# Regression Analysis for Probabilistic Cause-of-disease Assignment using Case-control Diagnostic Tests: A Hierarchical Bayesian Approach
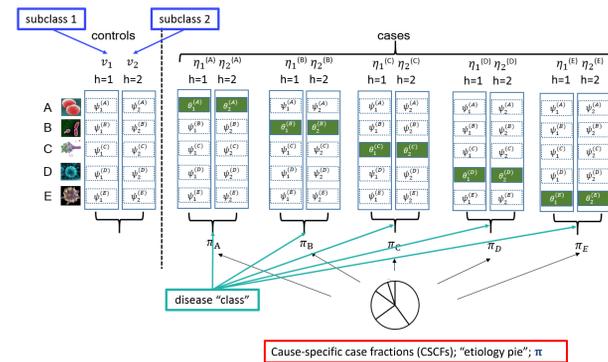
Irena Chen* and Zhenke Wu | Department of Biostatistics and Michigan Institute for Data Science | irena@umich.edu

## Why should you care?

- **Scientific Goal:** For a disease with multiple causes (*not directly observed*):
  1. Assess the effect of explanatory variables on *cause-specific case fractions (CSCFs)*, $\pi(X_i)$, for L causes
     - important for optimizing prevention and treatment strategies
  2. Assess the *overall* CSCFs ($\pi(X_i)$ averaged over the empirical distribution of covariates, $\pi^* = \int \pi(X_i)dG(X_i)$)

- **Data Setting:** *Case-control*, multiple *binary* diagnostic measurements ($M_i$) of disease causes (*with error*)
- **Current approaches to including covariates fall short:**
  - Fully stratified analysis breaks down for sparsely populated strata (Table 1)
  - Unable to **quantify how explanatory variables influence the probabilities** of the unobserved causes

## Existing nested partially-latent class models (npLCM)



- Estimate $\pi, \Theta, \psi$ (CSCFs, true and false positive rates) via **Markov chain Monte Carlo (MCMC)**

- **Proposed regression extension:** let $\nu_k, \eta_k$, (subclass weights) and $\pi$ (CSCFs) depend on observed covariates

## Childhood pneumonia etiology study

- $X_i$ = (age, gender, HIV status, disease severity, enrollment date); $W_i = X_i$ minus the disease severity (case-only)
- **Goal:** Evaluate $\pi_l(X_i)$ of seven single-pathogen and "Not Specified" (NoS) causes of lung infection using nasal pharyngeal polymerase chain reaction (NPPCR) tests.

**Table 1.** The observed counts (frequencies) of controls by age and HIV status; Case counts are further stratified by disease severity (1: yes; 0: no). The marginal case-control positive fractions for each covariate are shown at the bottom. Enrollment date ($t$) is not stratified upon here.

| Age $\geq 1$ | HIV + | # of controls Total: 964 (100) | very severe (case-only) | # of cases Total: 964 (100) |
|---|---|---|---|---|
| 0 | 0 | 548 (56.8) | 0 | 208 (40.2) |
| | | | 1 | 120 (23.2) |
| 1 | 0 | 280 (29.0) | 0 | 69 (13.3) |
| | | | 1 | 32 (6.2) |
| 0 | 1 | 85. (8.8) | 0 | 37 (7.1) |
| | | | 1 | 25 (4.8) |
| 1 | 1 | 51 (5.3) | 0 | 24 (4.6) |
| | | | 1 | 3 (0.6) |
| Case: 24.7% | 17.2% | | 34.7% | |
| Control: 34.3% | 14.1% | | --- | |

## Individual-level estimated probabilities of causes



(a) Cause: RSV    (b) Cause: NoS

**Figure 2.** Individual probability of cause estimates for RSV (left) and NoS (right) differ by age and season among HIV negative and severe pneumonia cases for whom the seven pathogens were all tested negative in the nasopharyngeal specimens.

## What's the regression model?

- **Our method:** A flexible Bayesian model for incorporating regression covariates in a latent class framework
- $W_i$ = vector of covariates that may influence controls ($\nu_k$) and cases ($\eta_k$)
- $X_i$ = vector of covariates that may influence CSCFs ($\pi$)

**Subclass Weight Regression:**
$\nu_k(W_i), \eta_k(W_i)$:

- $h_k(W_i; \Gamma_k^c) = \begin{cases} g(\alpha_{ik}^c)\Pi_{s<k}\{1 - g(\alpha_{is}^c)\} & k < K \\ \Pi_{s<k}\{1 - g(\alpha_{is}^c)\}, & k = K \end{cases}$

  - c designates control ($\nu$) or cases ($\eta$) subclass weights
  - **Stick-breaking parameterization**

- $\alpha_{ik}^c$ is obtained via **Generalized Additive Models:**

$\alpha_{ik}^c = \alpha_k^c(W_i; \Gamma_k^c) = \mu_{k0} + \sum_{j=1}^{q_1} f_{jk}(W_i; \beta_{kj}^c) + \widetilde{W_i}^T \gamma_k^c$

**CSCF regression:**

$\pi_l(X_i) = \frac{\exp(\phi_l(X_i))}{\sum_{l'=1}^{L} \exp(\phi_{l'}(X_i))}$; $\phi_l(X_i) - \phi_L(X_i)$ = log odds of case $i$ in disease class $l$ relative to disease class L

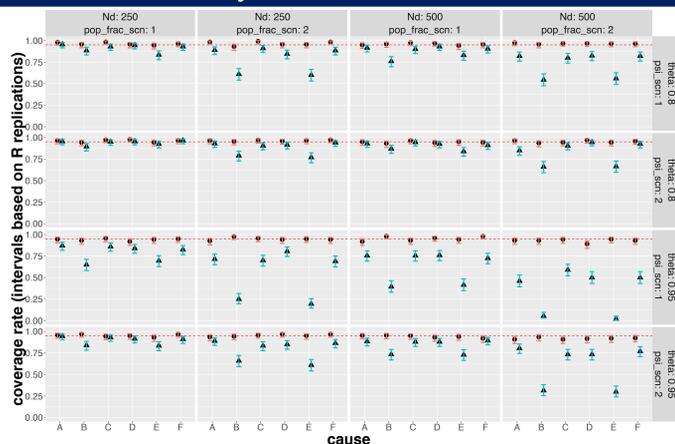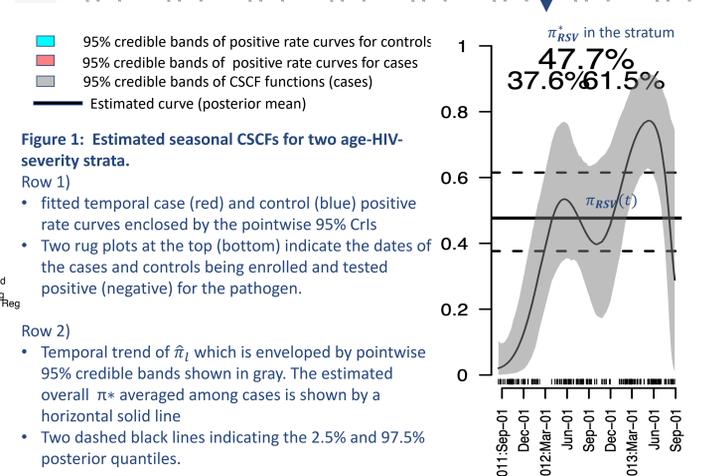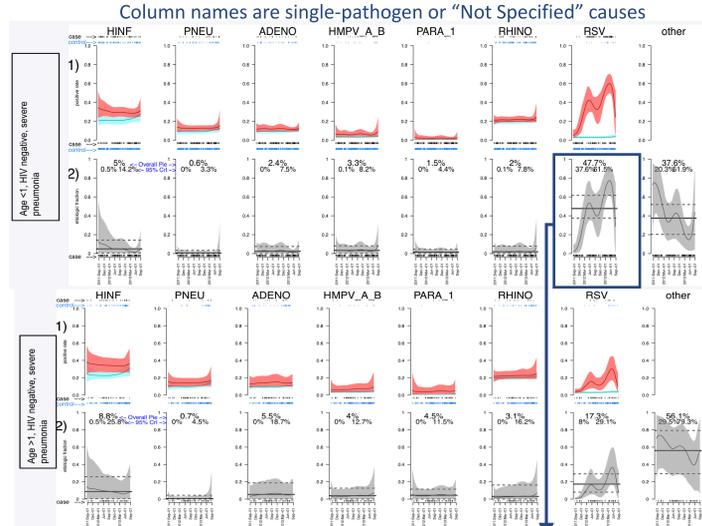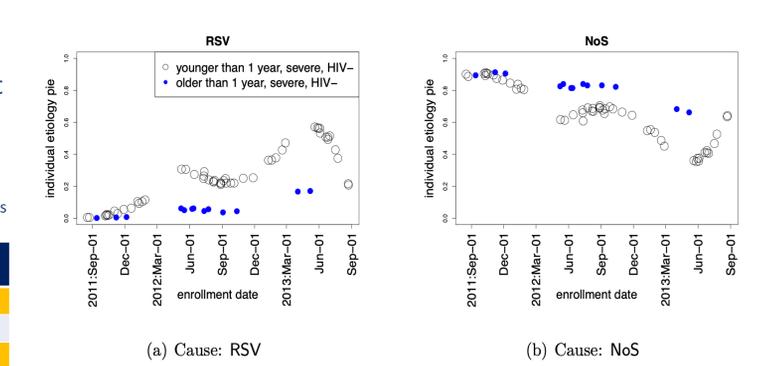Model $\phi_l(X_i)$ as **additive models:**
- $\phi_l(X_i; \Gamma_l^\pi) = \sum_{j=1}^{p_1} f_{lj}^\pi(X_i; \beta_{lj}^\pi) + \widetilde{X_i}^T \gamma_l^\pi$

**Specify shrinkage priors on $\mu_{k0}, f_{jk}, f_{lj}^\pi$ to encourage parsimonious regressions with few effective subclasses (not shown here)**

## Summary

- **Likelihood for controls:**
- $L_0^{reg} = \Pi_{i:\, Y_i=0} \sum_{k=1}^{K} \nu_{ik}(W_i; \Gamma_k^\nu)\Pi(m; \psi_k)$, where $\Pi(m; p)$ is the probability of observing $m$ for $J$ independent $m_j \sim Bernoulli(p_j)$ and $m \in \{0,1\}^J$

- **Likelihood for cases:**
- $L_1^{reg} = \Pi_{i:\, Y_i=1} \sum_{l=1}^{L} \pi_l(X_i; \Gamma_l^\pi) \sum_{k=1}^{K} \eta_{ik}(W_i; \Gamma_k^\eta)\Pi(m; \Theta_k, \psi_k)$

- **Unknown** parameters:
  - Etiology regression coefficients: $\Gamma_k^\pi$
  - Subclass weights: $\{\Gamma_k^\eta\}$ (cases), $\{\Gamma_k^\nu\}$ (controls)
  - True/false positive rates: $(\Theta = \{\theta_k^{(j)}\}, \Psi = \{\psi_k^{(j)}\})$
- Use **MCMC** to approximate posterior distribution

## Simulation: Improved coverage of 95% credible intervals (CrI) for $\pi_l^*$



## Population-level CSCF estimates for two strata

Column names are single-pathogen or "Not Specified" causes



- 95% credible bands of positive rate curves for controls
- 95% credible bands of positive rate curves for cases
- 95% credible bands of CSCF functions (cases)
- Estimated curve (posterior mean)

$\pi_{RSV}^*$ in the stratum
47.7%
37.6%  61.5%

**Figure 1:** Estimated seasonal CSCFs for two age-HIV-severity strata.
Row 1)
- fitted temporal case (red) and control (blue) positive rate curves enclosed by the pointwise 95% CrIs
- Two rug plots at the top (bottom) indicate the dates of the cases and controls being enrolled and tested positive (negative) for the pathogen.

Row 2)
- Temporal trend of $\hat{\pi}_l$ which is enveloped by pointwise 95% credible bands shown in gray. The estimated overall $\pi*$ averaged among cases is shown by a horizontal solid line
- Two dashed black lines indicating the 2.5% and 97.5% posterior quantiles.

## Why does this method matter?

- **$\pi$ regression**: specifies functional dependence of the CSCFs ($\pi$) upon important covariates
  - improves estimation stability for sparsely populated strata (see Table 1) using assumptions such as additivity

- **Utilizes case-control data**: estimate covariate-dependent reference distribution from controls
  - critical for assigning cause-specific probabilities to cases.

- **Correctly assess posterior uncertainty:** Uses informative priors ($\Theta$) only once in the elicited population
  - avoids overly-optimistic uncertainty estimates for $\pi$ (unlike stratified npLCM that reuses these priors)

## Future directions

- Explore more flexible regression models (e.g. Bayesian Additive Regression Trees, or BART)
- Extend to not-prespecified causes (combinatorial space)
- Applications to survey data such as *verbal autopsy*

## Open-source software (R package)

*baker:* **B**ayesian **A**nalysis **K**it for **E**tiology **R**esearch
https://github.com/zhenkewu/baker

**References:**
- **Wu Z. and Chen I** (2019+). Regression Analysis of Dependent Binary Data for Estimating Disease Etiology from Case-Control Studies. Submitted. https://doi.org/10.1101/672808
- **Wu Z** and Zeger SL (2018+). A Bayesian Approach to Restricted Latent Class Models for Scientifically-Structured Clustering of Multivariate Binary Outcomes. https://doi.org/10.1101/400192
- **Wu Z**, Deloria-Knoll M, Zeger S.L, Nested partially latent class models for dependent binary data; estimating disease etiology, Biostatistics, Volume 18, Issue 2, April 2017, Pages 200–213. https://doi.org/10.1093/biostatistics/kxw037
- O'Brien et al. (2019). Causes of severe pneumonia requiring hospital admission in children without HIV infection from Africa and Asia: the PERCH multi-country case-control study. The Lancet (2019): 394 (10200): 757-779. https://doi.org/10.1016/S0140-6736(19)30721-4