# Collaborative causal inference with a distributed data-sharing management

Mengtong Hu, Xu Shi, and Peter X.-K. Song
Department of Biostatistics, University of Michigan

## Abstract

Data sharing barriers are paramount challenges arising from multicenter clinical trials where multiple data sources are stored in a distributed fashion at different local study sites. Merging such data sources into a common data storage for a centralized statistical analysis requires a data use agreement, which is often time-consuming. Data merging may become more burdensome when causal inference is of primary interest because propensity score modeling involves combining many confounding variables, and systematic incorporation of this additional modeling in meta-analysis has not been thoroughly investigated in the literature. We propose a new causal inference framework that avoids the merging of subject-level raw data from multiple sites but needs only the sharing of summary statistics. The proposed collaborative inference enjoys maximal protection of data privacy and minimal sensitivity to unbalanced data distributions across data sources. We show theoretically and numerically that the new distributed causal inference approach has little loss of statistical power compared to the centralized method that requires merging the entire data. We present large-sample properties and algorithms for the proposed method. We illustrate its performance by simulation experiments and a real-world data example on a multicenter clinical trial of basal insulin treatment for reducing the risk of post-transplantation diabetes among kidney-transplant patients.

*K*eywords: Collaborative causal inference; Data privacy; Distributed inference; Meta-analysis; Multicenter study.

# 1 Introduction

Estimation of causal effects is the central interest in the analysis of data collected from a multicenter clinical trial (Hernán et al. 2002). This statistical task can be greatly challenged when serious data sharing barriers are present across multiple participating clinical centers and hard to resolve in the short run due to various logistic constraints, such as data security and privacy requirements, and tedious institutional IRB approval procedures (Carter & Ardery 2016, Mello et al. 2013, Coates et al. 2020). Such data-sharing obstacles may be significantly magnified in some trials involving international study sites. Holdups in data procurement can delay the publication of clinical findings, and consequently hold back the delivery of new therapeutics to patients.

Among several solutions available in the literature, meta-analysis is of great popularity. A meta-estimation of treatment effect may be calculated by an inverse-variance weighted average of site-specific treatment effects obtained from individual data sources respectively, termed the classical meta-analysis in this paper (for example, (Cochran 1954)). In this approach or others of the same meta-analytic nature, only site-specific summary statistics, rather than the full subject-level data, are involved in pooling site-specific treatment effects. However, this classical meta-analysis has two significant limitations. First, meta-analysis often concerns a single final estimator, while causal inference typically relies on intermediate steps such as building a propensity score model, the key weighting scheme of critical importance in clinical studies to balance covariate distributions across all sites (Rosenbaum & Rubin 1983). In fact, a typical meta-analysis has limited or no control over model specification at each participating site, and the pooled estimate from each site may suffer narrow interpretability, especially in the case of causal effect, due to lack of a well-defined common estimand as well as a uniform operation for the correction of confounding effects. Second, it may suffer from data attrition due to varying recruitment capacity across study sites leading to small sample sizes and low sample variability at some study sites, which can impair the statistical power of the analysis.

Several methods have been developed to improve the classical meta-analysis approach with different objectives other than estimating causal treatment effects. Jordan et al. (2018) developed a surrogate likelihood framework that communicates locally estimated gradients to update the likelihood function for estimation and inference, and the convergence rate for the estimators are at the order of the local sample size. The framework is later extended by Duan et al. (2018, 2019) for

distributed clinical datasets. A distributed empirical likelihood method designed for unbalanced datasets is recently proposed by Zhou et al. (2017). Recent work by Xiong et al. (2021) considers producing global propensity scores from the locally estimated gradients in a federated learning setting. Han et al. (2021) introduces an adaptive federated procedure of weighing the estimators locally estimated from source sites to augment average treatment effect estimate in a target site, meanwhile allowing potential heterogeneity in covariate distributions. Most of the newly proposed meta-analytic approaches adopt a divide-and-conquer strategy, similar to a parallelized operation that requires reliable local estimates and inferential quantities. Unfortunately, due to various reasons pointed out above, the demand for high-quality numerical results from all local sites is rarely satisfied in practice.

This paper focuses on the development of a more flexible and reliable meta-analysis methodology by overcoming the above-marked impediments to evaluating causal treatment effects through effective data-sharing management. Among multiple possible strategies concerning the calculation and utilization of propensity scores in the case of distributed data, we design four new procedures for effectual communication of summary statistics across study sites. Consequently, we can seamlessly integrate the propensity score calculation and causal treatment effect estimation using a systematic cross-site collaboration. We term this new approach as *collaborative operation of linked analysis (*COLA*)*.

The reason that the COLA method appears more flexible than existing meta-analysis is that it adopts a serial updating machinery (Luo & Song 2020), different from the currently popular parallelized operation, in the cross-site communication of summary statistics. The improved flexibility is achieved by different options of passing information in the COLA machinery to reach a desirable trade-off between information communication cost and statistical estimation efficiency. Figure 1 shows our proposed relays for summary data communications, which leads to a fully efficient inverse probability weighting estimation of treatment effect in the sense that its convergence rate is at the order of the cumulative sample size. We show both theoretically and numerically that the COLA has no loss of statistical power in comparison to the oracle estimation obtained by the centralized analysis that merges data from all sites.

To elucidate key elements and insights of the proposed COLA methodology, we choose a simple but practically important setting of logistic regression with binary outcomes, which is largely motivated by a multicenter clinical study of the Insulin Therapy for the Prevention of New-Onset
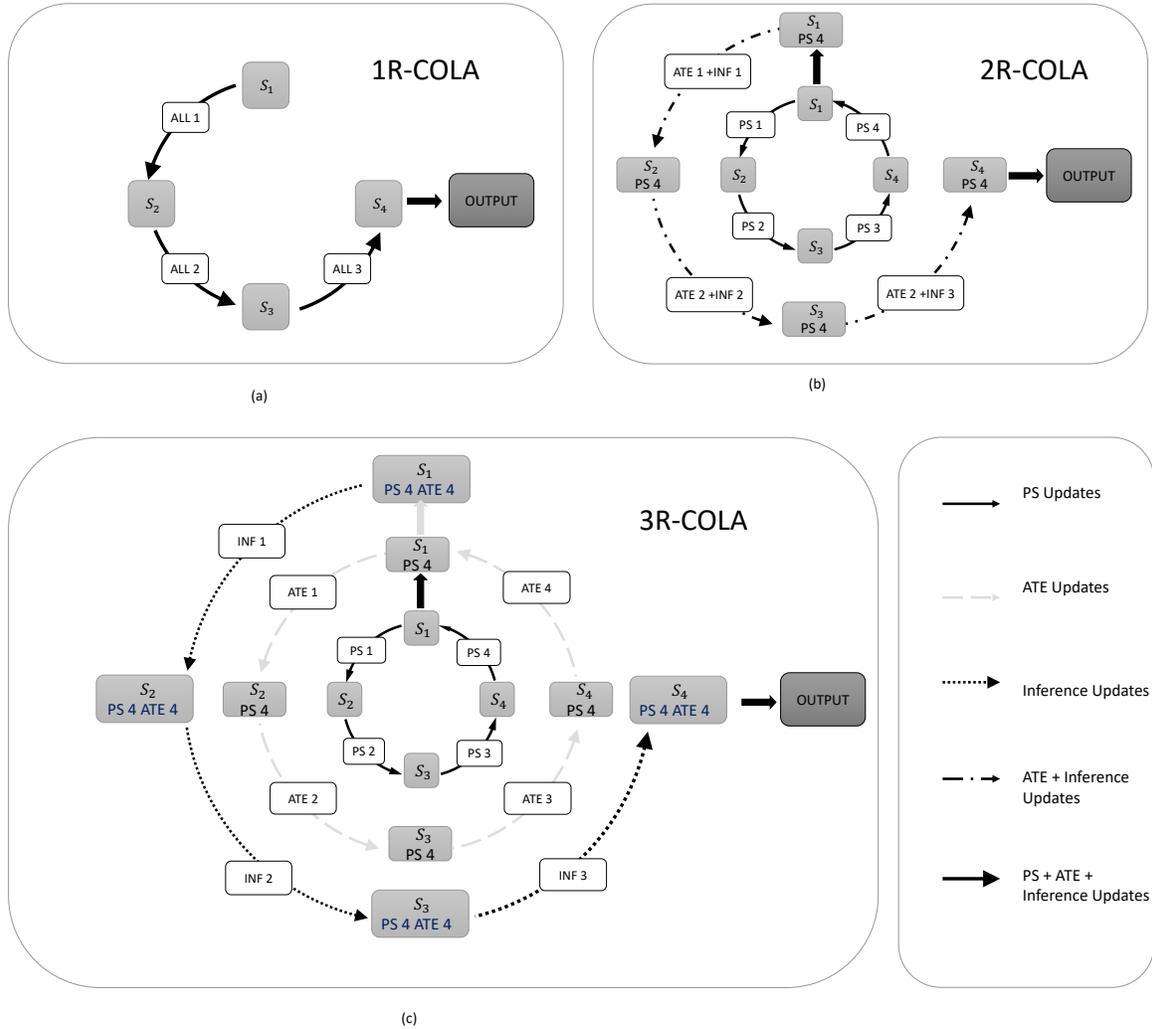
Figure 1: A diagram showing 1R-COLA, 2R-COLA, and 3R-COLA procedures of the collaborative causal inference framework. Different types of arrows indicate the major updates involved in each round as shown in the legend. ATE which is short for Average Treatment Effect denotes the causal effect estimate. Panel a) shows 1R-COLA evolves one round operation that updates PS, ATE, and inferential quantities simultaneously. Panel b) shows 2R-COLA involving two rounds of updates where the first round produces PS estimates, and the second round produces a causal effect and its variance. Panel c) shows 3R-COLA involving three-rounds of updates.

Diabetes after Transplantation (ITP-NODAT) trial. This cross-border clinical study is a multi-center randomized clinical trial conducted at four kidney transplant centers in Barcelona, Berlin, Graz, and Vienna, respectively. The goal of the study is to estimate the average treatment effect of basal insulin intervention on preventing overt post-transplantation diabetes mellitus (PTDM) using data collected from the four hospitals. According to Schwaiger et al. (2021), the collected data exhibit strong imbalances in some covariate distributions and logistic challenges in the creation of a centralized database. Meta-analysis is not the choice to bypass the need for data merging because zero disease cases are observed from the treatment group in Barcelona and likewise, from the control group in Graz. Thus, these two hospitals cannot provide local estimates, and the classical meta-analysis would use results only from two sites (Vienna and Berlin) and would thus be clearly underpowered. To expedite the delivery of clinical findings, it is of great interest to extend the capacity of existing meta-analysis for the needed flexibility of embracing rare outcomes and balancing covariate distributions to retain statistical power and produce an efficient and reliable estimation of the causal effect of the insulin therapy. With little effort, the COLA framework developed for the binary outcome may be extended to other continuous and categorical outcomes in the context of generalized linear models.

The organization of this paper is as follows. Section 2 begins with formulation and model assumptions for estimating causal effects and then introduces the proposed COLA framework. Section 3 presents different options for passing information in the COLA machinery. Section 4 establishes the theoretical guarantees for our proposed methods. We illustrate our proposed methods with simulation studies in Section 5 and an application to the ITP-NODAT data example in Section 6. We make some concluding remarks in Section 7.

## 2  Basic setup

Consider a random sample of $N$ individuals independently sampled from $K$ clinical sites, indexed by $j = 1, \cdots, K$, each site having a sample size of $n_j$. For each individual $i$, we observe an outcome $Y_i$ of interest , a binary treatment indicator $A_i \in \{0, 1\}$, and a vector of baseline covariates $X_i$. Let $I_j = \{n_{j-1} + 1, \ldots, n_{j-1} + n_j\}$ to be the index set for subjects from the $j$th site and we denote the $j$th site-specific data as $S_j = \{(Y_i, A_i, X_i) : i \in I_j\}$. We assume $\{(Y_i, A_i, X_i) : i \in \bigcup_{j=1}^{K} I_j\}$ are independent and identical observations under the same underlying population.

Following the potential outcome framework (Neyman 1923, Rubin 2005), we assume that there exists a pair of potential outcomes $\{Y_i(1), Y_i(0)\}$ for each individual $i$, representing the outcomes had the individual received treatment $(1)$ or control $(0)$. The causal effect is typically defined as a contrast between $\mu_1 = E\{Y(1)\}$ and $\mu_0 = E\{Y(0)\}$, such as the mean difference, $\Delta_D = \mu_1 - \mu_0$, the causal log risk ratio (logRR), $\Delta_{RR} = log\,(\mu_1/\mu_0)$, and the causal log odds ratio (logOR), $\Delta_{OR} = log[\{\mu_1/(1-\mu_1)\}/\{\mu_0/(1-\mu_0)\}]$. To proceed, we postulate the following four standard assumptions for the identification of the causal estimands:

**Assumption 1** (Causal effect identification). *For $a \in \{0,1\}$,*

*(a) Consistency: $Y = Y(a)$ almost surely when $A = a$ .*

*(b) Ignorability: $\{Y(0), Y(1)\} \perp\!\!\!\perp A \mid X$.*

*(c) Positivity: $0 < pr(A = a \mid X) < 1$ for all $a$ almost surely .*

Under Assumption 1, the mean potential outcome is identified as $E[Y(a)] = E[E(Y \mid A = a, X)] = E[I(A = a)Y/P(A = a \mid X)]$, and can be estimated through different methods including the inverse probability of treatment weighting (IPTW), G-computation, and augmented inverse propensity weighting (AIPTW) (Robins et al. 1994).

## 2.1 Inverse Propensity Treatment Weighting estimator

In this paper, we illustrate our method using IPTW estimation, although our method also applies to other causal estimation methods such as G-computation and AIPTW. The propensity score $e(X)$ is the probability of being assigned to the treatment group conditional on the covariates, i.e., $e(X) = pr(A = 1 \mid X)$. The propensity score model for $e(X)$ used in IPTW method needs to be correctly specified for consistent estimation of the causal effect. As such we introduce the following assumption.

**Assumption 2** (Propensity score model). *A propensity score model $e(X; \gamma)$ with a finite-dimensional parameter $\gamma$ is correctly specified for $e(X)$.*

Typically, the propensity score is modelled and estimated by the logistic regression, namely $e(X; \gamma) = \text{logit}(X^{\mathsf{T}}\gamma)$. We begin with a brief review of the classical estimation and inference method based on IPTW in the setting where data from all sites are available and analyzed in

a centralized fashion. We term this situation as the *oracle* setting in this paper. The IPTW method is a two-stage procedure. First, we estimate $\gamma$ as the solution to the logistic model equation $\sum_{i=1}^{N} \Psi^{ps}(A_i, X_i; \gamma) = 0$, denoted by $\hat{\gamma}$, where the kernel function is, suppressing index $i$, $\Psi^{ps}(A, X; \gamma) = X\{A - e(X; \gamma)\}$. Then, we fit a marginal structural model by solving $\sum_{i=1}^{N} \Psi^{\Delta}(A_i, X_i, Y_i; \hat{\gamma}, \beta) = 0$, with $\Psi^{\Delta}(A, X, Y; \gamma, \beta) = (1, A)^{\mathrm{T}} \omega(A, X; \gamma)\{Y - g(\beta_0 + \beta_A A)\}$, where the weight is $\omega(A, X; \gamma) = A/e(X; \gamma) + (1 - A)/\{1 - e(X; \gamma)\}$. The link function $g(.)$ is a known link function (McCullagh & Nelder 2019). For example, when the outcome is binary, one typically uses the logistic link function $g(x) = 1/(1 + e^{-x})$ and the slope parameter, i.e the coefficient for treatment $\beta_A = \Delta_{OR}$. To establish a valid inference, we jointly estimate all model parameters related to propensity scores and the causal effects. Let $\theta = (\gamma, \beta)^{\mathrm{T}}$ where true values are $\theta_0 = (\gamma_0, \beta_0)^{\mathrm{T}}$. Stacking $\Psi^{ps}$ and $\Psi^{\Delta}$, we form a joint estimating function

$$\Psi(\theta) = \Psi(A, X, Y; \theta) = \begin{pmatrix} \Psi^{ps}(A, X; \gamma) \\ \Psi^{\Delta}(A, X, Y; \gamma, \beta) \end{pmatrix}. \tag{1}$$

Under Assumption 1 and 2, the estimating function has mean zero when evaluated at the true value $\theta_0$, i,e., $E\{\Psi(\theta_0)\} = 0$ which will be needed for deriving the asymptotic results in Section 4.

Let $\hat{\theta}^{ora}$ denote the oracle estimator obtained from the centralized analysis which solves $\sum_{i=1}^{N} \Psi_i(\theta) = 0$, where $\Psi_i(\theta) = \Psi(A_i, X_i, Y_i; \theta)$. Since $\Psi(\theta)$ is an unbiased estimating function, under some mild regularity conditions, we have the asymptotic normality for $\hat{\theta}^{ora}$, namely $\sqrt{N}(\hat{\theta}^{ora} - \theta_0) \xrightarrow{d} N(0, J(\theta_0))$ where $J(\theta_0) = H(\theta_0)^{-1} V(\theta_0) \{H(\theta_0)^{-1}\}^{\mathrm{T}}$, where the variability matrix $V(\theta_0) = E_{\theta_0}\{\Psi(\theta_0)\Psi^{\mathrm{T}}(\theta_0)\}$, and the sensitivity matrix $H(\theta_0) = -E_{\theta_0}\{\partial \Psi(\theta_0)/\partial \theta^{\mathrm{T}}\}$. $J$ is the inverse of the Godambe information matrix also named as the sandwich covariance matrix (Godambe 1960, 1991, Stefanski & Boos 2002, Song 2007). We can estimate this asymptotic covariance using their sample counterparts given by $V(\hat{\theta}^{ora}) = \sum_{i=1}^{N} \Psi_i(\theta)\Psi_i^{\mathrm{T}}(\theta)\big|_{\theta=\hat{\theta}^{ora}}$, and $H(\hat{\theta}^{ora}) = \sum_{i=1}^{N} -\partial \Psi_i(\theta)/\partial \theta^{\mathrm{T}}\big|_{\theta=\hat{\theta}^{ora}}$. The resulting centralized oracle estimator $\hat{\theta}^{ora}$ and its covariance serve as the gold standard to contrast the performance of our proposed distributed methods.

## 2.2 Incremental causal effect estimator

We consider a situation of practical importance where pooling data from multiple sites is prohibited. To address this data-sharing challenge, we propose a regularized estimation method, termed Collaborative Operation of Linked Analysis (COLA), which does not require sharing individual-level data but only certain summary statistics across institutes. Note that COLA is not derived in a parallel computing paradigm, different from most of the existing solutions. Specifically, given an order of study sites, an incremental estimator obtained at site $k$, denoted by $\hat{\theta}_k$, is sequentially updated over the first $k$ sites, beginning with $\hat{\theta}_1$ at study site 1. Obviously, $\hat{\theta}_1$ is the same as a local estimator as a root of the estimating equation $\sum_{i \in I_1} \Psi_i(\theta) = 0$ with the local data of $S_1$. We define

$$H_j(\hat{\theta}_j) = \sum_{i \in I_j} -\frac{\partial \Psi_i(\theta)}{\partial \theta^{\mathrm{T}}}\Big|_{\theta = \hat{\theta}_j}, \ V_j(\hat{\theta}_j) = \sum_{i \in I_j} \Psi_i(\theta)\Psi_i^{\mathrm{T}}(\theta)\Big|_{\theta = \hat{\theta}_j}$$

as the sensitivity and variability matrix, respectively, evaluated at a local site $j \in \{1, \ldots, K\}$. The initial estimates of the sensitivity matrix and variability matrix $\{H_1(\hat{\theta}_1), V_1(\hat{\theta}_1)\}$ with the local data $S_1$, are also updated by COLA. After obtaining $\{\hat{\theta}_1, H_1(\hat{\theta}_1), V_1(\hat{\theta}_1)\}$ from site 1, COLA passes these summary statistics to site 2 where the triplet updates $\hat{\theta}_1$ to $\hat{\theta}_2$ by solving the following estimating equation (Luo & Song 2020):

$$\Psi_{n_2}(\hat{\theta}_2) + H_1(\hat{\theta}_1)(\hat{\theta}_1 - \hat{\theta}_2) = 0,$$

where $\Psi_{n_j}(\hat{\theta}_j) = \sum_{i \in I_j} \Psi_i(\hat{\theta}_j)$. Repeating this sequential updating, COLA can be carried out over a sequence of all sites to produce estimators and inferential quantities. In particular, when updating $\theta^{k-1}$ to $\theta^k$ at site $k$, we solve for a root of the following estimating equation:

$$\Psi_{n_k}(\hat{\theta}_k) + \sum_{j=1}^{k-1} H_j(\hat{\theta}_j)(\hat{\theta}_{k-1} - \hat{\theta}_k) = 0. \tag{2}$$

The estimating function in (2) consists of two parts: the first term $\Psi_{n_k}(\hat{\theta}_k)$ is based on the local data $S_k$ at site $k$, and the second term assembles the cumulative summary statistics preceding from all previous $k - 1$ sites. The Newton-Raphson algorithm is applied to numerically find a solution $\hat{\theta}_k$.

For statistical inference, we sequentially compute the cumulative sensitivity and variability

matrices over the first $k$ sites evaluated at a given point estimate. For example, if we update the sensitivity and variability matrices along with the incremental estimates, then at site $k$, the sensitivity matrix is given by $\sum_{j=1}^{k}\{H_j(\hat{\theta}_j)\}$ and the variability matrix is given by $\sum_{j=1}^{k}\{V_j(\hat{\theta}_j)\}$. Then we compute the estimated covariance matrix using the sandwich formula.

# 3   Implementation

This section illustrates three ways to implement the incremental causal estimator as defined in Section 2.2, where the main nuance lies in the trade-off between communication efficiency and finite-sample numerical accuracy. The outputs at each round of updates for the three procedures introduced below are summarized in Fig 2. For ease of illustration, we vary the numbers of "+" in the subscripts of the estimator to correspond to the number of rounds used in each implementation of COLA.

## 3.1   A three-round algorithm

An accurate estimate $\hat{\gamma}$ in the propensity score model is essential for accurate estimation and inference of the causal effect which is related to the performance of logistic regression in propensity score estimation. To implement COLA, we propose a three-round estimation algorithm, denoted by 3R-COLA, for the population-level causal effect estimation, as shown in Fig. 1(c).

Round 1: The first round fits the propensity score model using our sequential method and produces a "global" estimate of the model parameter, executing a full round of sequential updating through all $K$ sites. We output the coefficient estimate $\hat{\gamma}_K^{\ddagger}$ at the last site $K$. By Theorem 3 in Section 4, this $\hat{\gamma}_K^{\ddagger}$ approximates the oracle estimate at the order of $o_p(N^{-1/2})$.

Round 2: The second round estimates the causal effect by the same method. We communicate $\hat{\gamma}_K^{\ddagger}$ back to all local sites so that equation (2) can sequentially update the causal effect $\hat{\beta}_j^{\ddagger} = \hat{\beta}_j(\hat{\gamma}_K^{\ddagger})$. This round outputs the global estimate of causal effect, $\hat{\beta}_K^{\ddagger} = \hat{\beta}_K(\hat{\gamma}_K^{\ddagger})$, which is communicated back to all sites.

Round 3: The third round estimates the asymptotic covariance by updating the cumulative sums $\sum_{j=1}^{K} H_j(\hat{\gamma}_K^{\ddagger}, \hat{\beta}_K^{\ddagger})$ and $\sum_{j=1}^{K} V_j(\hat{\gamma}_K^{\ddagger}, \hat{\beta}_K^{\ddagger})$ sequentially over all $K$ sites.

## 3.2   A two-round algorithm

The 3R-COLA algorithm above may be simplified to a two-round algorithm, denoted by 2R-COLA illustrated in Fig. 1(b). It combines "Round 2" and "Round 3" of 3R-COLA for operation, with some details below. The pseudo-code for 2R-COLA is detailed in the supplementary material.

Round 1: The same "Round 1" of 3R-COLA is used to output $\hat{\gamma}_K^{\ddagger}$ which is the same as $\hat{\gamma}_K^{\ddagger}$.

Round 2: While updating $\hat{\beta}_j^{\ddagger} = \hat{\beta}_j(\hat{\gamma}_K^{\ddagger})$, the covariance is updated simultaneously using current $\hat{\beta}_j^{\ddagger}$, namely $\sum_{j=1}^{k} H_j(\hat{\gamma}_K^{\ddagger}, \hat{\beta}_j^{\ddagger})$ and $\sum_{j=1}^{k} V_j(\hat{\gamma}_K^{\ddagger}, \hat{\beta}_j^{\ddagger})$. Obviously, the current update $\hat{\beta}_j^{\ddagger}$ differs from the output from 3R-COLA $\hat{\beta}_K^{\ddagger}$. The numerical performance of the resulting inference would be different and is illustrated in simulation results.

2R-COLA and 3R-COLA produce the same point estimates for propensity scores model coefficients and causal effect, but the different covariance estimations. Using a fully updated $\hat{\beta}_K^{\ddagger}$ in the co-variance calculation may gain some numerical stability than that of the concurrent estimator $\hat{\beta}_j^{\ddagger}$ using preceding data information available in the incremental updating paradigm. An alternative two-round algorithm named 2R-COLA-INF estimates $\beta$ and $\gamma$ simultaneously in "Round 1" and only computes and updates covariances in a "Round 2". The details for 2R-COLA-INF are given in the supplementary material.


## 3.3   A one-round algorithm

Considering minimizing communication cost, we may further simplify the updating procedure by combining all three rounds into a one-round 1R-COLA algorithm shown in Fig. 1(a). It minimizes between-site communication, at the price of reduced numerical stability and finite-sample per-formance. Instead of communicating global $\hat{\gamma}_K^{\dagger}$ back into equation 2, concurrent estimate $\hat{\gamma}_j^{\dagger}$ is used to compute $\hat{\beta}_j^{\dagger}(\hat{\gamma}_j^{\dagger})$ within the same round. Concurrently, the covariance is updated through $\sum_{j=1}^{K} H_j(\hat{\gamma}_j^{\dagger}, \hat{\beta}_j^{\dagger}(\hat{\gamma}_j^{\dagger}))$, and $\sum_{j=1}^{K} V_j(\hat{\gamma}_j^{\dagger}, \hat{\beta}_j^{\dagger}(\hat{\gamma}_j^{\dagger}))$. At the last site, 1R-COLA produces point estimate $\hat{\gamma}_K^{\dagger}$ and $\hat{\beta}_K^{\dagger} = \hat{\beta}_K(\hat{\gamma}_K^{\dagger})$, with a different covariance estimate compared to 2R-COLA and 3R-COLA.

# 4 Large-sample Properties

Let $N_k$ be the cumulative sample size for the first $k$ sites, $N_k = \sum_{j=1}^{k} n_j$. We choose 1R-COLA to discuss the large sample proprieties of our incremental estimators as $N_k = \sum_{j=1}^{k} n_k \to \infty$ instead of $min_{j \in \{1, \cdots, k\}} n_j \to \infty$ in the parallel computing paradigm. This condition is satisfied when $n_k \to \infty$ at one of the sites, or when the number of sites $k \to \infty$ and the former is the focus of this paper. The asymptotic properties of COLA methods with more than one round can be minimally established with analytic effort. Thus, their proofs are omitted. Denote the $L^2$-norm of a vector $u$ by $\|u\|$. Let $N_\rho(\theta_0) = \{\theta : \|\theta - \theta_0\| \leq \rho\}$, $\rho > 0$ be a compact neighborhood of size $\rho$ around the true value $\theta_0$. In addition to Assumptions $1 - 2$ for causal effects identifiability and estimatability, we assume the following regularity conditions for the estimating function $\Psi$ given in equation (1) to establish some key asymptotic properties.

**Assumption 3** (Regularity Conditions). *We assume the following on estimating function in equation 1*

   *(a) The true value $\theta_0$ is the unique solution to $\lambda(\theta) = E\{\Psi(\theta)\} = 0$.*

   *(b) The estimating function $\Psi(\theta)$ is continuously differentiable for all $\theta$ in the neighborhood $N_\rho(\theta_0)$.*

   *(c) The sensitivity matrix $H(\theta)$ and the variability matrix $V(\theta)$ are positive definite for all $\theta \in N_\rho(\theta_0)$.*

   *(d) The sensitivity matrix $H(\theta)$ is Lipschitz continuous for all $\theta \in N_\rho(\theta_0)$.*

Conditions $3(a) - 3(d)$ are mild regularity conditions needed for legitimate asymptotic behaviors of the COLA estimator $\hat{\theta}_k$ under the classical theory of estimating functions (Song 2007, Tsiatis 2006). Condition 3(d) is satisfied usually for the generalized linear models (McCullagh & Nelder 2019).

**Theorem 1.** *Under the regularity conditions $3(a) - 3(d)$, the COLA estimator $\hat{\theta}_k$ is consistent for the true value $\theta_0$, i.e. $\hat{\theta}_k \xrightarrow{p} \theta_0$, as $N_k \to \infty$.*

**Theorem 2.** *Under the regularity conditions $3(a) - 3(d)$, the COLA estimator $\hat{\theta}_k$ is asymptotically normally distributed, i.e., $\sqrt{N_k}(\hat{\theta}_k - \theta_0) \xrightarrow{d} N(0, J(\theta_0))$, as $N_k \to \infty$.*

**Theorem 3.** *Under the regularity conditions 3(a) − 3(d), the* COLA *estimator* $\hat{\theta}_k$ *and the oracle estimator* $\hat{\theta}^{ora}$ *are asymptotically equivalent, in the sense that* $\|\hat{\theta}_k - \hat{\theta}^{ora}\|^2 = o_p(N_k^{-1})$ *as* $N_k \to \infty$.

Theorem 3 implies that the asymptotic difference between $\hat{\theta}_k$ and $\hat{\theta}^{ora}$ is $o_p(N_k^{-1/2})$, and thus they are stochastically equivalent in the sense that they have the same asymptotic normal distribution. This implies that the COLA estimator is fully efficient.

# 5  Simulation Experiments

We evaluate the finite-sample performance of the proposed collaborative causal inference method, comparing the above 3R-COLA, 2R-COLA, and 1R-COLA procedures with the classical meta-analysis (the inverse-variance weighted meta method (Cochran 1954)) and the oracle estimation (i.e., the gold standard obtained by the centralized analysis). To mimic the motivating data example of the ITP-NODAT trial, we consider a binary outcome and estimate the causal log odds ratio. Additional simulations for continuous and count outcomes are included in the supplementary material. We first generate the full data under an assumed model and then split the data into 5 subsets, one for a study site. We include three continuous variables $X_1, X_2$, and $X_3$ independently drawn from standard normal distribution $N(0, 1)$, and two independent binary covariates $X_4$ and $X_5$ from Bernoulli distribution with success probability of 0.5 and 0.6, respectively. We use $X$ to denote the vector of all five covariates. The treatment, $A$, follows Bernoulli distribution with $\mathrm{pr}(A = 1 \mid X) = \mathrm{expit}(0.5 + 0.3X_1 + 0.3X_2 + 0.5X_3 + 0.5X_4 + 0.3X_5)$, where $\mathrm{expit}(x) = 1/(1 + e^{-x})$. The outcome, $Y$, is drawn from a Bernoulli distribution with $\mathrm{pr}(Y = 1 \mid A, X) = \mathrm{expit}(-2.75 + 0.4A + 0.3X_1 + 0.5X_2 + 0.3X_3 + 0.3X_4 + 0.5X_5)$. The causal log odds ratio is estimated as $0.364$ using the Monte-Carlo simulation of $1,000,000$ random samples, and the proportion of cases ($Y = 1$) is approximately $30\%$. To simulate the scenarios of both unequal and equal proportions of cases across sites, we generate group indicators $I(j = 5)$ which follows the Bernoulli distribution with $\mathrm{pr}\{I(j = 5) \mid Y\} = \mathrm{expit}(a + bY)$, the probability that a sample belongs to the 5th site. The parameters $a$ and $b$ are predetermined such that we control the 5th site to have approximately 50 samples, of which $5\%$ or $30\%$ are cases. The sample size $n_5$ at the 5th site may not be exactly 50 because it is round to the next integer. We split the rest of the samples into sites 1 to 4 with sizes of $100, 80, 80$, and $100 - n_5$, respectively. We consider the following scenarios to illustrate the finite performances of the proposed procedures.

Table 1: Simulation results for both equal outcome prevalence across sites and unequal outcome prevalence across sites.

| Methods | FAILS(%) | CP(%) | ABIAS $\times 10^{-3}$ | MSE $\times 10^{-3}$ | ESE $\times 10^{-3}$ |
|---|---|---|---|---|---|
| | | $n = (100, 80, 80, 50, 50)$, $P(Y = 1) \approx 0.05$ at the 5th site | | | |
| Oracle | 0.00 | 94.5 | 214 | 262 | 270 |
| 3R-COLA | 0.00 | 94.5 | 213 | 261 | 268 |
| 2R-COLA | 0.01 | 93.0 | 213 | 244 | 268 |
| 2R-COLA-INF | 0.00 | 93.2 | 222 | 260 | 282 |
| 1R-COLA | 0.01 | 92.4 | 222 | 249 | 282 |
| Meta | 58.49 | 91.6 | 237 | 261 | 297 |
| | | $n = (100, 80, 80, 50, 50)$ $P(Y = 1) \approx 0.3$ at the 5th site | | | |
| Oracle | 0.00 | 94.7 | 217 | 264 | 273 |
| 3R-COLA | 0.00 | 94.7 | 216 | 262 | 271 |
| 2R-COLA | 0.05 | 94.8 | 216 | 262 | 271 |
| 2R-COLA-INF | 0.00 | 93.7 | 225 | 262 | 283 |
| 1R-COLA | 0.05 | 94.1 | 225 | 266 | 283 |
| Meta | 3.86 | 91.3 | 238 | 258 | 298 |

FAILS the number of non-convergence for incremental methods and the traditional meta-analysis method over $30,000$ replications; CP, coverage probability; ABIAS, average absolute bias; MSE, median estimated standard error of the estimates; ESE empirical standard error.

**Scenario 1.** *Cases are unequally distributed among all sites, and the 5th site has a small proportion of cases, namely 5%, while holding the overall cases rate constant.*

**Scenario 2.** *Cases are equally distributed among all sites at approximately 30% at each site.*

**Scenario 3.** *The number of sites $K$ varies to be $5, 10,$ and $15$ and $n_5 = 50$ to mimic the setting in the motivating example.*

Due to space limit, we report simulation results from a typical case under Scenarios 1 and 2 with sample sizes of $100, 80, 80, 50, 50$ at each of the five sites in Table 1. Results from scenario 3 are presented in the supplementary material. Simulation results are summarized over $30,000$ replications limited to those where all algorithms converge successfully in order to make a fair comparison to the classical meta and oracle estimations. We evaluate the estimation and inference performances for different methods of estimating the causal effects for Scenario 1 in the first sub-table of Table 1 and we make the following observations: (i) The classical meta-analysis suffers substantial numerical failures due to low proportions of cases at site 5, evident by the fact that $58.49\%$ of the

simulation replicates fail to reach convergence. In contrast, the rate of failures ($< 0.01\%$) of the proposed incremental methods is significantly less than the meta-analysis method. (ii) The oracle estimation and the 3R-COLA estimation produce very close coverage probability (CP) and average absolute bias (ABIAS). This fully confirms the theoretical results given in Section 4. By contrast, the CP of the meta-analysis estimation is at $91.6\%$, much lower than the $95\%$ nominal level, and the ABIAS and empirical standard error (ESE) are larger than those of the other competing methods. (iii) The average bias of 2R-COLA is the same as that of 3R-COLA due to the same fully updated estimate $\hat{\gamma}_K^{\ddagger}$, , equivalently $\hat{\gamma}_K^{\ddagger}$, being used to update the casual effect in "Round 2" of each method.

In Scenario 2, in which site 5 has a comparable case rate to other sites ($30\%$), in addition to the observations we made above, the coverage probability of both 2R-COLA and 1R-COLA is close to that of 3R-COLA and the oracle estimator. Overall, in both Scenarios 1 and 2, little statistical efficiency is lost by 3R-COLA in the estimation of causal log odds ratio compared to the oracle method even when the outcome is severely unevenly distributed.

The previous simulation results are conducted under a relatively ideal situation when we know which site deviates the most from the overall distribution and place the site with the worst data quality as the last site in the sequence of updating. To gain some insight into how the ordering of sites affects the results, we investigate a less desirable situation where we do not know *a priori* which site has the highest likelihood of failure in convergence due to either rare outcomes or skewed covariates. We consider three additional scenarios.

**Scenario 4.** *Vary $pr(X_2 = 1)$ at site 1 to be rare such that site 1 has a rare binary covariate.*

**Scenario 5.** *Rearrange the order of operation in Scenario 4.*

**Scenario 6.** *Vary $pr(Y = 1)$ at site 1 to be rare such that the first site has a rare outcome.*

When the binary $X_2$ is skewed at the first site as considered in Scenario 4, the local propensity score model may run into positivity assumption violations (Assumption 1(c)), and consequently, the PS model fitted within a local site may fail to converge or fail to generate stable and trustworthy results. The left panel in Fig 3 shows that COLA and meta fail to converge $13 \cdot 6\%$ and $16 \cdot 7\%$ of the time respectively when the starting site has poor quality data. To address this issue, one may rearrange the order of sites. In Scenario 5, we switch bad starting sites that violate the positivity

assumption with the latter good large sites that satisfy this assumption, so that the first "good" site after excluding the one that violates the positivity assumption is the new starting site. All COLA methods including 1R-COLA, 2R-COLA, and 3R-COLA perform well and reach CP close to the oracle after the rearrangement due to fully updated propensity score model estimates are used for causal effect estimation and inference. This reordering strategy is not applicable to the meta-analysis method because it does not yield viable local results.

In Scenario 6 in which the positive case rate is low at the first site, a local treatment effect estimate is likely to be unstable. It is interesting to note that, once the COLA updating procedure continues and incorporates data from more sites, the incrementally updated point estimate becomes close to the true value. 3R-COLA performs equally well in comparison to the oracle method as shown in the right panel of Fig 3.

The simulation results above lead to the practical guidelines summarized as follows. In general, when using COLA, we suggest starting with the largest site to take advantage of large sample properties. If the outcome is rare or binary covariates are unbalanced at the largest site, then one can choose 3R-COLA, or combine smaller sites with less severe distribution imbalances to create a new starting site and then proceed with either 2R-COLA or 1R-COLA. 3R-COLA is the top choice if the communication cost is allowed as it provides higher statistical accuracy.

# 6   Data Application Example

We apply the proposed collaborative inference method to analyze data from the Insulin Therapy for the Prevention of New-Onset Diabetes after Transplantation (ITP-NODAT) trial (Schwaiger et al. 2021). Two hundred and thirty-six kidney-transplantation patients are recruited in the study and randomized within each hospital to receive a diabetics preventive treatment, namely basal insulin injection right after kidney transplantation. The primary goal of the ITP-NODAT trial is to estimate the effectiveness of basal insulin intervention in preventing overt post-transplantation diabetes mellitus (PTDM) at month 12 after the randomization. Following the original analysis conducted by Schwaiger et al. (2021), we will estimate the intent-to-treat effect of the basal insulin treatment. The case is defined as $Y = 1$ if a patient receives antidiabetic therapy (which indicates the occurrence of diabetes), has 2-hour plasma glucose $\geq$ 200 mg/dL, or has Hemoglobin A1C (HbA1c) $\geq$ 6.5%; otherwise $Y = 0$. We adjust for the following covariates: age, gender, family

15

Table 2: Propensity score estimates for basal insulin treatment from combined data via centralized analysis and collaborative inference method.

| | "Gold-standard" from combined data | | | Collaborative inference method | | |
| --- | --- | --- | --- | --- | --- | --- |
| | Estimates | Std.Errors | p-value | Estimates | Std.Errors | p-value |
| Gender | 0.13 | 0.30 | 0.72 | 0.13 | 0.30 | 0.67 |
| Age | −0.01 | 0.01 | 0.39 | −0.01 | 0.01 | 0.46 |
| Family diabetes history | 0.72 | 0.44 | 0.10 | 0.72 | 0.44 | 0.10 |
| First transplant | 0.03 | 0.46 | 0.97 | 0.02 | 0.46 | 0.95 |
| Glomerular disease | −0.33 | 0.32 | 0.31 | −0.33 | 0.32 | 0.32 |
| Polycystic kidney disease | −0.41 | 0.36 | 0.24 | −0.42 | 0.35 | 0.25 |
| Living Doner | 0.86 | 0.45 | 0.05 | 0.87 | 0.45 | 0.05 |

We use 2R-COLA algorithm for estimating collaborative inference "global" PS parameters.

history of diabetes, whether the living donor or deceased donor, whether first-time transplantation or repeated transplantation, whether having polycystic kidney diseases, and whether having glomerular diseases. The proportion of missingness in the covariates is mild, thus we conduct a complete-case analysis excluding 42 dropouts and 8 participants following Schwaiger et al. (2021).

We conduct COLA implementing the 2R-COLA algorithm without requiring subject-level data sharing. This analysis is particularly meaningful for the ITP-NODAT trial because pooling data from the four hospitals took three years to complete due to various cross-country data-sharing barriers. We also perform a centralized analysis of the pooled data as the gold standard for benchmarking, as well as the classical meta-analysis based on site-specific estimates for comparison.

The biggest challenge in the data analysis pertains to the unequal proportions of PTDM cases across hospitals, as shown in Section 4 of the supplementary material. There were zero PTDM cases recorded in the treatment group at the Barcelona hospital and zero PTDM cases in the control group at the Graz hospital, which prevent us from getting any site-specific results at Barcelona and Graz, leading to an exclusion of two out of four site-specific estimates. This data attrition is undesirable in classical meta-analysis.

Table 2 presents the estimated coefficients in the propensity score model obtained from the COLA method based on summary statistics and the centralized method based on the pooled data. Then we obtain an inverse probability of treatment weighted estimate of the causal odds ratio and

its 95% confidence interval (CI). The estimated causal odds ratio is $0.37$ (96% CI: $0.15, 0.93$), which is very similar to the gold standard, $0.37$ (96% CI: $0.15, 0.91$). In contrast, the classical meta-analysis is based on site-specific estimates from two out of four study sites due to rare outcomes, and the meta-estimated odds ratio is $0.62$ (96% CI: $0.20, 1.93$) which overestimates the treatment effect by twice than the one from the centralized analysis.

It is evident from our analysis that the basal insulin treatment reduces occurrences of PTDM with an estimated odds ratio that is significantly less than one. This is in agreement with the findings reported in Schwaiger et al. (2021), which performed a standard clinical trial analysis with no considerations of propensity score weighting. Little loss of statistical power occurred in our collaborative inference approach, while thoroughly overcoming data sharing barriers and enjoying the maximal protection of data privacy. Had our method and analysis been available, these important clinical findings could be published a few years earlier to add a new clinical treatment that benefits transplant patients. It is also worth noting that our proposed causal inference approach is not affected by imbalanced distributions of disease cases across study sites, a striking advantage over the classical meta-analysis method.

# 7   Discussion

In this paper, we introduce a collaborative operation of linked analysis (COLA) framework that overcomes data-sharing barriers and provides an efficient population average treatment effect estimate. We show the desirable asymptotic properties of the proposed distributed inference method. We also investigate the finite-sample performance through numerical experiments with four algorithms to implement COLA at different levels of communication costs. We find that little statistical efficiency is lost compared to the centralized method when estimating causal log odds ratios incrementally via 3R-COLA and 2R-COLA procedures. Even when the outcome is severely rare, 3R-COLA achieves similar results as the oracle method. Although we focus our attention on binary outcomes, our COLA framework enjoys the same properties and performance in the numerical illustrations for other outcome types under the generalized linear models as shown in Section 4 of the supplementary material.

Meta-analytic types of causal inference methods in certain parallel computing diagrams can fail at two levels: local sites fail to converge and thus the pooled inference results fail to reflect the true

parameters of the underlying population. Convergence failures occur when some sites do not have enough variability in outcome measurements. In practice, our COLA methods only require the first site to have enough data variability which makes it an appealing method for multi-center clinical trials that involve small study sites. To facilitate COLA in practice, we provide an R package for data analysis and an interactive information communication platform that allows each site to run the R program independently and upload and download the summary statistics via our platform. Meanwhile, the emergence and development of federated learning which allows individual sites to collaboratively conduct data analysis while mitigating data privacy risks can be another possible solution to help us build an automated privacy-preserving software (Li, Sahu, Talwalkar & Smith 2020, Kairouz et al. 2021, Li, Fan, Tse & Lin 2020).

The current COLA framework is developed under a set of reasonable assumptions for identifiability and large-sample properties. Future work can potentially relax those assumptions and focus on examining the model's robustness in estimation by flagging out incompatible or outlying local data sites and reducing communication burdens between sites by allowing a varying control of data privacy in different parallel problems. In recent causal inference method development, a line of research is primarily aimed at combining multiple datasets collected by different designs from potentially heterogeneous populations (Yang & Ding 2020, Wang et al. 2020, Bareinboim & Pearl 2016, Shi et al. 2021). Most data fusion methods estimate causal treatment effect by incorporating patient-level data from auxiliary data sources into the main data source without consideration of data privacy issues. We plan to take advantage of the privacy-preserving nature of COLA and extend it to data fusion problems. Another interesting potential extension of our method is to transport our COLA estimation which targets the population underlying the current multi-center clinical trials to a new population (Dahabreh et al. 2020, Han et al. 2021).
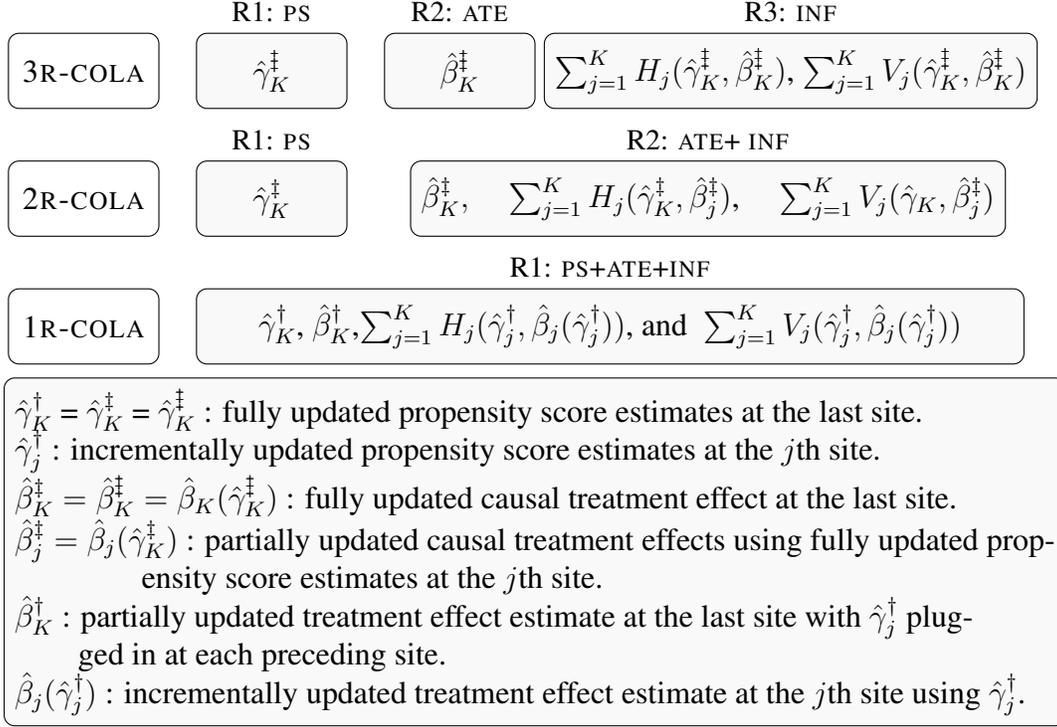
|  | R1: PS | R2: ATE | R3: INF |
|---|---|---|---|

**3R-COLA** $\quad$ $\hat{\gamma}_K^{\ddagger}$ $\quad$ $\hat{\beta}_K^{\ddagger}$ $\quad$ $\sum_{j=1}^{K} H_j(\hat{\gamma}_K^{\ddagger}, \hat{\beta}_K^{\ddagger}), \sum_{j=1}^{K} V_j(\hat{\gamma}_K^{\ddagger}, \hat{\beta}_K^{\ddagger})$

R1: PS $\qquad\qquad$ R2: ATE+ INF

**2R-COLA** $\quad$ $\hat{\gamma}_K^{\ddagger}$ $\quad$ $\hat{\beta}_K^{\ddagger}, \quad \sum_{j=1}^{K} H_j(\hat{\gamma}_K^{\ddagger}, \hat{\beta}_j^{\ddagger}), \quad \sum_{j=1}^{K} V_j(\hat{\gamma}_K, \hat{\beta}_j^{\ddagger})$

R1: PS+ATE+INF

**1R-COLA** $\quad$ $\hat{\gamma}_K^{\dagger}, \hat{\beta}_K^{\dagger}, \sum_{j=1}^{K} H_j(\hat{\gamma}_j^{\dagger}, \hat{\beta}_j(\hat{\gamma}_j^{\dagger})),$ and $\sum_{j=1}^{K} V_j(\hat{\gamma}_j^{\dagger}, \hat{\beta}_j(\hat{\gamma}_j^{\dagger}))$

$\hat{\gamma}_K^{\dagger} = \hat{\gamma}_K^{\ddagger} = \hat{\gamma}_K^{\ddagger}$ : fully updated propensity score estimates at the last site.
$\hat{\gamma}_j^{\dagger}$ : incrementally updated propensity score estimates at the $j$th site.
$\hat{\beta}_K^{\ddagger} = \hat{\beta}_K^{\ddagger} = \hat{\beta}_K(\hat{\gamma}_K^{\ddagger})$ : fully updated causal treatment effect at the last site.
$\hat{\beta}_j^{\ddagger} = \hat{\beta}_j(\hat{\gamma}_K^{\ddagger})$ : partially updated causal treatment effects using fully updated propensity score estimates at the $j$th site.
$\hat{\beta}_K^{\dagger}$ : partially updated treatment effect estimate at the last site with $\hat{\gamma}_j^{\dagger}$ plugged in at each preceding site.
$\hat{\beta}_j(\hat{\gamma}_j^{\dagger})$ : incrementally updated treatment effect estimate at the $j$th site using $\hat{\gamma}_j^{\dagger}$.

Figure 2: Final outputs obtained after each round of update by four COLA methods

Rare binary Covariates, Common incidence Rate, order 1,2,3,4,5 | Rare binary Covariates, Common incidence Rate, with a revised order | Common binary Covariates, Rare incidence Rate, order 1,2,3,4,5

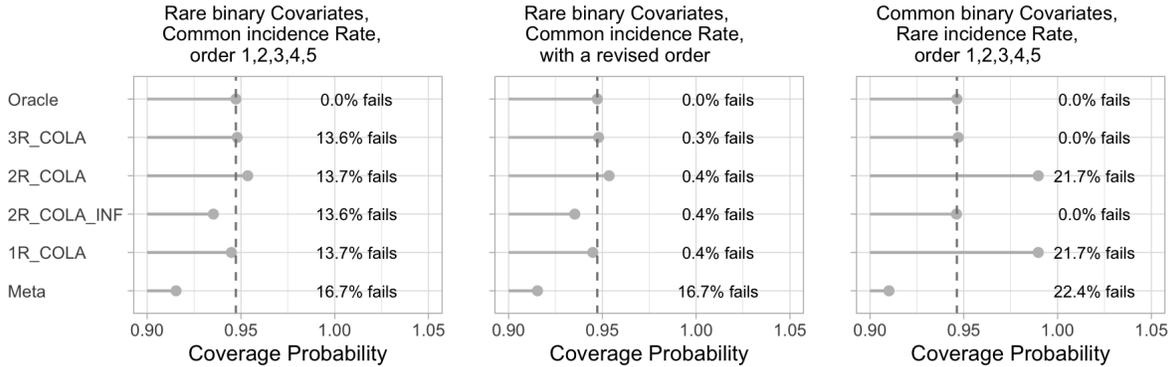| Method | Panel 1 | Panel 2 | Panel 3 |
|---|---|---|---|
| Oracle | 0.0% fails | 0.0% fails | 0.0% fails |
| 3R_COLA | 13.6% fails | 0.3% fails | 0.0% fails |
| 2R_COLA | 13.7% fails | 0.4% fails | 21.7% fails |
| 2R_COLA_INF | 13.6% fails | 0.4% fails | 0.0% fails |
| 1R_COLA | 13.7% fails | 0.4% fails | 21.7% fails |
| Meta | 16.7% fails | 16.7% fails | 22.4% fails |

Coverage Probability (0.90, 0.95, 1.00, 1.05)

Figure 3: A comparison of coverage probabilities for all methods when the binary covariates or outcome incidence are rare based on 30000 replications.

# References

Bareinboim, E. & Pearl, J. (2016), 'Causal inference and the data-fusion problem', *Proceedings of the National Academy of Sciences* **113**(27), 7345–7352.

Carter, B. L. & Ardery, G. (2016), 'Avoiding pitfalls with implementation of randomized controlled multicenter trials: strategies to achieve milestones', *Journal of the American Heart Association* **5**(12), e004432.

Coates, E. C., Mann-Salinas, E. A., Caldwell, N. W. & Chung, K. K. (2020), 'Challenges associated with managing a multicenter clinical trial in severe burns', *Journal of Burn Care & Research* **41**(3), 681–689.

Cochran, W. G. (1954), 'The combination of estimates from different experiments', *Biometrics* **10**(1), 101–129.

Dahabreh, I. J., Petito, L. C., Robertson, S. E., Hernán, M. A. & Steingrimsson, J. A. (2020), 'Toward causally interpretable meta-analysis: Transporting inferences from multiple randomized trials to a new target population', *Epidemiology* **31**(3), 334–344.

Duan, R., Boland, M. R., Moore, J. H. & Chen, Y. (2018), Odal: A one-shot distributed algorithm to perform logistic regressions on electronic health records data from multiple clinical sites, *in* 'BIOCOMPUTING 2019: Proceedings of the Pacific Symposium', World Scientific, pp. 30–41.

Duan, R., Ning, Y. & Chen, Y. (2019), 'Heterogeneity-aware and communication-efficient distributed statistical inference', *arXiv preprint arXiv:1912.09623* .

Godambe, V. P. (1960), 'An optimum property of regular maximum likelihood estimation', *The Annals of Mathematical Statistics* **31**(4), 1208–1211.

Godambe, V. P. (1991), *Estimating functions*, Oxford University Press.

Han, L., Hou, J., Cho, K., Duan, R. & Cai, T. (2021), 'Federated adaptive causal estimation (face) of target treatment effects', *arXiv preprint arXiv:2112.09313* .

Hernán, M. A., Brumback, B. A. & Robins, J. M. (2002), 'Estimating the causal effect of zidovudine on cd4 count with a marginal structural model for repeated measures', *Statistics in medicine* **21**(12), 1689–1709.

Jordan, M. I., Lee, J. D. & Yang, Y. (2018), 'Communication-efficient distributed statistical inference', *Journal of the American Statistical Association* .

Kairouz, P., McMahan, H. B., Avent, B., Bellet, A., Bennis, M., Bhagoji, A. N., Bonawitz, K., Charles, Z., Cormode, G., Cummings, R. et al. (2021), 'Advances and open problems in federated learning', *Foundations and Trends® in Machine Learning* **14**(1–2), 1–210.

Li, L., Fan, Y., Tse, M. & Lin, K.-Y. (2020), 'A review of applications in federated learning', *Computers & Industrial Engineering* p. 106854.

Li, T., Sahu, A. K., Talwalkar, A. & Smith, V. (2020), 'Federated learning: Challenges, methods, and future directions', *IEEE Signal Processing Magazine* **37**(3), 50–60.

Luo, L. & Song, P. X.-K. (2020), 'Renewable estimation and incremental inference in generalized linear models with streaming data sets', *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **82**(1), 69–97.

McCullagh, P. & Nelder, J. A. (2019), *Generalized linear models*, Routledge.

Mello, M. M., Francer, J. K., Wilenzick, M., Teden, P., Bierer, B. E. & Barnes, M. (2013), 'Preparing for responsible sharing of clinical trial data'.

Neyman, J. S. (1923), 'On the application of probability theory to agricultural experiments. essay on principles. section 9.(tlanslated and edited by dm dabrowska and tp speed, statistical science (1990), 5, 465-480)', *Annals of Agricultural Sciences* **10**, 1–51.

Robins, J. M., Rotnitzky, A. & Zhao, L. P. (1994), 'Estimation of regression coefficients when some regressors are not always observed', *Journal of the American statistical Association* **89**(427), 846–866.

Rosenbaum, P. R. & Rubin, D. B. (1983), 'The central role of the propensity score in observational studies for causal effects', *Biometrika* **70**(1), 41–55.

Rubin, D. B. (2005), 'Causal inference using potential outcomes: Design, modeling, decisions', *Journal of the American Statistical Association* **100**(469), 322–331.

Schwaiger, E., Krenn, S., Kurnikowski, A., Bergfeld, L., Pérez-Sáez, M. J., Frey, A., Topitz, D., Bergmann, M., Hödlmoser, S., Bachmann, F. & et al. (2021), 'Early postoperative basal insulin therapy versus standard of care for the prevention of diabetes mellitus after kidney transplantation: A multicenter randomized trial', *Journal of the American Society of Nephrology* **32**(8), 2083–2098.

Shi, X., Pan, Z. & Miao, W. (2021), 'Data integration in causal inference', *arXiv preprint arXiv:2110.01106* .

Song, P. X.-K. (2007), *Correlated data analysis: modeling, analytics, and applications*, Springer Science & Business Media.

Stefanski, L. A. & Boos, D. D. (2002), 'The calculus of m-estimation', *The American Statistician* **56**(1), 29–38.

Tsiatis, A. A. (2006), 'Semiparametric theory and missing data'.

Wang, C., Lu, N., Chen, W.-C., Li, H., Tiwari, R., Xu, Y. & Yue, L. Q. (2020), 'Propensity score-integrated composite likelihood approach for incorporating real-world evidence in single-arm clinical studies', *Journal of Biopharmaceutical Statistics* **30**(3), 495–507.

Xiong, R., Koenecke, A., Powell, M., Shen, Z., Vogelstein, J. T. & Athey, S. (2021), 'Federated causal inference in heterogeneous observational data', *CoRR* **abs/2107.11732**.
    **URL:** *https://arxiv.org/abs/2107.11732*

Yang, S. & Ding, P. (2020), 'Combining multiple observational data sources to estimate causal effects', *Journal of the American Statistical Association* **115**(531), 1540–1554.

Zhou, M., Wang, S. V., Leonard, C. E., Gagne, J. J., Fuller, C., Hampp, C., Archdeacon, P., Toh, S., Iyer, A., Woodworth, T. S. & et al. (2017), 'Sentinel modular program for propensity score-matched cohort analyses: Application to glyburide, glipizide, and serious hypoglycemia', *Epidemiology (Cambridge, Mass.)* **28**(6), 838–846.