

Advancing Interpretable Regression Analysis for Binary Data: A Novel Distributed Algorithm Approach

Jiayi (Jessie) Tong, Assistant Professor

Department of Biostatistics, Johns Hopkins Bloomberg School of Public Health

December 3, 2024

OHDSI Community Call

Background

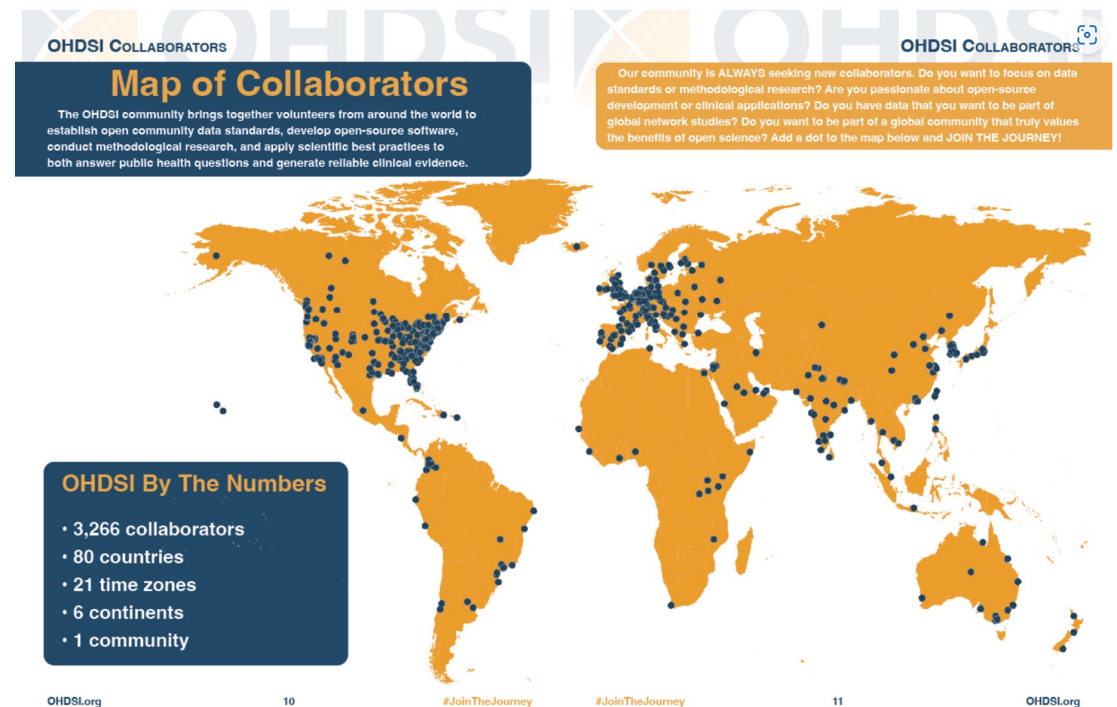
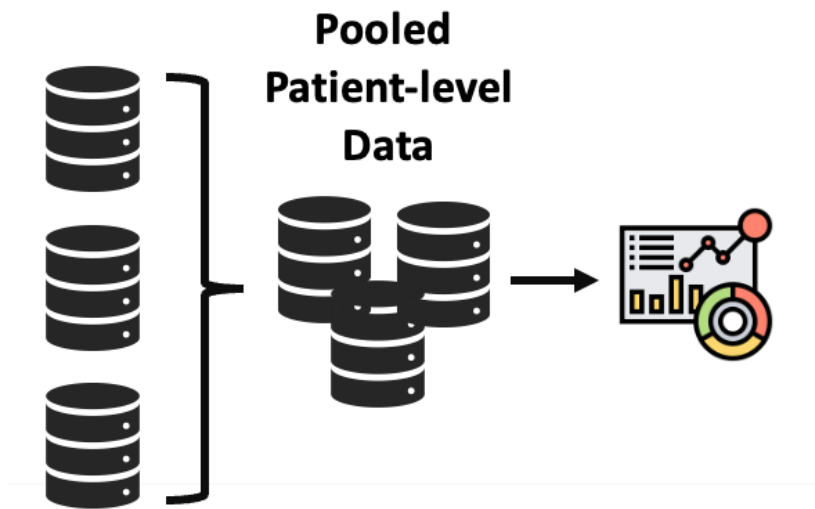
- **Rare diseases -> rare binary outcome**, the sparse data problem is a significant challenge.
- The **lack of sufficient cases** (e.g., patients with disease) in the data leads to **biased estimates** of the effect of a treatment or a medication in observational studies.



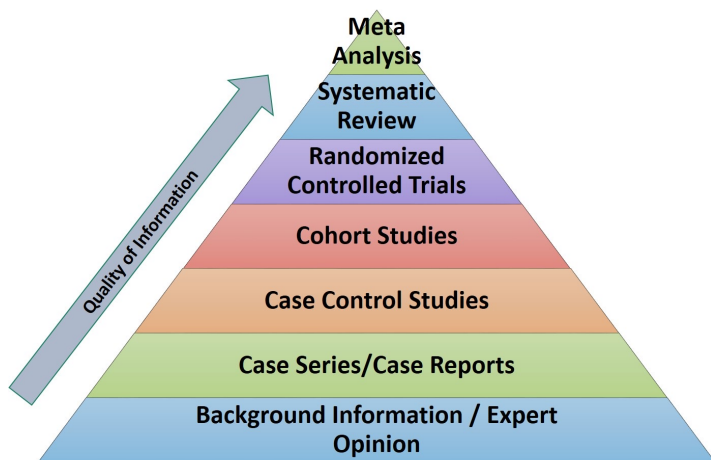
Background

- Multi-site individual patient-level data (IPD) analysis increases the number of cases
- However, it is not feasible when the individual-level data from different studies cannot be shared.

IPD Analysis



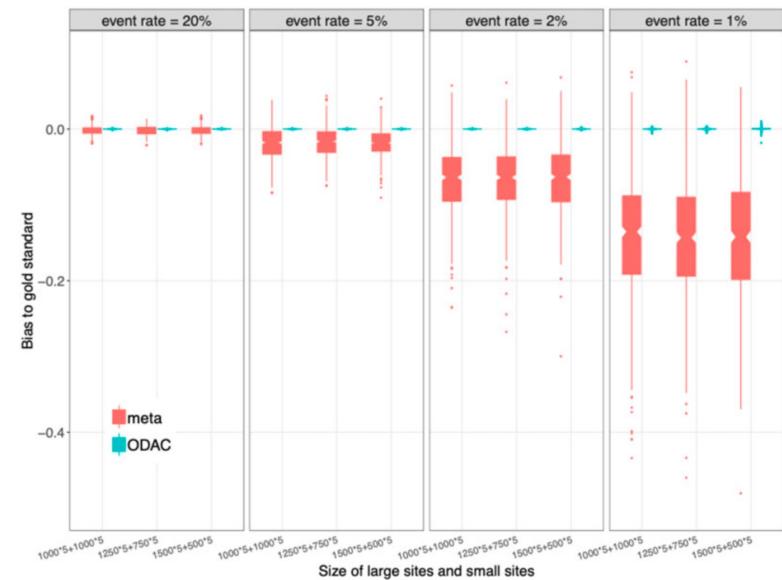
Meta-analysis Method



Meta-analysis (Divide and conquer)



Meta-analysis has low accuracy for rare data



Federated Learning Algorithms

Privacy-Preserving Federated Learning Algorithms

- Enables fitting statistical models in a federated manner
- Requires summary statistics, instead of IPD
- Ensures data privacy and security



Federated Learning Algorithms for Rare Binary Data

Existing Federated Learning Algorithms for Binary Data

Logistic Regression Model -> Odds Ratio (OR)

Aggregated

Aggregated

The choice between **Relative Risk (RR)** and OR has been debated in the literature, with RR being preferred in most prospective studies due to its collapsibility and better interpretation.

Poisson Regression Model -> Relative Risk (RR)

Research

Research and Applications

Shared models without sharing data

Yuan Wu, Xiaoqian Jiang, Jihoon Kim, Lucila Ohno-Machado

distributed algorithm

Rui Duan¹, Mary Regina Boland¹, Zixuan Liu², Yue Liu³, Howard H Chang⁴, Hua Xu⁵, Haitao Chu⁶, Christopher H Schmid⁷, Christopher B Forrest⁸, John H Holmes¹, Martijn J Schuemie⁹, Jesse A Berlin⁹, Jason H Moore¹ and Yong Chen¹



American Journal of Epidemiology
Copyright © 2004 by the Johns Hopkins Bloomberg School of Public Health
All rights reserved

Vol. 159, No. 7
Printed in U.S.A.
DOI: 10.1093/aje/kwh090

A Modified Poisson Regression Approach to Prospective Studies with Binary Data

Guangyong Zou^{1,2}

¹ Robarts Clinical Trials, Robarts Research Institute, London, Ontario, Canada.

² Department of Epidemiology and Biostatistics, University of Western Ontario, London, Ontario, Canada.

Proposed Method

ODAP-B: One-shot Distributed Algorithm of Modified Poisson Regression for Binary Data

- The surrogate likelihood (SL) approach (Jordan et al 2018 JASA)

- For some given initial value $\bar{\beta}$, consider Taylor expansion for the **multi-site likelihood**

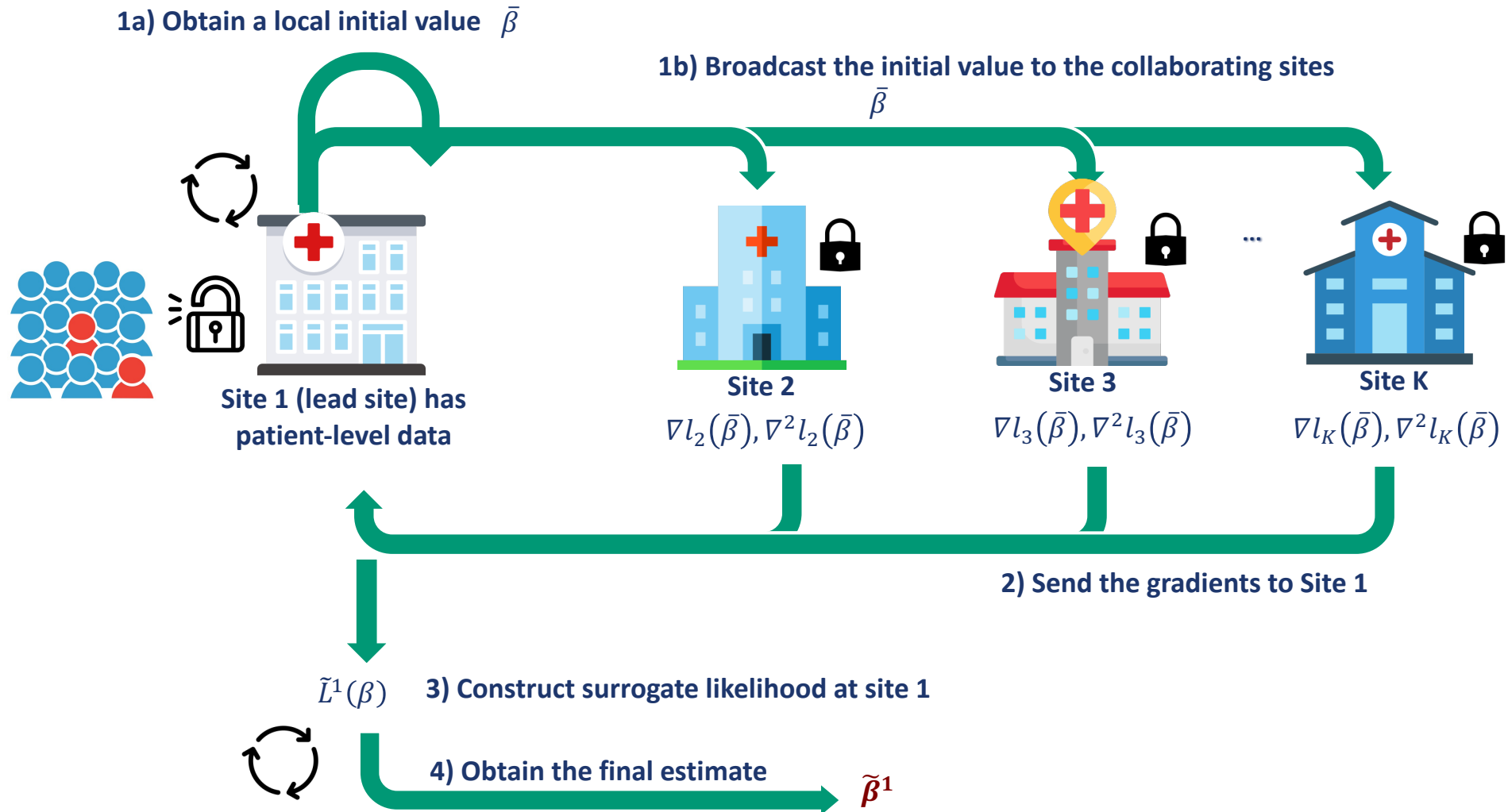
$$L(\beta) = L(\bar{\beta}) + \nabla L(\bar{\beta})^T (\beta - \bar{\beta}) + \sum_{t=2}^{\infty} \frac{1}{t!} \nabla^t L(\bar{\beta}) (\beta - \bar{\beta})^{\otimes t}$$

- For some given initial value $\bar{\beta}$, consider Taylor expansion for the **local likelihood**

$$L_1(\beta) = L_1(\bar{\beta}) + \nabla L_1(\bar{\beta})^T (\beta - \bar{\beta}) + \sum_{t=2}^{\infty} \frac{1}{t!} \nabla^t L_1(\bar{\beta}) (\beta - \bar{\beta})^{\otimes t}$$

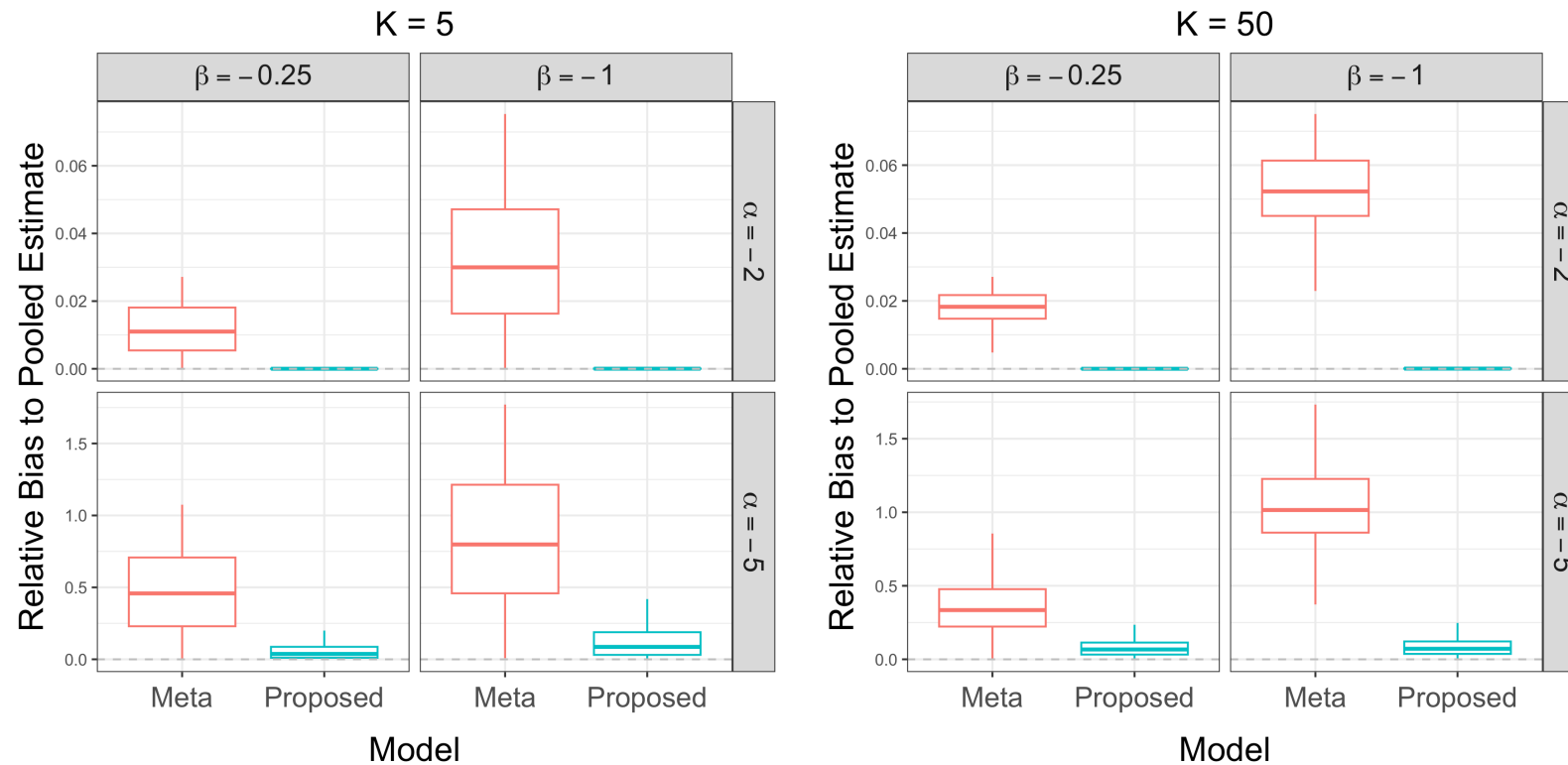
Surrogate likelihood $\tilde{L}^1(\beta) = L_1(\beta) + \{\nabla L(\bar{\beta}) - \nabla L_1(\bar{\beta})\}^T \beta$

$$\nabla L(\bar{\beta}) = \frac{1}{K} \sum \nabla L_j(\bar{\beta}); \quad \tilde{\beta} = \operatorname{argmax}_{\beta} \tilde{L}(\beta)$$



Simulation Study

Number of sites (K)



Real-world Data Application

Investigate the relationship between **COVID-19 viral** (SARS-CoV-2 polymerase chain reaction [PCR] or antigen) **test positivity** and the **symptoms and conditions associated with Long-COVID** in children

- Use case 1: Centralized data from PEDSnet
 - Nine children's hospitals across the nation
 - Sample size: 184,501
- Use case 2: Decentralized datasets
 - 12 sites: PEDSnet (9), OHDSI (2), OneFlorida (3)
 - Sample size: 452,160

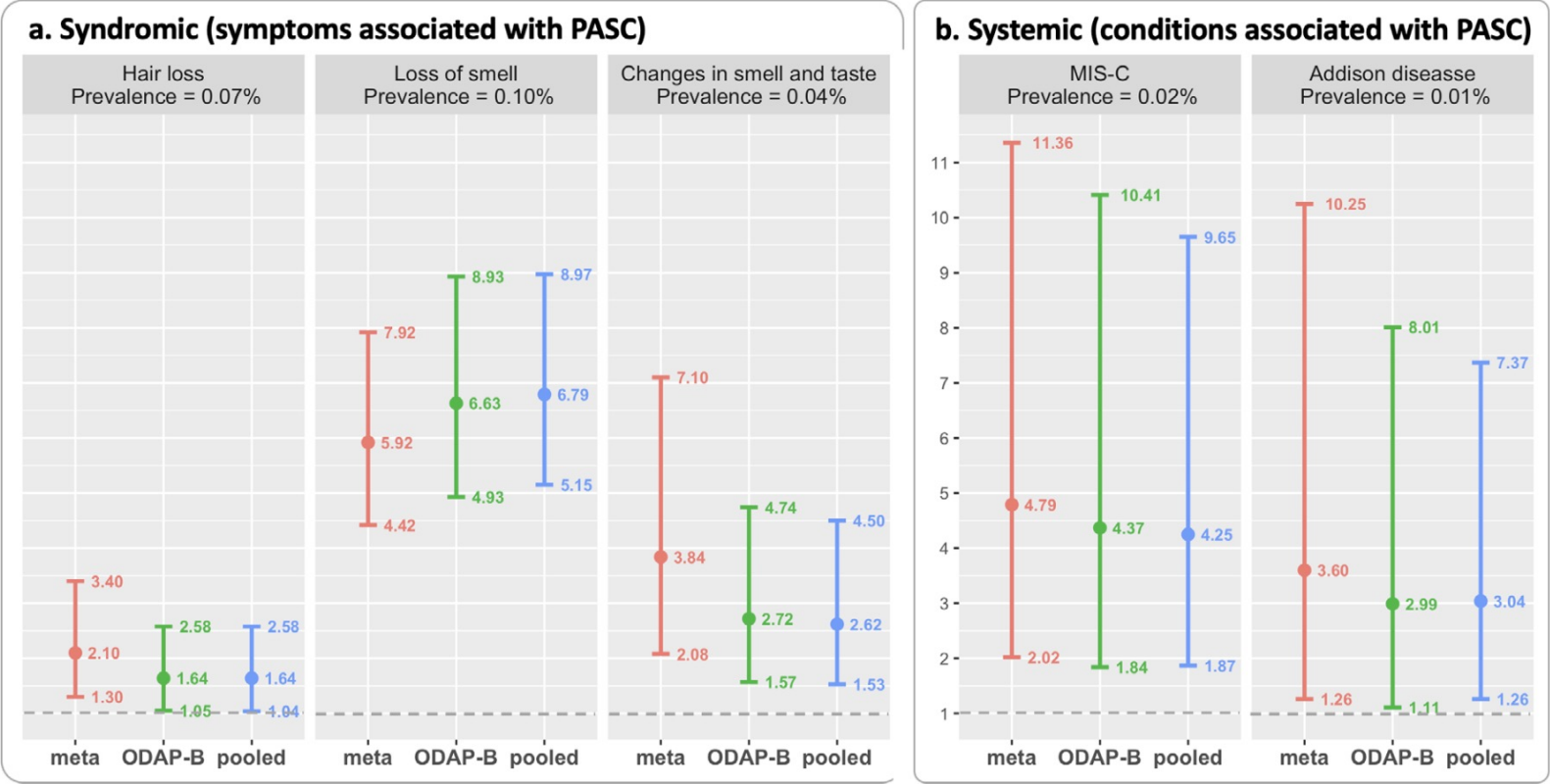


Confounding variables adjusted in the regression model:

age at cohort entrance, sex (male vs. female), race, COVID-19 testing location, diagnosis date of the outcome (i.e., PASC conditions), Pediatric Medical Complexity Algorithm (PMCA).

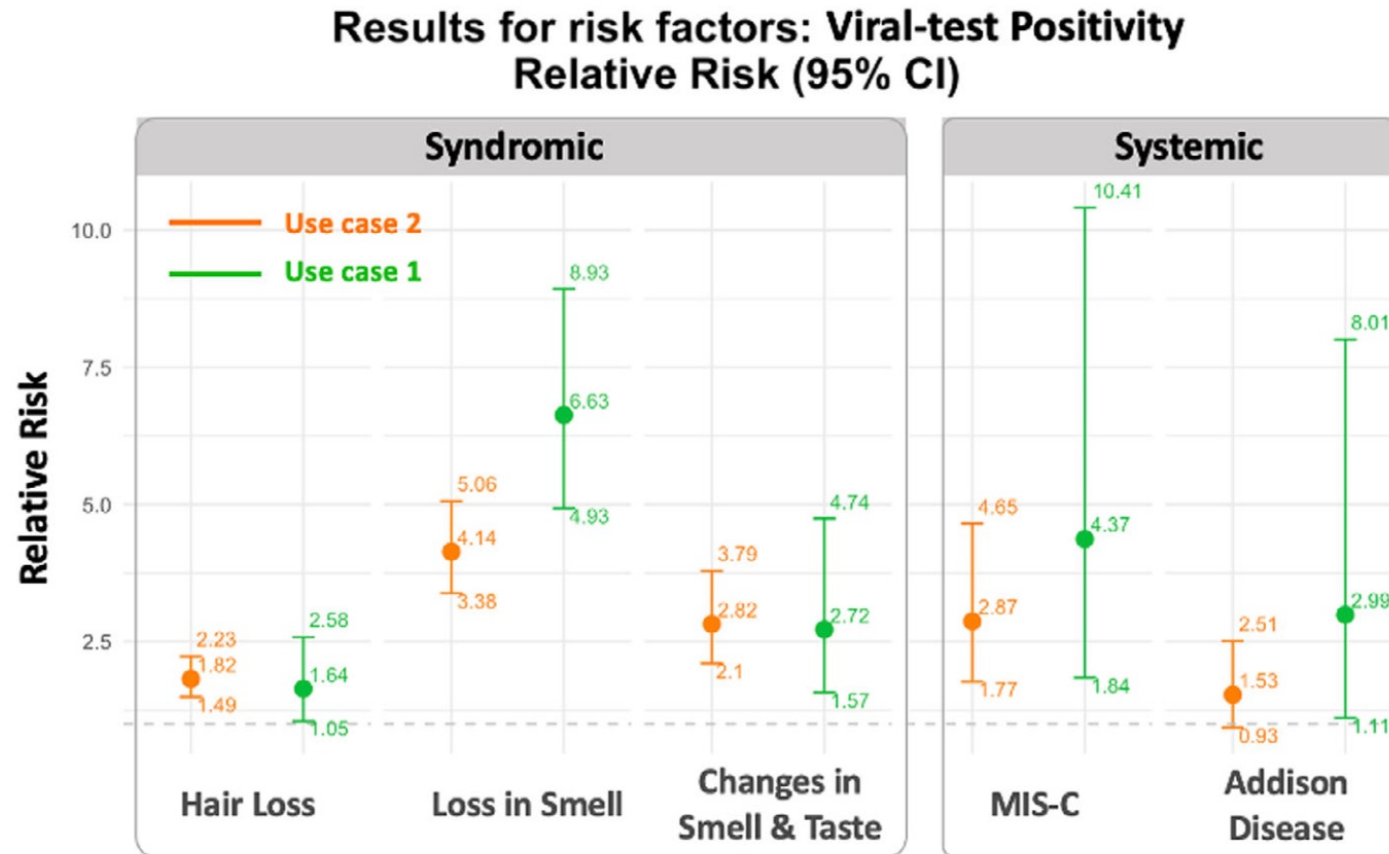
Use Case 1 Result

Results for risk factors: Viral-test Positivity Relative Risk (95% CI)



MIS-C: a multisystem inflammatory syndrome in children

Use Case 1 & Use Case 2 Results Comparison



Summary

- ODAP-B is an effective federated learning algorithm for **Poisson regression** to study **rare binary outcome**.
- ODAP-B provides inference on **adjusted relative risk** with robust variance estimator.
- ODAP-B is **easy to implement** and **applicable** to analyze multi-site data



PDA R Package: 13300+ downloads since 2020



PDA website: <https://pdamethods.org/>



PDA Github Page: <https://github.com/Penncil/pda>



PDA-OTA

Penn security office certified

PDA-OTA: <https://pda-ota.pdamethods.org/>



Statistics in Medicine

WILEY

Statistics
in Medicine

RESEARCH ARTICLE **OPEN ACCESS**

Advancing Interpretable Regression Analysis for Binary Data: A Novel Distributed Algorithm Approach

Jiayi Tong^{1,2}  | Lu Li^{1,3} | Jenna Marie Reps^{4,5,6}  | Vitaly Lorman⁷ | Naimin Jing⁸ | Mackenzie Edmondson⁸ | Xiwei Lou⁹ | Ravi Jhaveri¹⁰ | Kelly J. Kelleher¹¹ | Nathan M. Pajor¹² | Christopher B. Forrest⁷ | Jiang Bian⁹ | Haitao Chu¹³ | Yong Chen^{1,2,14,15,16,17}



Thank you