invalid. To address this issue,  
11 we propose a novel Covariate Shift Corrected Pearson Chi-squared Conditional  
12 Randomization (csPCR) test. This test adapts to covariate shifts by integrating  
13 importance weights and employing the control variates method to reduce variance  
14 in the test statistics and thus enhance power. Theoretically, we establish that the  
15 csPCR test controls the Type-I error asymptotically. Empirically, through simu-  
16 lation studies, we demonstrate that our method not only maintains control over  
17 Type-I errors but also exhibits superior power, confirming its efficacy and practical  
18 utility in real-world scenarios where covariate shifts are prevalent. Finally, we  
19 apply our methodology to a real-world dataset to assess the impact of a COVID-19  
20 treatment on the 90-day mortality rate among patients.

## 21 1 Introduction

22 Conditional independence tests are important across diverse fields for determining whether an  
23 outcome variable  $Y$  is independent of a treatment variable  $X$ , conditioning on a potentially high-

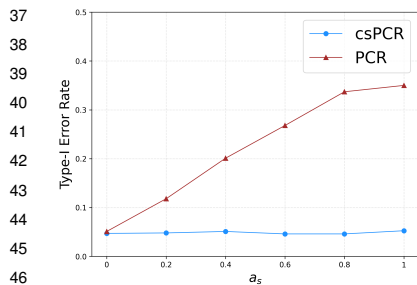
---

\*These authors contributed equally to this work.

<sup>†</sup>Corresponding author. To whom correspondence should be addressed.

24 dimensional vector of confounding variables  $Z$ . This type of testing is critical for understanding the  
 25 complex relationships among variables. For instance, scientists may hope to understand whether  
 26 a specific genetic feature influences disease outcomes, whether a particular treatment effectively  
 27 extends life expectancy, or whether certain demographic factors impact college admissions.

28 Traditionally, these conditional testing problems are approached by modeling  $Y$  against  $X$  and  $Z$   
 29 through some parametric or semiparametric model. However, this strategy has been criticized due  
 30 to potential model misspecification and limited observations of  $Y$ . As an alternative strategy, the  
 31 model-X framework and Conditional Randomization Test (CRT) propose testing for the general  
 32 conditional independence hypothesis  $H_0 : X \perp\!\!\!\perp Y \mid Z$ , free of any specific effect parameters [2].  
 33 The CRT assumes the distribution of  $X \mid Z$  to be known and can control the type-I error exactly,  
 34 allowing for the choice of any flexible, black-box test statistic. This strategy is particularly useful  
 35 when there is either strong and reliable scientific knowledge of the distribution of  $X \mid Z$  or an  
 36 auxiliary dataset of  $(X, Z)$  of large sample size, known as the semi-supervised setting.



37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
Figure 1: Type-I Error rates of our proposed csPCR and the source-only PCR on a simulated example. The Type-I error inflation of PCR demonstrates that source analysis is not valid or generalizable on the target due to covariate shift.

In practice, testing for conditional independence frequently involves using data from a source population to draw conclusions about a target population. This situation presents challenges due to potential differences in the distribution of variables between the two populations. For example, economists may be interested in whether college admission ( $Y$ ) is independent of family income ( $X$ ), conditioning on variables such as GPA, extracurricular activities, geographic location, and other demographics ( $Z$ ). In the source population, the relationship might be influenced by factors like wealthy parents investing in SAT preparation, which boosts admission rates—a relationship that may not exist in a target population where such preparation is less common. Although  $Y$  may not appear independent of  $X$  given  $Z$  in the source population, the conclusion could vary significantly in the target population. This discrepancy underscores the need for a robust and flexible testing procedure that can adapt to shifts in distributions.

54 More specifically, we address the *covariate shift* scenario, where the distributions of the treatment  
 55 variables  $X$ , the confounding variables  $Z$ , and some surrogate or auxiliary variables  $V$  (e.g., SAT  
 56 scores) may differ between the source and target populations. However, the conditional distribution  
 57 of  $Y$  given  $X, Z$ , and  $V$  remains the same between them. In such scenarios, our goal is to leverage  
 58 information from the source to accurately test for conditional independence in the target population  
 59 without the observation of  $Y$  on target. In the scenario we consider, the presence of  $V$  and potential  
 60 differences in  $P(V \mid X, Z)$  between the source and target populations may lead to the conditional  
 61 independence  $X \perp\!\!\!\perp Y \mid Z$  not holding simultaneously in the two populations. Specifically, because

$$P(Y \mid X, Z) = \int P(Y \mid X, Z, V)P(V \mid X, Z) dV,$$

62 the conditional distribution of  $Y$  given  $X$  and  $Z$  can vary between populations. This underscores  
 63 why the problem is non-trivial.

64 See Figure 1 for an example of the consequences of such covariate shift.

65 In this paper, we propose a novel conditional independence test suitable for covariate shift scenarios.  
 66 Our method builds upon the Pearson Chi-Squared Conditional Randomization (PCR) test, a powerful  
 67 model-X testing procedure that effectively addresses a broader range of alternative  $p$ -value distri-  
 68 butions than the vanilla CRT [5]. Methodologically, we make two major contributions. First, we  
 69 introduce importance weights into the label counting steps of the original PCR test, making the new  
 70 test valid under covariate shift. These weights adjust the importance of each sample according to  
 71 its density ratio, effectively rebalancing the source data to match the target population’s distribution.  
 72 Second, we introduce a power enhancement method that employs the control variates method to  
 73 reduce variance in the test statistics. Although importance weights can increase the variance in test  
 74 statistics, especially when the density ratio can become extremely high, potentially reducing power,  
 75 our power enhancement method effectively addresses this issue. Together, these innovations enable  
 76 us to develop a PCR test that is both powerful and valid under covariate shifts.

77 The rest of the paper is organized as follows: In Section 2, we provide a formal introduction to  
78 the problem setup. In Section 3, we introduce the proposed Covariate Shift Corrected Pearson  
79 Chi-squared Conditional Randomization (csPCR) test and establish that the proposed csPCR test  
80 controls the Type-I error asymptotically. In Section 4, we demonstrate the empirical performance  
81 of the csPCR test through simulation studies. In Section 5, we apply the proposed csPCR test to a  
82 real-world dataset to assess the impact of a COVID-19 treatment on the 90-day mortality rate among  
83 patients.

## 84 1.1 Related Work

85 Our work builds upon the model-X framework and the conditional randomization test proposed by  
86 Candes et al. [2]. The particular method we develop is based on a variant of the vanilla CRT, the  
87 Pearson Conditional Randomization (PCR) test [5]. Recent advances in the CRT include improving  
88 computation time [7, 10], studying robustness [6, 11], and examining statistical power [19]. The  
89 focus of this paper, different from the above, is on how to build a valid CRT procedure when there  
90 is covariate shift. The paper is also complementary to the above literature: for example, we hope  
91 that future work can conduct theoretical power analysis for our procedure or develop a double robust  
92 version of the procedure just like in [6]. Finally, we note that surrogate variables play a crucial role in  
93 this paper: because the distribution of the surrogate variables is different in the source and the target  
94 population, naively testing the conditional independence hypothesis in the source population can  
95 yield invalid conclusions for the target population. A surrogate or silver standard label is a variable  
96 that is more feasible and accessible than  $Y$  in data collection and can be viewed as a noisy measure  
97 of  $Y$ . For example, tumor response rate is often used as an early endpoint surrogate for the long-term  
98 survival outcome [3], and blood pressure is commonly used as a surrogate for heart attacks. Surrogate  
99 variables are also commonly used in environmental studies and economics. Surrogate variables also  
100 play an important role in the paper by [6], albeit in a different way, where the surrogate variables are  
101 used to learn the distribution of  $Y \mid X, Z$  and to further improve the robustness of the CRT procedure.

102 Statistical learning and inference under covariate shift has been extensively studied over the past years.  
103 As a seminal work in addressing covariate shift bias, [4] proposed a density ratio weighting approach  
104 using kernel mean matching to characterize the adjusting weights. Their key idea of importance  
105 (re)weighting is intrinsically connected with early work in broader contexts like importance sampling  
106 [15, e.g.] and semiparametric inference [13, e.g.]. [8] extended this idea to a doubly robust framework  
107 accommodating surrogate variables like  $V$  and being more robust to the misspecification or poor  
108 quality of the density ratio models. [18] handled a more challenging scenario with severe shift and  
109 poor overlap between the source and target populations. Among this track of literature, [17] is the  
110 most closely related to our work as they also considered conditional independence testing under  
111 distributional shifts and proposed a general testing procedure base on importance sampling (IS)  
112 allowing for the use of CRT. Different from us, their work does not accommodate the covariate shifts  
113 of some surrogate or auxiliary  $V$ . Moreover, as will be shown in our numerical studies, their general  
114 IS testing strategy can encounter the loss of effective sample sizes and be less powerful than ours.

## 115 2 Problem Setup

### 116 2.1 Conditional Independence Testing under Covariate Shift

117 Let  $Y \in \mathbb{R}$  denote the outcome variable,  $X \in \mathbb{R}$  the treatment variable,  $Z \in \mathbb{R}^p$  a vector of  
118 confounding variables, and  $V \in \mathbb{R}^d$  a vector of surrogate variables. To make the problem more  
119 concrete, consider the following two examples:

120 **Example 1** (College Admission).  $Y$  is college admission,  $X$  is family income,  $Z$  includes a number  
121 of factors such as GPA, extracurricular activities, geographic location, and demographic information,  
122  $V$  is the SAT score. In this case,  $V$  is easier to collect compared to  $Y$  as the college admission  
123 requires individual-level surveys.

124 **Example 2** (Health Outcome).  $Y$  is a long-term health outcome,  $X$  is a medical treatment,  $Z$   
125 includes factors such as age, gender, and health history,  $V$  includes surrogate variables like blood  
126 pressure, BMI, and duration of hospital stays post the treatment, which can be measured within a  
127 much shorter term than  $Y$ .

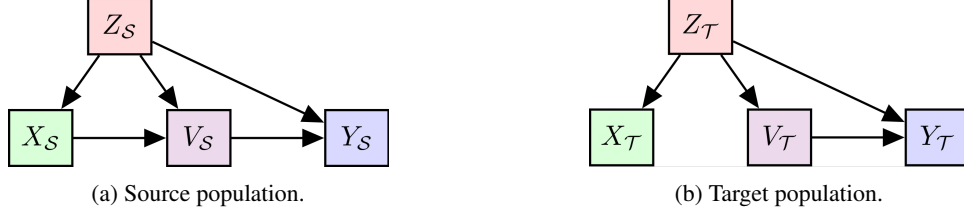


Figure 2: Direct acyclic graphs illustrating possible differences between the source and the target populations.

128 Consider a scenario involving two distinct populations: the source population  $\mathcal{S}$  and the target  
 129 population  $\mathcal{T}$ . We collect data from the source population with the goal of making inferences about  
 130 the target population. The source data contains  $n$  independent and identically distributed samples of  
 131  $(Y_i, X_i, Z_i, V_i)$  for  $i = 1, \dots, n$ . Let  $\mathbf{y} = (Y_1, Y_2, \dots, Y_n)^\top \in \mathbb{R}^n$ ,  $\mathbf{x} = (X_1, X_2, \dots, X_n)^\top \in \mathbb{R}^n$ ,  
 132  $\mathbf{Z} = (Z_1, Z_2, \dots, Z_n)^\top \in \mathbb{R}^{n \times p}$ , and  $\mathbf{V} = (V_1, V_2, \dots, V_n)^\top \in \mathbb{R}^{n \times d}$ . We are interested in  
 133 testing the following conditional independence hypothesis in the target population:

$$\mathcal{H}_0 : X \perp\!\!\!\perp Y \mid Z. \quad (1)$$

134 We assume that the conditional distribution of  $Y \mid X, Z, V$  is the same in both populations; however,  
 135 the distribution of  $(X, Z, V)$  varies between  $\mathcal{S}$  and  $\mathcal{T}$ . More precisely, the joint distribution of  
 136  $Y, X, Z, V$  can be described as follows:

$$\begin{aligned} P_{\mathcal{S}}(Y, X, Z, V) &= P_{\mathcal{S}}(X, Z, V)P(Y|X, Z, V) \quad \text{on } \mathcal{S}, \\ P_{\mathcal{T}}(Y, X, Z, V) &= P_{\mathcal{T}}(X, Z, V)P(Y|X, Z, V) \quad \text{on } \mathcal{T}. \end{aligned} \quad (2)$$

137 This situation is referred to as the *covariate shift* scenario because the distribution of the covariates  
 138  $X, Z$ , and  $V$  in the source population  $\mathcal{S}$  does not match that in the target population  $\mathcal{T}$ .

139 Let's understand the above assumption and its implications through the two examples above. In the  
 140 college admissions example, it is plausible to assume that the rate of college admissions remains  
 141 consistent across the two populations when conditioned on the SAT score, family income, and other  
 142 confounding variables. However, the joint distribution of  $X, V$  and  $Z$  can differ: in the source  
 143 population, if wealthy parents frequently invest in SAT preparation, boosting admission rates, this  
 144 relationship may not hold in a target population where such preparation is uncommon. In such cases,  
 145 it is thus possible that  $X \not\perp\!\!\!\perp Y \mid Z$  in the source population but  $X \perp\!\!\!\perp Y \mid Z$  in the target population  
 146 (see Figure 2 for such an example). In the health outcomes example, it is again plausible that the  
 147 conditional distribution of long-term health outcomes given the treatment variable, confounding  
 148 variables, and surrogates remains the same across the two populations. However, the assignment of  
 149 the treatment may depend differently on the surrogate variables across the two populations. Therefore,  
 150 it's possible that  $X \perp\!\!\!\perp Y \mid Z$  in one population, but not in the other.

151 In both examples, we can see that the result of naively applying a valid conditional independence test  
 152 on the source population cannot guarantee a valid conclusion for testing  $\mathcal{H}_0$  in the target population.  
 153 Therefore, we need to develop new tools for addressing covariate shifts in conditional independence  
 154 tests.

## 155 2.2 Model-X Framework

156 In this paper, we operate within the model-X framework, as described by Candes et al. [2], which  
 157 assumes that the joint distributions of covariates  $X, V, Z$  are perfectly known in both the source and  
 158 target populations. This framework is particularly suited for scenarios where: (1) there is substantial  
 159 prior domain knowledge about the covariates  $X, V$ , and  $Z$ , or (2) there is a significant amount of  
 160 unsupervised data for these covariates in both populations, in addition to  $n$  labeled observations in  
 161 the source population, characterizing a semi-supervised setting.

162 An example of the first scenario can be seen in genetics, where researchers have well-established  
 163 models for the joint distributions of single nucleotide polymorphisms (SNPs). For the second scenario,  
 164 consider our earlier example involving health outcomes. Here, the outcome variable  $Y$  represents a  
 165 long-term health outcome that is more costly or sensitive to measure compared to the shorter-term

166 variables  $X$ ,  $V$ , and  $Z$ . In such cases, the variables  $X$ ,  $V$ , and  $Z$  are typically easier and less costly to  
 167 collect, frequently resulting in a semi-supervised setting in these health-related studies.

### 168 **3 Method: Covariate Shift Corrected PCR Test**

#### 169 **3.1 Incorporating the Density Ratio into the PCR Test**

170 In Section 2.1, we discussed how naively applying conditional independence tests to the source  
 171 data cannot guarantee valid conclusions for the target population. To address this issue, we must  
 172 incorporate information about the differences between the two populations into our testing procedure.  
 173 In particular, we will make use of the density ratio defined as:

$$e(X, Z, V) = \frac{P_{\mathcal{T}}(X, Z, V)}{P_{\mathcal{S}}(X, Z, V)}. \quad (3)$$

174 This ratio measures the relative likelihood of observing each combination of variables  $(X, Z, V)$   
 175 in the target population compared to the source population. By reweighting the data points in the  
 176 source population using this density ratio, we effectively transform the source distribution to match  
 177 the distribution of the target population, thereby addressing the covariate shift problem.

178 More specifically, we build our method upon the recently proposed Pearson Chi-Squared Conditional  
 179 Randomization (PCR) test [5]. Compared to the vanilla CRT, the PCR test is designed to be more  
 180 powerful across a broader range of alternative  $p$ -value distributions. At a high level, the PCR test  
 181 assigns a label to each data point following a counterfeit sampling step and a subsequent score  
 182 computation step. Under the null hypothesis that  $X \perp\!\!\!\perp Y \mid Z$ , the distribution of these labels  
 183 should be uniform across all possible labels. The PCR test then rejects the null hypothesis if  
 184 the empirical distribution of the labels deviates significantly from uniformity, as determined by a  
 185 Pearson’s chi-squared test.

186 Under distributional shift, if the data points were sampled from the target population, then the  
 187 distribution of the labels would be uniform. However, since the data points are actually sampled from  
 188 the source population, they must be reweighted using the density ratio. More specifically, in the final  
 189 step of the PCR test, where the Pearson’s chi-squared test is applied, we consider not the count of  
 190 data points for each label, but the sum of the density ratios of the data points for each label instead.  
 191 Under the null hypothesis, each sum should approximate  $n/L$ , where  $L$  is the total number of labels.  
 192 Consequently, we modify the Pearson’s chi-squared test to determine whether these weighted sums  
 193 deviate significantly from  $n/L$ .

194 Based on the above intuition, we propose the Covariate Shift Corrected PCR (csPCR) Test, as outlined  
 195 in Algorithm 1.

196 In Algorithm 1, lines 1-7 correspond to those in the original PCR test. These lines initiate the test by  
 197 generating counterfeit samples  $\tilde{X}_j^{(m)}$ . Assuming the source and target populations were identical, under  
 198 the null hypothesis, the random variables  $(X_j, Y_j, Z_j), (\tilde{X}_j^{(1)}, Y_j, Z_j), \dots, (\tilde{X}_j^{(M)}, Y_j, Z_j)$  would  
 199 be exchangeable. Consequently, the rank  $R_j$  would be uniformly distributed over  $\{1, \dots, M + 1\}$  in  
 200 the absence of ties, leading to a uniform distribution of the labels as well.

201 Lines 8-10 in Algorithm 1 address the covariate shift by incorporating density ratios as importance  
 202 weights into  $W_j$ . Due to this redefinition of  $W_\ell$ , the null distribution of the final test statistic  
 203  $U_{n,L}$  is also different. Therefore, we also adjust the rejection threshold from the quantile of a chi-  
 204 squared distribution, as in the original PCR test, to the quantile of the weighted sum of chi-squared  
 205 distributions.

#### 206 **3.2 Power Enhancement**

207 To effectively address covariate shift, incorporating density ratios as importance weights into the  
 208 PCR test is essential. However, when these ratios become large, they can increase the variance of the  
 209 statistics  $W_\ell$ . This elevated variance can diminish the test’s power. Therefore, developing methods to  
 210 reduce this variance is crucial for maintaining the power of the test.

---

**Algorithm 1** Covariate Shift Corrected PCR (csPCR) Test.

**Input:** Data  $D_{\mathcal{T}} = (\mathbf{y}, \mathbf{x}, \mathbf{Z}, \mathbf{V})$ , the density ratio  $e$ , the test statistics  $T$ , integers  $K, L \geq 1$ , and the significance level  $\alpha$ .

1: Take  $M = KL - 1$ .

2: **for** each data point  $j = 1$  to  $n$  **do**

3: Draw  $M$  i.i.d samples  $\tilde{X}_j^{(1)}, \dots, \tilde{X}_j^{(M)}$  from  $P_{\mathcal{T}}(X | \mathbf{Z})$ .

4: Use  $T$  to score the initial data point  $(X_j, Y_j, Z_j)$  and its  $M$  counterfeits  $(\tilde{X}_j^{(1:M)}, Y_j, Z_j)$

$$\begin{aligned} T_j &= T(X_j, Y_j, Z_j) \\ \tilde{T}_j^{(i)} &= T(\tilde{X}_j^{(i)}, Y_j, Z_j), \text{ for } i \in \{1, \dots, M\}. \end{aligned} \quad (4)$$

5: Let  $R_j$  denote the rank of  $T_j$  among  $\{T_j, \tilde{T}_j^{(1)}, \dots, \tilde{T}_j^{(M)}\}$ , with ties broken randomly.

6: Partition  $\{1, \dots, M + 1\} = S_1 \cup \dots \cup S_L$  with  $S_\ell := \{(\ell - 1)K + 1, \dots, \ell K\}$ . Assign label  $\ell_j \in \{1, 2, \dots, L\}$  to sample  $j$  if  $R_j \in S_{\ell_j}$ .

7: **end for**

8: Let  $w_j = e(X_j, Z_j, V_j)$  for each  $j \in \{1, 2, \dots, n\}$ .

9: **for** each label  $\ell \in \{1, 2, \dots, L\}$ : **do**

10: Let  $W_\ell$  be the sum of  $\ell$ -labeled importance weights:  $W_\ell = \sum_{j=1}^n w_j \cdot \mathbb{1}\{\ell_j = \ell\}$ .

11: Let  $D_\ell$  be the sum of  $\ell$ -labeled squared importance weights:  $D_\ell = \sum_{j=1}^n w_j^2 \cdot \mathbb{1}\{\ell_j = \ell\}$ .

12: **end for**

13: Let  $\hat{\Omega}_n = \frac{L}{n} \text{diag}(D_1, D_2, \dots, D_L) - \frac{1}{L} \cdot \mathbf{1}_{L \times L}$ .

14: Calculate the test statistic  $U_{n,L}$  as follows  $U_{n,L} = \frac{L}{n} \sum_{\ell=1}^L (W_\ell - \frac{n}{L})^2$ .

**Output:** Reject the null hypothesis if  $U_{n,L} \geq \theta_{\hat{\Omega}_n, \alpha}$ ; otherwise, accept the null hypothesis. Here,  $\theta_{\hat{\Omega}_n, \alpha}$  is the  $1 - \alpha$  quantile of the distribution  $\chi_{\hat{\Omega}_n}^2$ , where  $A \sim \chi_{\Omega}^2$  denotes that  $A = x^\top x$  for  $x \sim \mathcal{N}(0, \Omega)$ .

---

211 To this end, we introduce a control variate function  $a$ , allowing  $a(X, Z, V)$  to serve as a control  
 212 variate in reducing variance in  $W_\ell$  [14]. Specifically, for a chosen  $\gamma_\ell$ , we define

$$\tilde{W}_\ell = \sum_{j=1}^n w_j \cdot [\mathbb{1}\{\ell_j = \ell\} - \gamma_\ell a(X_j, Z_j, V_j)] + n\gamma_\ell \mathbb{E}_{\mathcal{T}} [a(X, Z, V)]. \quad (5)$$

213 We can then use  $\tilde{W}_\ell$  instead of  $W_\ell$  in our algorithm.

214 We note that for any arbitrary choice of the function  $a$  and the parameter  $\gamma_\ell$ , the expectation of  $\tilde{W}_\ell$   
 215 would be the same as that of  $W_\ell$ :

$$\begin{aligned} \mathbb{E} [\tilde{W}_\ell] &= \sum_{j=1}^n \mathbb{E} [w_j \mathbb{1}\{\ell_j = \ell\}] - \sum_{j=1}^n \gamma_\ell \mathbb{E} [w_j a(X_j, Z_j, V_j)] + n\gamma_\ell \mathbb{E}_{\mathcal{T}} [a(X, Z, V)] \\ &= \mathbb{E} [W_\ell] - n\gamma_\ell \left( \mathbb{E}_{\mathcal{S}} [e(X, Z, V) a(X, Z, V)] - \mathbb{E}_{\mathcal{T}} [a(X, Z, V)] \right) = \mathbb{E} [W_\ell]. \end{aligned} \quad (6)$$

216 Therefore, even if we make a sub-optimal choice of the function  $a$  and the parameter  $\gamma_\ell$  in practice,  
 217 the resulting test (under certain assumptions) will still remain asymptotically valid (see Section 3.3  
 218 for more details).

219 However, for effective variance reduction, it is preferable to have the control covariates  $a(X, Z, V)$   
 220 well-correlated with the outcome (See Section 4 for practical discussions on choices of the function  
 221  $a$ ). This is quite feasible, especially since the surrogate variable  $V$  is likely to be predictive of  $Y$ .

---

**Algorithm 2** Covariate Shift Corrected PCR Test with Power Enhancement.

---

**Input:** Data  $D_{\mathcal{T}} = (\mathbf{y}, \mathbf{x}, \mathbf{Z}, \mathbf{V})$ , the density ratio  $e$ , the test statistics  $T$ , the control variate function  $a$ , integers  $K, L \geq 1$ , and the significance level  $\alpha$ .

- 1: **for** each data point  $j = 1$  to  $n$  **do**
- 2:   Compute the labels  $\ell_j$  as in Algorithm 1.
- 3: **end for**
- 4: Let  $w_j = e(X_j, Z_j, V_j)$  for each  $j \in \{1, 2, \dots, n\}$ .
- 5: **for** each label  $\ell \in \{1, 2, \dots, L\}$ : **do**
- 6:   Compute  $\hat{\gamma}_\ell$ , the regression coefficient obtained by a weighted linear regression of the indicator function  $\{\mathbb{1}\{\ell_j = \ell\}\}_{j=1}^n$  on the control variate  $\{a(X_j, Z_j, V_j)\}_{j=1}^n$  with weights  $\{w_j\}_{j=1}^n$ .
- 7:   Compute the augmented version of  $W_\ell$  as

$$\widetilde{W}_\ell = \sum_{j=1}^n w_j \cdot [\mathbb{1}\{\ell_j = \ell\} - \hat{\gamma}_\ell a(X_j, Z_j, V_j)] + n\hat{\gamma}_\ell \mathbb{E}_{\mathcal{T}} [a(X, Z, V)].$$

- 8: **end for**
- 9: Let  $\mathbf{W} = \left( w_j \cdot [\mathbb{1}\{\ell_j = \ell\} - \hat{\gamma}_\ell a(X_j, Z_j, V_j)] + \hat{\gamma}_\ell \mathbb{E}_{\mathcal{T}} [a(X, Z, V)] \right)_{\ell, j}$  for  $1 \leq \ell \leq L, 1 \leq j \leq n$ .
- 10: Calculate the sample covariance matrix  $\widetilde{\Omega}_n = \frac{L}{n} (\mathbf{W} - \frac{1}{L} \cdot \mathbf{1}_{L \times n}) (\mathbf{W} - \frac{1}{L} \cdot \mathbf{1}_{L \times n})^\top$ .
- 11: Calculate the test statistic  $U_{n,L}$  as follows  $\widetilde{U}_{n,L} = \frac{L}{n} \sum_{\ell=1}^L \left( \widetilde{W}_\ell - \frac{n}{L} \right)^2$ .

**Output:** Reject the null hypothesis if  $\widetilde{U}_{n,L} \geq \theta_{\widetilde{\Omega}_n, \alpha}$ ; otherwise, accept the null hypothesis. Here,  $\theta_{\widetilde{\Omega}_n, \alpha}$  is the  $1 - \alpha$  quantile of the distribution  $\chi_{\widetilde{\Omega}_n}^2$ , where  $A \sim \chi_{\Omega}^2$  denotes that  $A = x^\top x$  for  $x \sim \mathcal{N}(0, \Omega)$ .

---

222 We would also like to discuss the choice of  $\gamma_\ell$ . According to the control covariate literature, with a  
 223 fixed function  $a$ , the optimal choice of  $\gamma_\ell$  that minimizes variance is given by:

$$\gamma_\ell = \frac{\text{Cov} \left[ w_j \mathbb{1}\{\ell_j = \ell\}, w_j a(X_j, Z_j, V_j) \right]}{\text{Var} \left[ w_j a(X_j, Z_j, V_j) \right]}. \quad (7)$$

224 This coefficient is also the same as that obtained from a linear regression [14]. Thus, when implement-  
 225 ing the algorithm, we take  $\gamma_\ell$  to be the regression coefficient obtained by running a weighted linear  
 226 regression of the indicator function  $\{\mathbb{1}\{\ell_j = \ell\}\}_{j=1}^n$  on the control variate  $\{a(X_j, Z_j, V_j)\}_{j=1}^n$  with  
 227 weights  $\{w_j\}_{j=1}^n$ .

228 We have outlined the new csPCR test, including this power enhancement step, in Algorithm 2.

### 229 3.3 Theoretical Properties

230 In this section, we establish that the proposed tests control the type-I error asymptotically. Further-  
 231 more, we show that the power enhancement step effectively reduces the variance of the statistics  $W_\ell$ ,  
 232 which can typically improve the power.

233 **Assumption 1** (Fourth moment). *The fourth moment of the density ratio  $e(X, Z, V)$  is finite:*  
 234  $\mathbb{E}_{\mathcal{S}} [e(X, Z, V)^4] < \infty$ . *Furthermore, the fourth moment of product of the density ratio and the*  
 235 *control variate function is also finite:*  $\mathbb{E}_{\mathcal{S}} [e(X, Z, V)^4 a(X, Z, V)^4] < \infty$ .

236 **Theorem 1** (Valid Tests). *Under Assumption 1, assume that the null hypothesis of  $X \perp\!\!\!\perp Y \mid Z$  holds*  
 237 *in the target population, then*

$$\lim_{n \rightarrow \infty} \mathbb{P} [\text{Algorithm 1 rejects}] = \alpha. \quad (8)$$

238 
$$\lim_{n \rightarrow \infty} \mathbb{P} [\text{Algorithm 2 rejects}] = \alpha. \quad (9)$$

239 **Theorem 2** (Variance Reduction). *Let  $W_l$  be the statistics computed in line 10 in Algorithm 1, and*  
 240  *$\widetilde{W}_l$  be the statistics computed in line 7 in Algorithm 2. Under Assumption 1,*

$$\limsup_{n \rightarrow \infty} \left( \text{Var} [\widetilde{W}_l] / \text{Var} [W_l] \right) \leq 1. \quad (10)$$

## 241 4 Numerical Simulation

242 In this section, we present simulation studies to assess the performance of our proposed csPCR  
 243 method and its power enhancement version denoted csPCR(pe), and compare them to a benchmark  
 244 method. The benchmark method adopted is an importance-resampling based method [17], denoted as  
 245 the IS method. For fair comparison, we used the same PCR statistic as our method for the testing  
 246 with IS. We use a significance level of  $\alpha = 0.05$ .

### 247 4.1 Simulation Setup

248 We consider a semi-supervised setting where we have a large volume of unlabeled data of  $(X_j, Z_j, V_j)$   
 249 from both the source and target populations. In addition, we have a small number of labeled data of  
 250  $(Y_j, X_j, Z_j, V_j)$  from the source population.

251 We separate confounding variables  $Z$  into two sets:  $Z = (Z_r, Z_{\text{null}})$ , where  $Z_r$  is the relevant set and  
 252  $Z_{\text{null}}$  is the null set. The relevant confounding variables  $Z_r$  are generated as i.i.d. multivariate normal,  
 253 with mean 0 for the source population and 1 for the target population to simulate the distributional  
 254 shift in  $Z$ , where  $Z_r \in \mathbb{R}^p$  and we set  $p = 5$ . Null confounding variables  $Z_{\text{null}}$  are generated  
 255 independently with no correlation to other variables, modeled as  $\mathcal{N}(0.1, I_q)$  with  $q = 50$  for sparse  
 256 high-dimensional settings in both populations.

257 The treatment variable  $X$  and the surrogate variable  $V$  are conditionally generated based on  $Z$ .  
 258 Specifically,  $X$  is modeled identically across both the source and target populations as  $\mathcal{N}(u^\top Z_r, 1)$ ,  
 259 where  $u$  is a predefined parameter vector that remains the same for both populations.

260 For  $V$ , it is modeled differently in the two populations, represented as  $\mathcal{N}(v_{\mathcal{S}/\mathcal{T}}^\top Z_r + (1 - \theta)a_{\mathcal{S}/\mathcal{T}}X +$   
 261  $\theta a_{\mathcal{S}/\mathcal{T}} \sin(X), 1)$ . Here,  $v_{\mathcal{S}}$  and  $v_{\mathcal{T}}$  are predefined parameter vectors for the source and target  
 262 populations, respectively. The parameter  $a$  varies between populations ( $a_{\mathcal{S}}$  for the source and  $a_{\mathcal{T}}$  for  
 263 the target), controlling the effect of  $X$  on  $V$ , modeling the indirect effect. The factor  $\theta$  modulates the  
 264 nonlinear component of this relationship.

265 The outcome variable  $Y$  is generated for both populations using the same conditional model over  
 266  $(X, Z, V)$ :

$$Y|(X, Z, V)_{\mathcal{S}/\mathcal{T}} \sim \mathcal{N}((v^\top Z_r)^2 + \beta V + \gamma X, 1),$$

267 where  $\beta$  and  $\gamma$  control the effects of  $V$  (indirect) and  $X$  (direct) on  $Y$ , respectively.

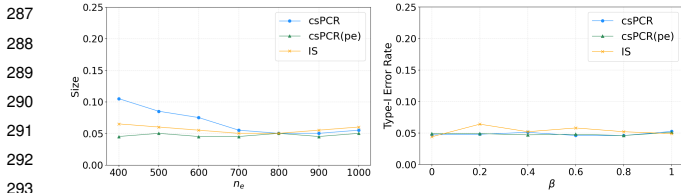
268 We generate 1000 unlabeled source and target samples to estimate the density ratio and generate 500  
 269 labeled source samples for testing. Moreover, in the simulation, we assume we have full knowledge  
 270 of the joint distribution of  $(X, Z)$  and estimate  $V|X, Z$  using an Elastic net regression model with  
 271 5-fold cross-validation [20]. For the test statistic  $T$  in the algorithm, we choose a simple function  
 272  $T(\tilde{X}, Z, V, Y) = Y \cdot \tilde{X}$ . For each parameter iteration, we conduct 1000 Monte Carlo simulations  
 273 to estimate the Type-I error and power. We estimate the covariance matrix of the sequence of  $W_i$ 's  
 274 using the Monte Carlo method and use the momentchi2 package [1] for calculating the  $p$ -value.  
 275 Additionally, we empirically choose the best hyperparameter  $L = 3$  for all our experiments through  
 276 additional experiments shown in Appendix B.2.

### 277 4.2 Simulation Results

278 In Figure 3, we choose  $a_{\mathcal{S}} = 1$  and  $a_{\mathcal{T}} = 0$  to compare the Type-I error control of our methods  
 279 with the benchmark. The left panel shows the Type-I error rate as the sample size of the data used  
 280 to estimate the density ratio,  $n_e$ , varies from small to large. There appears to be a slight Type-I  
 281 error inflation for all three methods when the sample size  $n_e$  is small, but the Type-I error quickly  
 282 converges to the ideal level of 0.05 as  $n_e$  grows larger. Moreover, our methods show more stable  
 283 Type-I error control than the benchmark method when the estimation sample size is low. The right



284 panel shows that when the density ratio is well approximated, all three methods attain good Type-I  
 285 error control regardless of the change in  $\beta$ , i.e., the strength of the indirect effect, but the csPCR and  
 286 csPCR(pe) methods have more stable control.



294 Figure 3: Comparison of Type-I error control across three methods.

295 effect size  $\beta$ . For example, when  $\beta = 1.4$ , the benchmark IS method has a power of 0.33, the csPCR  
 296 method has a power of 0.44, and the csPCR(pe) method can attain a power of 0.8.

298 When we fix the indirect effect  $\beta = 2$  and vary the direct effect of  $X$  ( $\gamma$ ), as shown in Figure 4b, our  
 299 methods still exceed the benchmark, and the power enhancement significantly improves the original  
 300 version of the test. For example, when  $\gamma = 1$ , the benchmark IS method has a power of 0.4, the  
 301 csPCR method has a power of 0.62, and the csPCR(pe) method can attain a power of 0.86.

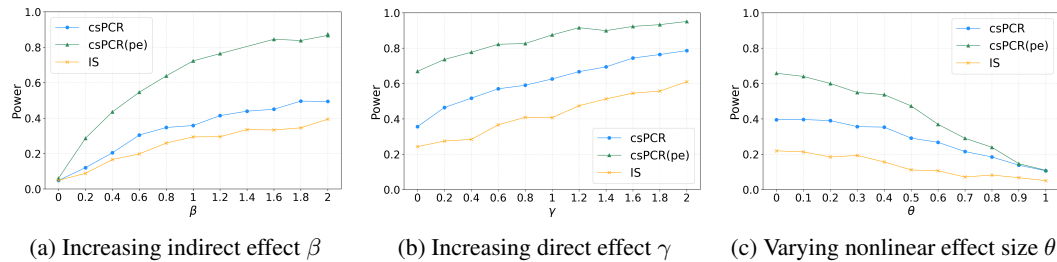


Figure 4: Comparison of statistical power of the three methods as the effect size varies: (a) indirect effect  $\beta$ , (b) direct effect  $\gamma$ , and (c) nonlinear effect size  $\theta$ .

302 We also test how adding a nonlinear component to the indirect effect affects the power when we  
 303 assume a linear model of  $V | Z, X$  in the estimation stage. This can be helpful in assessing the  
 304 performance of our methods under model misspecification. As Figure 4c indicates, as the nonlinear  
 305 effect increases, the power of all three methods decreases, though our methods still significantly  
 306 exceed the benchmark. Interestingly, we observe that as  $\theta \rightarrow 1$ , i.e., there is a full nonlinear  
 307 component without a linear component, the advantage of the power-enhanced version over the  
 308 original csPCR test disappears. This occurs because when the  $V | X, Z$  model is misspecified and  
 309 the density ratio estimation is inaccurate, the variance reduction in the control variates step reduces  
 310 variance in the “wrong” direction, and thus does not improve the power of the original method.

### 311 4.3 Effective Sample Size

312 We notice a series of work in measuring the effective sample size (ESS) of the density ratio reweighting  
 313 approaches [9]. Among them, one of the most common measure is  $n_{\text{eff}} = (\sum_{i=1}^n w_i)^2 / \sum_{i=1}^n w_i^2$ .  
 314 When the covariate shift between the source and target becomes stronger, the variance of the  
 315 importance weight  $w_i$  tends to be large and  $n_{\text{eff}}$  will become smaller, which could result in lower  
 316 power. We carry out simulation studies on the relationship between the power of csPCR and the ESS  
 317 determined by the degree of covariate shift as discussed in Appendix B.3.

## 318 5 Real-World Application

319 The COVID-19 pandemic has presented unprecedented challenges to global health systems, with  
 320 high variability in outcomes based on demographic and clinical characteristics. Early identification  
 321 of patients at high risk for severe outcomes, such as mortality within 90 days of hospital admission,  
 322 is crucial for timely and effective treatment interventions. This study leverages extensive hospital  
 323 data to develop models predicting 90-day mortality following hospital admission due to COVID-19.

324 For this study, we extract patient data spanning from January 2020 to December 2023 from Duke  
 325 University Health System (DUHS), focusing on individuals admitted with COVID-19. This period  
 326 encompasses multiple waves of the pandemic, influenced by various circulating variants.

327 Our dataset comprises patient records for a total of  $N = 3,057$  individuals admitted with COVID-19.  
 328 The outcome  $Y$  is defined as mortality within 90 days since hospital admission due to COVID-19.  
 329 The treatment variable  $X$  is defined as binary, where 1 indicates the administration of any COVID-19  
 330 specific medication (explained in Appendix C) and 0 otherwise. The covariates  $Z$  include comorbidity  
 331 indices (renal disease, diabetes without complication, diabetes with complication, local tumor, and  
 332 metastatic tumor), age, gender, and race, which are critical for adjusting the risk models due to  
 333 their known influence on COVID-19 outcomes. The length of hospitalization, denoted as  $V$ , is  
 334 standardized to follow a standard normal distribution (with a mean of zero and a standard deviation  
 335 of one), facilitating comparisons and integration into predictive models regardless of original scale or  
 336 distribution.

337 The dataset is segmented into two distinct groups based on the date of hospital admission to align  
 338 with pivotal changes in virus strain predominance and public health guidelines. The source data  
 339 comprises COVID-19 admissions prior to November 30, 2021, with a sample size of  $N_1 = 1,131$   
 340 patients. The target data includes admissions from November 30, 2021, through December 2023,  
 341 totaling  $N_2 = 792$  patients. This temporal division allows for the analysis of trends and outcomes  
 342 associated with the evolving pandemic landscape. Prevalence of the 90-day mortality outcome within  
 343 the source data is 14.3%, reflecting the impact of earlier virus strains and treatment protocols, while  
 344 in the target data, the prevalence is substantially lower at 3.7%, possibly indicating the effect of  
 345 improved treatments and vaccines, as well as the influence of different virus variants over time.

Table 1:  $p$ -values of different methods on COVID-19 dataset

Method	csPCR	csPCR(pe)	IS
$p$ -value	0.025	0.032	0.663

346 For the analysis, we divide 50% of the source data, comprising 565 individuals, alongside the entirety  
 347 of the target data, to estimate the density ratio. Density ratios of  $X, Z$  are estimated using probabilistic  
 348 classification method [12], while the density ratio of  $V|X, Z$  is determined through Elastic Net  
 349 regression. For all three methods, the test statistic  $T$  is chosen to be  $T(\tilde{X}, Z, V, Y) = Y \cdot \tilde{X}$ . As  
 350 indicated in Table 1, both csPCR and csPCR(pe) give statistically significant results, whereas the  
 351 IS method does not. The statistically significant results are consistent with biomedical literature.  
 352 For example, through systematic review and meta-analysis, [21] reported that Bamlanivimab is  
 353 effective in reducing the mortality rates of COVID patients. In a cohort study, [16] also found similar  
 354 effectiveness for Nirmatrelvir–ritonavir.

355 These results align with our findings from the simulation study and demonstrate that our method has  
 356 increased power compared with the benchmark IS method.

357 **References**

- 358 [1] Dean A Bodenham and Niall M Adams. A comparison of efficient approximations for a  
359 weighted sum of chi-squared random variables. *Statistics and Computing*, 26(4):917–928, 2016.
- 360 [2] Emmanuel Candes, Yingying Fan, Lucas Janson, and Jinchi Lv. Panning for gold: ‘model-  
361 x’ knockoffs for high dimensional controlled variable selection. *Journal of the Royal Statistical  
362 Society Series B: Statistical Methodology*, 80(3):551–577, 2018.
- 363 [3] Emerson Y Chen, Vikram Raghunathan, and Vinay Prasad. An overview of cancer drugs  
364 approved by the us food and drug administration based on the surrogate end point of response  
365 rate. *JAMA Internal Medicine*, 179(7):915–921, 2019.
- 366 [4] Jiayuan Huang, Arthur Gretton, Karsten Borgwardt, Bernhard Schölkopf, and Alex Smola.  
367 Correcting sample selection bias by unlabeled data. *Advances in neural information processing  
368 systems*, 19, 2006.
- 369 [5] Adel Javanmard and Mohammad Mehrabi. Pearson chi-squared conditional randomization test.  
370 *arXiv preprint arXiv:2111.00027*, 2021.
- 371 [6] Shuangning Li and Molei Liu. Maxway crt: improving the robustness of the model-x inference.  
372 *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 85(5):1441–1470,  
373 2023.
- 374 [7] Molei Liu, Eugene Katsevich, Lucas Janson, and Aaditya Ramdas. Fast and powerful conditional  
375 randomization testing via distillation. *Biometrika*, 109(2):277–293, 2022.
- 376 [8] Molei Liu, Yi Zhang, Katherine P Liao, and Tianxi Cai. Augmented transfer regression learning  
377 with semi-non-parametric nuisance models. *Journal of Machine Learning Research*, 24(293):  
378 1–50, 2023.
- 379 [9] Luca Martino, Víctor Elvira, and Francisco Louzada. Effective sample size for importance  
380 sampling based on discrepancy measures. *Signal Processing*, 131:386–401, 2017.
- 381 [10] Binh T Nguyen, Bertrand Thirion, and Sylvain Arlot. A conditional randomization test for  
382 sparse logistic regression in high-dimension. *Advances in Neural Information Processing  
383 Systems*, 35:13691–13703, 2022.
- 384 [11] Ziang Niu, Abhinav Chakraborty, Oliver Dukes, and Eugene Katsevich. Reconciling  
385 model-x and doubly robust approaches to conditional independence testing. *arXiv preprint  
386 arXiv:2211.14698*, 2022.
- 387 [12] Jing Qin. Inferences for case-control and semiparametric two-sample density ratio models.  
388 *Biometrika*, 85(3):619–630, 1998.
- 389 [13] James M Robins, Andrea Rotnitzky, and Lue Ping Zhao. Estimation of regression coefficients  
390 when some regressors are not always observed. *Journal of the American statistical Association*,  
391 89(427):846–866, 1994.
- 392 [14] Sheldon M Ross. *Simulation*. academic press, 2022.
- 393 [15] Donald B Rubin. The calculation of posterior distributions by data augmentation: Comment: A  
394 noniterative sampling/importance resampling alternative to the data augmentation algorithm for  
395 creating a few imputations when fractions of missing information are modest: The sir algorithm.  
396 *Journal of the American Statistical Association*, 82(398):543–546, 1987.
- 397 [16] Kevin L Schwartz, Jun Wang, Mina Tadrous, Bradley J Langford, Nick Daneman, Valerie  
398 Leung, Tara Gomes, Lindsay Friedman, Peter Daley, and Kevin A Brown. Population-based  
399 evaluation of the effectiveness of nirmatrelvir–ritonavir for reducing hospital admissions and  
400 mortality from covid-19. *Cmaj*, 195(6):E220–E226, 2023.
- 401 [17] Nikolaj Thams, Sorawit Saengkyongam, Niklas Pfister, and Jonas Peters. Statistical testing under  
402 distributional shifts. *Journal of the Royal Statistical Society Series B: Statistical Methodology*,  
403 85(3):597–663, 2023.

- 404 [18] Kaizheng Wang. Pseudo-labeling for kernel ridge regression under covariate shift. *arXiv*  
405 *preprint arXiv:2302.10160*, 2023.
- 406 [19] Wenshuo Wang and Lucas Janson. A high-dimensional power analysis of the conditional  
407 randomization test and knockoffs. *Biometrika*, 109(3):631–645, 2022.
- 408 [20] Hui Zou and Trevor Hastie. Regularization and variable selection via the elastic net. *Journal of*  
409 *the Royal Statistical Society Series B: Statistical Methodology*, 67(2):301–320, 2005.
- 410 [21] Ling Zuo, Guangyu Ao, Yushu Wang, Ming Gao, and Xin Qi. Bamlanivimab improves  
411 hospitalization and mortality rates in patients with covid-19: a systematic review and meta-  
412 analysis. *The Journal of infection*, 84(2):248, 2022.

413 **A Proofs**

414 **A.1 Preliminaries**

415 Throughout this section, we write  $S(x_j, z_j, v_j) = s_j$  as the label assigned to sample  $j$  in Algorithms  
416 1 and 2, instead of using  $\ell_j$ . This notation helps avoid confusion between different label choices.

417 **Proposition 1.** Assume that the conditional independence  $X \perp\!\!\!\perp Y \mid Z$  holds on the target population  
418  $\mathcal{T}$ . Let  $e(x_j, z_j, v_j)$  denote the density ratio. For any integer  $\ell \in [1, L]$ , the following holds:

$$\mathbb{E}_{\mathcal{S}}[e(x_j, z_j, v_j) \cdot \mathbb{1}\{S_{\mathcal{T}}(x_j, z_j, v_j) = \ell\}] = \frac{1}{L}.$$

419 *Proof of Proposition 1.* For simplicity, denote  $w_j = e(X_j, Z_j, V_j)$  and  $s_j = S_{\mathcal{T}}(X_j, Z_j, V_j)$ .

$$\begin{aligned} & \mathbb{E}_{\mathcal{S}}[e(X_j, Z_j, V_j) \cdot \mathbb{1}\{S_{\mathcal{T}}(X_j, Z_j, V_j) = \ell\}] \\ &= \mathbb{E}_{\mathcal{S}} \left[ \mathbb{E} [w_j \cdot \mathbb{P}(s_j = \ell \mid Y_j, Z_j, X_j, V_j)] \mid Z_j, X_j, V_j \right] \\ &= \mathbb{E}_{\mathcal{S}} \left[ w_j \cdot \mathbb{E} [\mathbb{P}(s_j = \ell \mid Y_j, Z_j, X_j, V_j)] \mid Z_j, X_j, V_j \right] \\ &= \int_{(Z_j, X_j, V_j)} w_j \cdot p_{\mathcal{S}}(Z_j, X_j, V_j) \cdot \mathbb{E} [\mathbb{P}(s_j = \ell \mid Y_j, Z_j, X_j, V_j)] dZ_j dX_j dV_j \\ &= \int_{(Z_j, X_j, V_j)} p_{\mathcal{T}}(Z_j, X_j, V_j) \cdot \mathbb{E} [\mathbb{P}(s_j = \ell \mid Y_j, Z_j, X_j, V_j)] dZ_j dX_j dV_j \\ &= \mathbb{E}_{\mathcal{T}} \left[ \mathbb{E} [\mathbb{P}(s_j = \ell \mid Y_j, Z_j, X_j, V_j)] \mid Z_j, X_j, V_j \right] \\ &= \mathbb{E}_{\mathcal{T}} [\mathbb{P}(s_j = \ell \mid Y_j, Z_j, X_j, V_j)] \\ &= \frac{1}{L}. \end{aligned} \tag{11}$$

420 The last equation follows from results in the non-covariate-shift scenario, e.g., from [5].  $\square$

421 **A.2 Proof of Theorem 1**

422 **A.2.1 Results for Algorithm 1**

423 Let  $(W_\ell)_{\ell=1, \dots, L}$  be the sum of weights and  $\hat{\Omega}_n$  be the sample covariance matrix in Algorithm 1. By  
424 Proposition 1, we have that

$$\mathbb{E}(W_\ell) = n \cdot \mathbb{E}[w_j \cdot \mathbb{1}\{\ell_j = \ell\}] = \frac{n}{L}.$$

425 Note that the  $W_\ell$ 's are sums of i.i.d. random variables, and thus by the Central Limit Theorem, as  
426  $n \rightarrow \infty$ ,

$$\mathbf{A}_n = \sqrt{\frac{L}{n}} \left( W_1 - \frac{n}{L}, W_2 - \frac{n}{L}, \dots, W_L - \frac{n}{L} \right) \xrightarrow{d} \mathcal{N}_L(0, \Omega),$$

427 where for any  $\ell, \ell^* \in \{1, \dots, L\}$

$$\begin{aligned} \Omega_{\ell, \ell^*} &= LCov(w_1 \mathbb{1}\{s_1 = \ell\}, w_1 \mathbb{1}\{s_1 = \ell^*\}) = L\mathbb{E}_{\mathcal{S}} \left[ w_1^2 \cdot \mathbb{1}\{s_1 = \ell\} \mathbb{1}\{s_1 = \ell^*\} \right] - \frac{1}{L} \\ &= L\mathbb{E}_{\mathcal{S}} \left[ w_1^2 \cdot \mathbb{1}\{s_1 = \ell\} \right] \mathbb{1}\{\ell = \ell^*\} - \frac{1}{L}. \end{aligned}$$

428 Therefore,

$$U_{n,L} = \mathbf{A}_n^\top \mathbf{A}_n \xrightarrow{d} \chi_{\Omega}^2.$$

429 Next, we will focus on the variance estimation part. We will show that  $\hat{\Omega}_n \xrightarrow{p} \Omega$  as  $n \rightarrow \infty$ . For any  
 430  $\ell, \ell^* \in \{1, \dots, L\}$ ,

$$\begin{aligned} \hat{\Omega}_{n,\ell,\ell^*} &= \mathbb{1}\{\ell = \ell^*\} \frac{L}{n} D_l - \frac{1}{L} = \mathbb{1}\{\ell = \ell^*\} \frac{L}{n} \sum_{j=1}^n w_j^2 \cdot \mathbb{1}\{\ell_j = \ell\} - \frac{1}{L} \\ &\xrightarrow{p} \mathbb{1}\{\ell = \ell^*\} L \mathbb{E}_{\mathcal{S}} \left[ w_1^2 \cdot \mathbb{1}\{s_1 = \ell\} \right] - \frac{1}{L} = \Omega_{\ell,\ell^*}. \end{aligned}$$

431 Up til now, we have that

$$U_{n,L} = \mathbf{A}_n^\top \mathbf{A}_n \xrightarrow{d} \chi_\Omega^2, \quad \text{and} \quad \theta_{\hat{\Omega}_n, \alpha} \xrightarrow{p} \theta_{\Omega, \alpha}.$$

432 Therefore,

$$\mathbb{P}(\text{Algorithm 1 rejects}) = \mathbb{P}(U_{n,L} \geq \theta_{\hat{\Omega}_n, \alpha}) \rightarrow \mathbb{P}(\chi_\Omega^2 \geq \theta_{\Omega, \alpha}) = \alpha.$$

### 433 A.2.2 Results for Algorithm 2

434 Let  $(\tilde{W}_\ell)_{\ell=1, \dots, L}$  and  $\tilde{\Omega}_n$  be the sum of weights and the sample covariance matrix in Algorithm 2.  
 435 Let  $\hat{\gamma}_\ell$  be the estimated coefficient in Algorithm 2.

436 Recall that in (7), we have identified the optimal choice of  $\gamma_\ell$ . We will start by working with this  
 437 optimal choice and show that the  $\hat{\gamma}_\ell$  is close to it. Define

$$\begin{aligned} \tilde{W}_\ell &= \sum_{j=1}^n w_j (\mathbb{1}\{\ell_j = \ell\} - \gamma_\ell a(X_j, Z_j, V_j)) + n \gamma_\ell \mathbb{E}_{\mathcal{T}}[a(X, Z, V)], \\ K_{\ell,j}(\gamma) &= w_j (\mathbb{1}\{\ell_j = \ell\} - \gamma a(X_j, Z_j, V_j)) + \gamma \mathbb{E}_{\mathcal{T}}[a(X, Z, V)], \text{ and} \\ H_j &= w_j a(X_j, Z_j, V_j) - \mathbb{E}_{\mathcal{T}}[a(X, Z, V)]. \end{aligned}$$

438 Therefore, we have  $\tilde{W}_\ell = \sum_j K_{\ell,j}(\gamma_\ell)$  and  $\tilde{W}_\ell = \sum_j K_{\ell,j}(\hat{\gamma}_\ell) = \tilde{W}_\ell - (\hat{\gamma}_\ell - \gamma_\ell) \sum_j H_j$ .

439 Note that by (6),  $\mathbb{E}(H_j) = 0$ . By Proposition 1 and (6), we have  $\mathbb{E}(\tilde{W}_\ell) = \frac{n}{L}$ . Furthermore, because  
 440  $\tilde{W}_\ell$  is a sum of i.i.d. random variables, we have that as  $n \rightarrow \infty$ ,

$$\check{\mathbf{A}}_n = \sqrt{\frac{L}{n}} \left( \tilde{W}_1 - \frac{n}{L}, \tilde{W}_2 - \frac{n}{L}, \dots, \tilde{W}_L - \frac{n}{L} \right) \xrightarrow{d} \mathcal{N}_L(0, \check{\Omega}), \quad (12)$$

441 where for each  $\ell, \ell^* \in \{1, \dots, L\}$ ,

$$\check{\Omega}_{\ell,\ell^*} = LCov \{K_{\ell,j}(\gamma_\ell), K_{\ell^*,j}(\gamma_{\ell^*})\}.$$

442 Therefore,

$$\check{U}_{n,L} = \check{\mathbf{A}}_n^\top \check{\mathbf{A}}_n \xrightarrow{d} \chi_{\check{\Omega}}^2.$$

443 Next, we will show that the actual statistic  $\tilde{U}_{n,L}$  is close to  $\check{U}_{n,L}$ , and that the estimated variance  
 444 matrix is also close to  $\check{\Omega}$ . We start with noting that the estimator  $\hat{\gamma}_\ell$  from linear regression is close to  
 445 the optimal choice  $\gamma_\ell$  defined in (7): by the Central Limit Theorem,  $\hat{\gamma}_\ell = \gamma_\ell + \mathcal{O}_p(1/\sqrt{n})$ . And thus

$$\begin{aligned} \tilde{W}_\ell &= \sum_{j=1}^n w_j (\mathbb{1}\{\ell_j = \ell\} - \hat{\gamma}_\ell a(X_j, Z_j, V_j)) + n \hat{\gamma}_\ell \mathbb{E}_{\mathcal{T}}[a(X, Z, V)] \\ &= \sum_{j=1}^n w_j \mathbb{1}\{\ell_j = \ell\} - \hat{\gamma}_\ell \left( \sum_{j=1}^n w_j a(X_j, Z_j, V_j) - n \mathbb{E}_{\mathcal{T}}[a(X, Z, V)] \right) \\ &= \sum_{j=1}^n w_j \mathbb{1}\{\ell_j = \ell\} - \gamma_\ell \left( \sum_{j=1}^n w_j a(X_j, Z_j, V_j) - n \mathbb{E}_{\mathcal{T}}[a(X, Z, V)] \right) + \mathcal{O}_p(1) \\ &= \tilde{W}_\ell + \mathcal{O}_p(1). \end{aligned}$$

446 The second-to-last line is because  $\hat{\gamma}_\ell = \gamma_\ell + \mathcal{O}_p(1/\sqrt{n})$  and the terms inside the parenthesis,  $\sum_j H_j$ ,  
 447 is a sum of  $n$  independent mean-zero random variables.

448 Therefore, together with (12), by Slutsky's Theorem, we have that

$$\tilde{\mathbf{A}}_n = \sqrt{\frac{L}{n}} \left( \tilde{W}_1 - \frac{n}{L}, \tilde{W}_2 - \frac{n}{L}, \dots, \tilde{W}_L - \frac{n}{L} \right) \xrightarrow{d} \mathcal{N}_L(0, \tilde{\Omega}),$$

449 and thus,

$$\tilde{U}_{n,L} = \tilde{\mathbf{A}}_n^\top \tilde{\mathbf{A}}_n \xrightarrow{d} \chi_{\tilde{\Omega}}^2.$$

450 We will work on sample covariance matrix now. Recall that the sample covariance matrix  $\tilde{\Omega}_n =$   
 451  $\frac{L}{n} (\mathbf{W} - \frac{1}{L} \cdot \mathbf{1}_{L \times n}) (\mathbf{W} - \frac{1}{L} \cdot \mathbf{1}_{L \times n})^\top$ , where  $\mathbf{W}_{\ell,j} = w_j \cdot [\mathbb{1}\{\ell_j = \ell\} - \hat{\gamma}_\ell a(X_j, Z_j, V_j)] +$   
 452  $\hat{\gamma}_\ell \mathbb{E}_{\mathcal{T}} [a(X, Z, V)] = K_{\ell,j}(\hat{\gamma}_\ell)$ . Let's start with  $\mathbf{W}\mathbf{W}^\top$ . For any  $\ell, \ell^* \in \{1, \dots, L\}$ ,

$$\begin{aligned} (\mathbf{W}\mathbf{W}^\top)_{\ell,\ell^*} &= \sum_j K_{\ell,j}(\hat{\gamma}_\ell) K_{\ell^*,j}(\hat{\gamma}_{\ell^*}) \\ &= \sum_j (K_{\ell,j}(\gamma_\ell) - (\hat{\gamma}_\ell - \gamma_\ell) H_j) (K_{\ell^*,j}(\gamma_{\ell^*}) - (\hat{\gamma}_{\ell^*} - \gamma_{\ell^*}) H_j) \\ &= \sum_j K_{\ell,j}(\gamma_\ell) K_{\ell^*,j}(\gamma_{\ell^*}) - (\hat{\gamma}_\ell - \gamma_\ell) \sum_j H_j K_{\ell^*,j}(\gamma_{\ell^*}) - (\hat{\gamma}_{\ell^*} - \gamma_{\ell^*}) \sum_j H_j K_{\ell,j}(\gamma_\ell) \\ &\quad + (\hat{\gamma}_\ell - \gamma_\ell)(\hat{\gamma}_{\ell^*} - \gamma_{\ell^*}) \sum_j H_j^2 \\ &= \sum_j K_{\ell,j}(\gamma_\ell) K_{\ell^*,j}(\gamma_{\ell^*}) + \mathcal{O}_p(\sqrt{n}) \end{aligned}$$

453 Therefore, by the law of large numbers,

$$\frac{L}{n} (\mathbf{W}\mathbf{W}^\top)_{\ell,\ell^*} = \frac{L}{n} \sum_j K_{\ell,j}(\gamma_\ell) K_{\ell^*,j}(\gamma_{\ell^*}) + \mathcal{O}_p(1/\sqrt{n}) = L \mathbb{E} [K_{\ell,1}(\gamma_\ell) K_{\ell^*,1}(\gamma_{\ell^*})] + \mathcal{O}_p(1/\sqrt{n}).$$

454 Similarly, for  $\mathbf{W}\mathbf{1}^\top$ , we have that for any  $\ell, \ell^* \in \{1, \dots, L\}$ ,

$$(\mathbf{W}\mathbf{1}^\top)_{\ell,\ell^*} = \sum_j K_{\ell,j}(\hat{\gamma}_\ell) = \sum_j K_{\ell,j}(\gamma_\ell) - (\hat{\gamma}_\ell - \gamma_\ell) H_j = \sum_j K_{\ell,j}(\gamma_\ell) + \mathcal{O}_p(\sqrt{n}).$$

455 Therefore, again by the law of large numbers,

$$\frac{L}{N} (\mathbf{W}\mathbf{1}^\top)_{\ell,\ell^*} = \frac{L}{N} \sum_j K_{\ell,j}(\gamma_\ell) + \mathcal{O}_p(1/\sqrt{n}) = L \mathbb{E} [K_{\ell,j}(\gamma_\ell)] + \mathcal{O}_p(1/\sqrt{n}) = 1 + \mathcal{O}_p(1/\sqrt{n}).$$

456 Combining the above results gives,

$$\begin{aligned} \tilde{\Omega}_{n,\ell,\ell^*} &= \frac{L}{n} \left[ (\mathbf{W} - \frac{1}{L} \cdot \mathbf{1}_{L \times n}) (\mathbf{W} - \frac{1}{L} \cdot \mathbf{1}_{L \times n})^\top \right]_{\ell,\ell^*} \\ &= L \mathbb{E} [K_{\ell,1}(\gamma_\ell) K_{\ell^*,1}(\gamma_{\ell^*})] - L \mathbb{E} [K_{\ell,j}(\gamma_\ell)] \mathbb{E} [K_{\ell^*,j}(\gamma_{\ell^*})] + \mathcal{O}_p(1/\sqrt{n}) \\ &= L \text{Cov} [K_{\ell,1}(\gamma_\ell), K_{\ell^*,1}(\gamma_{\ell^*})] + \mathcal{O}_p(1/\sqrt{n}) \\ &= \tilde{\Omega}_{\ell,\ell^*} + \mathcal{O}_p(1/\sqrt{n}). \end{aligned}$$

457 Therefore,  $\tilde{\Omega}_n \xrightarrow{p} \tilde{\Omega}$ .

458 To summarize, we have that

$$\tilde{U}_{n,L} = \tilde{\mathbf{A}}_n^\top \tilde{\mathbf{A}}_n \xrightarrow{d} \chi_{\tilde{\Omega}}^2, \quad \text{and} \quad \theta_{\tilde{\Omega}_n, \alpha} \xrightarrow{p} \theta_{\tilde{\Omega}, \alpha}.$$

459 Therefore,

$$\mathbb{P}(\text{Algorithm 2 rejects}) = \mathbb{P}(\tilde{U}_{n,L} \geq \theta_{\tilde{\Omega}_n, \alpha}) \rightarrow \mathbb{P}(\chi_{\tilde{\Omega}}^2 \geq \theta_{\tilde{\Omega}, \alpha}) = \alpha.$$

460 **A.3 Proof of Theorem 2**

461 Similar to the proof of Theorem 1, we define

$$\begin{aligned}\widetilde{W}_\ell &= \sum_{j=1}^n w_j (\mathbb{1}\{\ell_j = \ell\} - \gamma_\ell a(X_j, Z_j, V_j)) + n\gamma_\ell \mathbb{E}_{\mathcal{T}}[a(X, Z, V)], \\ K_{\ell,j}(\gamma) &= w_j (\mathbb{1}\{\ell_j = \ell\} - \gamma a(X_j, Z_j, V_j)) + \gamma \mathbb{E}_{\mathcal{T}}[a(X, Z, V)], \text{ and} \\ H_j &= w_j a(X_j, Z_j, V_j) - \mathbb{E}_{\mathcal{T}}[a(X, Z, V)].\end{aligned}$$

462 Therefore, we have  $\widetilde{W}_\ell = \sum_j K_{\ell,j}(\gamma_\ell)$  and  $\widetilde{W}_\ell = \sum_j K_{\ell,j}(\hat{\gamma}_\ell) = \widetilde{W}_\ell - (\hat{\gamma}_\ell - \gamma_\ell) \sum_j H_j$ .

463 We know from the literature that  $\gamma_\ell$  is the optimal choice of  $\gamma$  and thus  $\text{Var}[\widetilde{W}_\ell] \leq \text{Var}[W_\ell]$ . We  
464 will then move on to show that  $\text{Var}[\widetilde{W}_\ell]$  is close to  $\text{Var}[W_\ell]$  and thus asymptotically no greater  
465 than  $\text{Var}[W_\ell]$ .

466 To this end, note that

$$\begin{aligned}\text{Var}[\widetilde{W}_\ell] &= \text{Var}\left[\widetilde{W}_\ell - (\hat{\gamma}_\ell - \gamma_\ell) \sum_j H_j\right] \\ &= \text{Var}[\widetilde{W}_\ell] + 2 \text{Cov}\left[\widetilde{W}_\ell, (\hat{\gamma}_\ell - \gamma_\ell) \sum_j H_j\right] + \text{Var}\left[(\hat{\gamma}_\ell - \gamma_\ell) \sum_j H_j\right] \\ &\leq \text{Var}[\widetilde{W}_\ell] + 2\sqrt{\text{Var}[\widetilde{W}_\ell] \mathbb{E}\left[\left((\hat{\gamma}_\ell - \gamma_\ell) \sum_j H_j\right)^2\right]} + \mathbb{E}\left[\left((\hat{\gamma}_\ell - \gamma_\ell) \sum_j H_j\right)^2\right].\end{aligned}$$

467 But we also know from the proof of Theorem 1 that  $\hat{\gamma}_\ell - \gamma_\ell \xrightarrow{P} 0$ . Then, because of the bounded  
468 fourth moment assumption, by the Dominated Convergence Theorem, we have that

$$\frac{1}{n} \mathbb{E}\left[\left((\hat{\gamma}_\ell - \gamma_\ell) \sum_j H_j\right)^2\right] \rightarrow 0.$$

469 Therefore,

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \left( \text{Var}[\widetilde{W}_\ell] - \text{Var}[W_\ell] \right) \leq 0.$$

470 Finally, we note that  $\text{Var}[W_\ell] = \Omega(n)$ , and hence

$$\limsup_{n \rightarrow \infty} (\text{Var}[\widetilde{W}_\ell] / \text{Var}[W_\ell]) \leq 1.$$

471 **B Additional Simulation Results**

472 **B.1 Running time**

473 All experiments run on a Macbook Pro 2022 M2.

474 **Artificial dataset:** Regarding running time for one iteration including density ratio estimation and  
475  $X|Z$  model fitting (on average), csPCR took 5.12s, csPCR(pe) took 14.95s, IS method took 1.5s,  
476 PCR took 1.25s.

477 **Real-world application:** Regarding running time for one test procedure, csPCR took 3.41s,  
478 csPCR(pe) took 11.32s, IS method took 0.81s.

479 **B.2 Finding optimal hyperparameter  $L$**

480 We find the optimal  $L$  value for the testing algorithm by performing numerical simulations, evaluating  
481 its Type-I error control and power. We adopt the same numerical simulation setup as in the main text  
482 Section 4. We first choose  $a_S = 1$  and  $a_{\mathcal{T}} = 0$  and also fix  $\beta = 1$  to compare the Type-I error rate  
483 for different choice of  $L$  of the csPCR method. We perform experiments with both true density ratio  
484 and estimated density ratio. The results are shown in Table 2.



Table 2: Type-I Error Rates at Different Levels of L of csPCR Method

L	2	3	5	10	15	20
True Density Ratio	0.05125	0.05000	0.04575	0.03675	0.02825	0.02425
Estimated Density Ratio	0.04620	0.05025	0.04425	0.03905	0.02725	0.02175

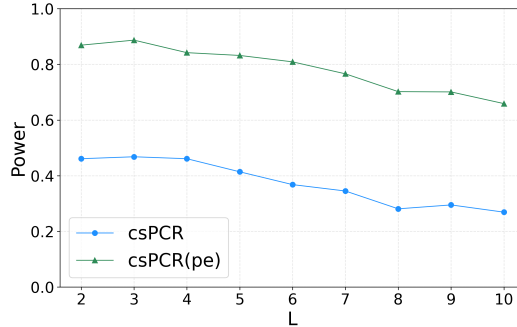


Figure 5: Comparison of statistical power of the three methods as the parameter  $L$  varies.

485 We also test the power of the csPCR and csPCR(pe) method with different choices of  $L$  value. We  
 486 choose  $a_S = 0$  and  $a_T = 2$  and fix  $\beta = 2$ .

487 As Table 2 and Figure 5 shows, as  $L$  value increases, the csPCR method become more conservative  
 488 with more tight Type-I error control and lower power. We can observe that when we set  $L = 3$ , the  
 489 csPCR method can achieve most stable Type-I error rate control and also highest power empirically.  
 490 Therefore, in our simulation experiments and real world data experiments, we fix  $L = 3$ .

### 491 B.3 Role of effective sample size

492 We notice a series of work in measuring the effective sample size (ESS) of importance weight or  
 493 sampling in the statistical computation literature, e.g., [Martino, et al, 2017] and others. Among them,  
 494 one of the most common ways is to use the ratio  $n_{eff} = \frac{(\sum_{i=1} w_i)^2}{\sum_{i=1} w_i^2}$  to approximate the ESS. When  
 495 the covariate shift between the source and target becomes stronger, the variance of the importance  
 496 weight  $w_i$  tends to be large and  $n_{eff}$  will become smaller, which can result in lower power. Our  
 497 power enhancement method based on control variate could potentially alleviate this issue with  
 498 properly specified control functions.  
 499

500 In the simulation study, we varied only  $\mu_z$ , the mean of the confounding variables  $Z_T$ . A higher  $\mu_z$   
 501 signifies a stronger covariate shift between the source and target populations. From Figure 6, it is  
 502 evident that as  $\mu_z$  increases, the Effective Sample Size (ESS) required significantly decreases, while  
 503 the power of the csPCR method concurrently declines. These results suggest that increasing covariate  
 504 shift leads to a reduction in ESS and a corresponding decrease in statistical power.

### 505 B.4 Instability of the Importance Resampling (IS) method

506 In this section, we will use numerical simulations to illustrate that the performance if the IS method  
 507 is subject to the resample size heavily. IS method performs resampling without replacement and  
 508 typically has to sample a much smaller subset (theoretically, in the order of  $o(\sqrt{n})$ ) of the source data  
 509 to approximate the target. Consequently, the power of IS is substantially lower than our approach. If  
 510 the resample size of IS is overly increased, it may fail to control the Type-I error due to excessive  
 511 similarity between the resampled data and the original source data.  
 512

513 To further illustrate, we conducted additional experiments with varied resample sizes in IS to assess  
 514 its effect on Type-I error control and power. From Figure 7, one can observe that IS starts to show  
 515 high Type-I error inflation when its resample size increases to 400 but still shows much lower power

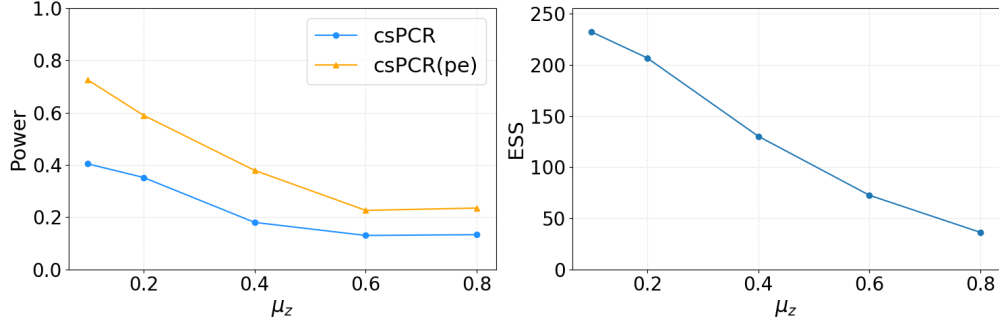


Figure 6: The left panel shows the comparison of statistical power of csPCR and csPCR(pe) method as the covariate shift gets stronger. The right panel illustrates how the Effective Sample Size(ESS) changes as covariate shift scale becomes larger.

516 (by around 0.4) than our method with this resample size (or even larger ones). This indicates that our  
 517 method achieves better statistical efficiency than IS (DRPL).

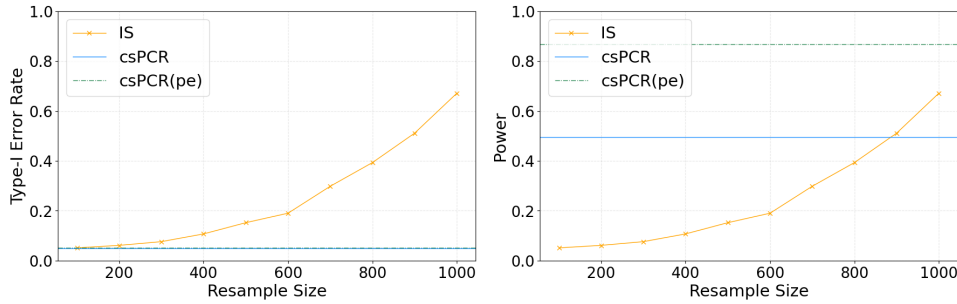


Figure 7: Detailed comparison of Type-I error rate and power of csPCR and the IS method. With the source sample size  $n_s = 1000$ , we gradually increase the resample size for the IS method from 100 to 1000. The two horizontal lines represent the Type-I error rate and power, respectively, of the csPCR and csPCR(pe) methods (they do not change with the tuning of IS).

517

## 518 B.5 Choice of test statistic

519 In this section we explore the effect of test statistics on the algorithm performance. The main  
 520 principle of choosing the test statistic is to characterize the conditional dependency between  $X$  and  
 521  $Y$  under the alternative hypothesis. The test statistic  $YX$  may not be the optimal choice and that  
 522 using  $(Y - \hat{E}[Y | Z])(X - E[X | Z])$  could remove the confounding effect of  $Z$ .

523

524 Inspired by this, we used  $Y(X - E[X|Z])$  as the test statistic to conduct additional simulations. As  
 525 illustrated in Figure8, we find that  $Y(X - E[X|Z])$  and  $YX$  produce nearly the same power for  
 526 both csPCR and csPCR(pe) with the change of effect size.

## 527 C Real-World Application

528 The specific medication indicated by the treatment variable  $X$  includes Ritonavir, Bamlanivimab,  
 529 Casirivimab-Imdevimab, Remdesivir, Ritonavir Nirmatrelvir, Sotrovimab, Bamlanivimab Etesevimab.  
 530 For simplicity,  $X = 1$  indicates any of these specific medication and  $X = 0$  otherwise.

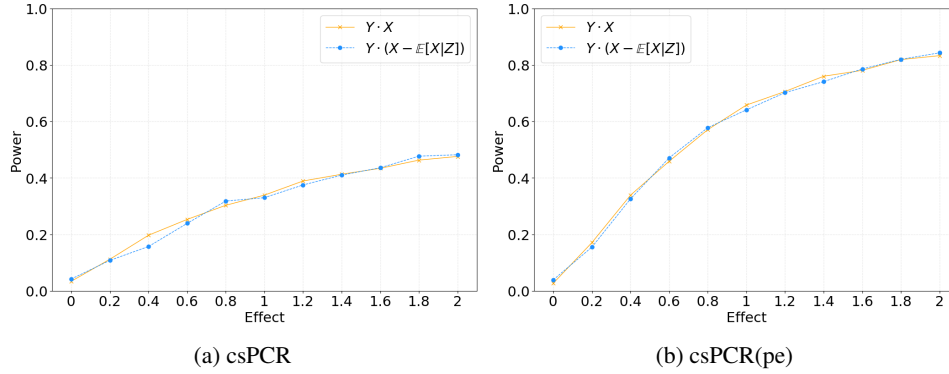


Figure 8: Power against effect size for csPCR and csPCR(pe) with two different test statistics  $XY$  and  $(X - E[X | Z])Y$ . We observe that the power is very similar with the two different test statistics.

531 **C.1 Different outcome**

532 In our real data experiment part, the outcome variable  $Y$  is defined as mortality within 90 days since  
 533 hospital admission due to COVID-19. In addition, we also analyzed mortality within 30 days since  
 534 hospital admission. As shown in Table 3, both csPCR and csPCR(pe) methods give significant results,  
 aligning with biomedical literature.

Table 3:  $p$ -values of different methods on COVID-19 dataset (mortality 30

Method	csPCR	csPCR(pe)
$p$ -value	0.029	0.013

535

536 **NeurIPS Paper Checklist**

537 **1. Claims**

538 Question: Do the main claims made in the abstract and introduction accurately reflect the  
539 paper's contributions and scope?

540 Answer: [\[Yes\]](#)

541 Justification: In the abstract and introduction, we mentioned that we propose a novel  
542 Covariate Shift Corrected Pearson Chi-squared Conditional Randomization (csPCR) test  
543 and discussed our methodological, theoretical, and empirical contributions.

544 Guidelines:

- 545 • The answer NA means that the abstract and introduction do not include the claims  
546 made in the paper.
- 547 • The abstract and/or introduction should clearly state the claims made, including the  
548 contributions made in the paper and important assumptions and limitations. A No or  
549 NA answer to this question will not be perceived well by the reviewers.
- 550 • The claims made should match theoretical and experimental results, and reflect how  
551 much the results can be expected to generalize to other settings.
- 552 • It is fine to include aspirational goals as motivation as long as it is clear that these goals  
553 are not attained by the paper.

554 **2. Limitations**

555 Question: Does the paper discuss the limitations of the work performed by the authors?

556 Answer: [\[Yes\]](#)

557 Justification: We examine the performance of the algorithms under model misspecification  
558 in Section 4.

559 Guidelines:

- 560 • The answer NA means that the paper has no limitation while the answer No means that  
561 the paper has limitations, but those are not discussed in the paper.
- 562 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 563 • The paper should point out any strong assumptions and how robust the results are to  
564 violations of these assumptions (e.g., independence assumptions, noiseless settings,  
565 model well-specification, asymptotic approximations only holding locally). The authors  
566 should reflect on how these assumptions might be violated in practice and what the  
567 implications would be.
- 568 • The authors should reflect on the scope of the claims made, e.g., if the approach was  
569 only tested on a few datasets or with a few runs. In general, empirical results often  
570 depend on implicit assumptions, which should be articulated.
- 571 • The authors should reflect on the factors that influence the performance of the approach.  
572 For example, a facial recognition algorithm may perform poorly when image resolution  
573 is low or images are taken in low lighting. Or a speech-to-text system might not be  
574 used reliably to provide closed captions for online lectures because it fails to handle  
575 technical jargon.
- 576 • The authors should discuss the computational efficiency of the proposed algorithms  
577 and how they scale with dataset size.
- 578 • If applicable, the authors should discuss possible limitations of their approach to  
579 address problems of privacy and fairness.
- 580 • While the authors might fear that complete honesty about limitations might be used by  
581 reviewers as grounds for rejection, a worse outcome might be that reviewers discover  
582 limitations that aren't acknowledged in the paper. The authors should use their best  
583 judgment and recognize that individual actions in favor of transparency play an impor-  
584 tant role in developing norms that preserve the integrity of the community. Reviewers  
585 will be specifically instructed to not penalize honesty concerning limitations.

586 **3. Theory Assumptions and Proofs**

587 Question: For each theoretical result, does the paper provide the full set of assumptions and  
588 a complete (and correct) proof?

589  
590  
591  
592  
593  
594  
595  
596  
597  
598  
599  
600  
601  
602  
603  
604  
605  
606  
607  
608  
609  
610  
611  
612  
613  
614  
615  
616  
617  
618  
619  
620  
621  
622  
623  
624  
625  
626  
627  
628  
629  
630  
631  
632  
633  
634  
635  
636  
637  
638  
639  
640  
641  
642

Answer: [\[Yes\]](#)

Justification: The assumptions are listed in Section 3.3 and the proofs are provided in Appendix appendix:proof.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

#### 4. Experimental Result Reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [\[Yes\]](#)

Justification: Details of the simulation studies and real-data application are included in Sections 4, 5 and Appendix B.2.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
  - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
  - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
  - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
  - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

#### 5. Open access to data and code

643 Question: Does the paper provide open access to the data and code, with sufficient instruc-  
644 tions to faithfully reproduce the main experimental results, as described in supplemental  
645 material?

646 Answer: [Yes]

647 Justification: Replication code for our simulation studies is submitted as supplementary  
648 material. It will also be made publicly available on GitHub once our paper is accepted.  
649 The COVID data set used for the real example in our paper is not publicly available due to  
650 privacy constraints.

651 Guidelines:

- 652 • The answer NA means that paper does not include experiments requiring code.
- 653 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/  
654 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 655 • While we encourage the release of code and data, we understand that this might not be  
656 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not  
657 including code, unless this is central to the contribution (e.g., for a new open-source  
658 benchmark).
- 659 • The instructions should contain the exact command and environment needed to run to  
660 reproduce the results. See the NeurIPS code and data submission guidelines ([https:  
661 //nips.cc/public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 662 • The authors should provide instructions on data access and preparation, including how  
663 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 664 • The authors should provide scripts to reproduce all experimental results for the new  
665 proposed method and baselines. If only a subset of experiments are reproducible, they  
666 should state which ones are omitted from the script and why.
- 667 • At submission time, to preserve anonymity, the authors should release anonymized  
668 versions (if applicable).
- 669 • Providing as much information as possible in supplemental material (appended to the  
670 paper) is recommended, but including URLs to data and code is permitted.

## 671 6. Experimental Setting/Details

672 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-  
673 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the  
674 results?

675 Answer: [Yes]

676 Justification: Details of the simulation studies and real-data application are included in  
677 Sections 4, 5 and Appendix B.2.

678 Guidelines:

- 679 • The answer NA means that the paper does not include experiments.
- 680 • The experimental setting should be presented in the core of the paper to a level of detail  
681 that is necessary to appreciate the results and make sense of them.
- 682 • The full details can be provided either with the code, in appendix, or as supplemental  
683 material.

## 684 7. Experiment Statistical Significance

685 Question: Does the paper report error bars suitably and correctly defined or other appropriate  
686 information about the statistical significance of the experiments?

687 Answer: [Yes]

688 Justification: Type-I errors, power, and p-values are provided in Sections 4 and 5.

689 Guidelines:

- 690 • The answer NA means that the paper does not include experiments.
- 691 • The authors should answer "Yes" if the results are accompanied by error bars, confi-  
692 dence intervals, or statistical significance tests, at least for the experiments that support  
693 the main claims of the paper.

- 694
- 695
- 696
- 697
- 698
- 699
- 700
- 701
- 702
- 703
- 704
- 705
- 706
- 707
- 708
- 709
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
  - The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
  - The assumptions made should be given (e.g., Normally distributed errors).
  - It should be clear whether the error bar is the standard deviation or the standard error of the mean.
  - It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
  - For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
  - If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

## 710 8. Experiments Compute Resources

711 Question: For each experiment, does the paper provide sufficient information on the com-  
712 puter resources (type of compute workers, memory, time of execution) needed to reproduce  
713 the experiments?

714 Answer: [Yes]

715 Justification: We record relevant information in Appendix B.1.

716 Guidelines:

- 717
- 718
- 719
- 720
- 721
- 722
- 723
- 724
- The answer NA means that the paper does not include experiments.
  - The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
  - The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
  - The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

## 725 9. Code Of Ethics

726 Question: Does the research conducted in the paper conform, in every respect, with the  
727 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

728 Answer: [Yes]

729 Justification: Yes, our research conforms to the NeurIPS Code of Ethics in every respect,  
730 including fairness, transparency, privacy, and social responsibility.

731 Guidelines:

- 732
- 733
- 734
- 735
- 736
- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
  - If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
  - The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

## 737 10. Broader Impacts

738 Question: Does the paper discuss both potential positive societal impacts and negative  
739 societal impacts of the work performed?

740 Answer: [Yes]

741 Justification: We have discussed the impact of the paper on the fields of healthcare and  
742 social sciences.

743 Guidelines:

- 744
- The answer NA means that there is no societal impact of the work performed.

- 745
- 746
- 747
- 748
- 749
- 750
- 751
- 752
- 753
- 754
- 755
- 756
- 757
- 758
- 759
- 760
- 761
- 762
- 763
- 764
- 765
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
  - Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
  - The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
  - The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
  - If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

766 **11. Safeguards**

767 Question: Does the paper describe safeguards that have been put in place for responsible  
768 release of data or models that have a high risk for misuse (e.g., pretrained language models,  
769 image generators, or scraped datasets)?

770 Answer: [NA]

771 Justification: We believe the paper poses no such risks.

772 Guidelines:

- 773
- 774
- 775
- 776
- 777
- 778
- 779
- 780
- 781
- 782
- The answer NA means that the paper poses no such risks.
  - Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
  - Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
  - We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

783 **12. Licenses for existing assets**

784 Question: Are the creators or original owners of assets (e.g., code, data, models), used in  
785 the paper, properly credited and are the license and terms of use explicitly mentioned and  
786 properly respected?

787 Answer: [Yes]

788 Justification: We have cited the relevant papers and packages.

789 Guidelines:

- 790
- 791
- 792
- 793
- 794
- 795
- 796
- The answer NA means that the paper does not use existing assets.
  - The authors should cite the original paper that produced the code package or dataset.
  - The authors should state which version of the asset is used and, if possible, include a URL.
  - The name of the license (e.g., CC-BY 4.0) should be included for each asset.
  - For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.



- 797
- 798
- 799
- 800
- 801
- 802
- 803
- 804
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets) has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
  - For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
  - If this information is not available online, the authors are encouraged to reach out to the asset's creators.

### 805 13. **New Assets**

806 Question: Are new assets introduced in the paper well documented and is the documentation  
807 provided alongside the assets?

808 Answer: [Yes]

809 Justification: We have provided detailed documentation of the newly proposed algorithm.

810 Guidelines:

- 811
- 812
- 813
- 814
- 815
- 816
- 817
- 818
- The answer NA means that the paper does not release new assets.
  - Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
  - The paper should discuss whether and how consent was obtained from people whose asset is used.
  - At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

### 819 14. **Crowdsourcing and Research with Human Subjects**

820 Question: For crowdsourcing experiments and research with human subjects, does the paper  
821 include the full text of instructions given to participants and screenshots, if applicable, as  
822 well as details about compensation (if any)?

823 Answer: [NA]

824 Justification: The paper does not involve crowdsourcing nor research with human subjects.

825 Guidelines:

- 826
- 827
- 828
- 829
- 830
- 831
- 832
- 833
- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
  - Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
  - According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

### 834 15. **Institutional Review Board (IRB) Approvals or Equivalent for Research with Human 835 Subjects**

836 Question: Does the paper describe potential risks incurred by study participants, whether  
837 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)  
838 approvals (or an equivalent approval/review based on the requirements of your country or  
839 institution) were obtained?

840 Answer: [NA]

841 Justification: The paper does not involve crowdsourcing nor research with human subjects.

842 Guidelines:

- 843
- 844
- 845
- 846
- 847
- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
  - Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.

848  
849  
850  
851  
852

- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.