

Session 5 of Module 16: Methods for Assessing Immunological Correlates of Risk and Optimal Surrogate Endpoints

Peter Gilbert

Summer Institute in Statistics and Modeling in Infectious Diseases

U of W July 24–26, 2017

Outline of Module 16: Evaluating Vaccine Efficacy

Session 1 (Gabriel)	Introduction to Study Designs for Evaluating VE
Session 2 (Follmann)	Introduction to Vaccinology Assays and Immune Response
Session 3 (Gilbert)	Introduction to Frameworks for Assessing Surrogate Endpoints/Immunological Correlates of VE
Session 4 (Follmann)	Additional Study Designs for Evaluating VE
Session 5 (Gilbert)	Methods for Assessing Immunological Correlates of Risk and Optimal Surrogate Endpoints
Session 6 (Gilbert)	Effect Modifier Methods for Assessing Immunological Correlates of VE (Part I)
Session 7 (Gabriel)	Effect Modifier Methods for Assessing Immunological Correlates of VE (Part II)
Session 8 (Sachs)	Tutorial for the R Package <i>pseval</i> for Effect Modifier Methods for Assessing Immunological Correlates of VE
Session 9 (Gilbert)	Introduction to Sieve Analysis of Pathogen Sequences, for Assessing How VE Depends on Pathogen Genomics
Session 10 (Follmann)	Methods for VE and Sieve Analysis Accounting for Multiple Founders

Outline of Session 5

- ① Traditional CoR methods: Inverse probability weighted Cox model
- ② Key issues
 - Marker sampling design
 - Marker measurement error
- ③ Improved CoR methods (Breslow et al., 2009; Rose and van der Laan, 2011)
- ④ Estimated optimal surrogate (Price, Gilbert, van der Laan, 2017)

Prospective Cohort Study Sub-Sampling Design Nomenclature

- Terms used: case-cohort, nested case-control, 2-phase sampling
 - Case-cohort sampling originally meant taking a Bernoulli random sample of subjects at study entry for marker measurements (the “sub-cohort”), and also measuring the markers in all disease cases (Prentice, 1986, *Biometrika*)
 - Nested case-control sampling is Bernoulli or without replacement sampling done separately within disease cases and controls (retrospective sampling)
 - 2-phase sampling is the generalization of nested case-control sampling that samples within discrete levels of a covariate as well as within case and control strata (Breslow et al., 2009, *AJE, Stat Biosciences*)
 - Source of confusion: Some papers allow case-cohort to include retrospective sampling
- We restrict case-cohort to its original meaning

The Cox Model with a Sub-Sampling Design

- Cox proportional hazards model

$$\lambda(t|Z) = \lambda_0(t) \exp \left\{ \beta_0^T Z(t) \right\}$$

- $\lambda(t|Z)$ = conditional failure hazard given covariate history until time t
- β_0 = unknown vector-valued parameter
- $\lambda_0(t) = \lambda(t|0)$ = unspecified baseline hazard function
 - Z are “expensive” covariates only measured on failures and subjects in a random sub-sample
 - i.e., Z = immune response biomarkers, measured at fixed time τ post-randomization or at longitudinal visits

Notation and Set-Up (Matches Kulich and Lin, 2004, *JASA*)

- T = failure time (e.g., time to HIV infection diagnosis)
- C = censoring time
- $X = \min(T, C), \Delta = I(T \leq C)$
- $N(t) = I(X \leq t, \Delta = 1)$
- $Y(t) = I(X \geq t)$
- Cases are subjects with $\Delta = 1$
- Controls are subjects with $\Delta = 0$

Notation and Set-Up (Matches Kulich and Lin, 2004, *JASA*)

- Consider a prospective cohort of N subjects, who are stratified by a variable V with K categories
- $\epsilon =$ indicator of whether a subject is selected for measurement of immune responses Z (and they are measured)
 - $\alpha_k = Pr(\epsilon = 1|V = k)$, where $\alpha_k > 0$
- $(X_{ki}, \Delta_{ki}, Z_{ki}(t), 0 \leq t \leq \tau, V_{ki}, \epsilon_{ki} \equiv 1)$ observed for all marker subcohort subjects
- At least $(X_{ki}, \Delta_{ki} \equiv 1, Z_{ki}(X_{ki}))$ observed for all cases

Estimation of β_0

- With full data, β_0 may be estimated by the MPLE, defined as the root of the score function

$$U_F(\beta) = \sum_{i=1}^n \int_0^\tau \{Z_i(t) - \bar{Z}_F(t, \beta)\} dN_i(t), \quad (1)$$

where

$$\bar{Z}_F(t, \beta) = S_F^{(1)}(t, \beta) / S_F^{(0)}(t, \beta);$$

$$S_F^{(1)}(t, \beta) = n^{-1} \sum_{i=1}^n Z_i(t) \exp \{ \beta^T Z_i(t) \} Y_i(t)$$

$$S_F^{(0)}(t, \beta) = n^{-1} \sum_{i=1}^n \exp \{ \beta^T Z_i(t) \} Y_i(t)$$

Estimation of β_0

- Due to missing data (1) cannot be calculated under the sub-sampling design
- Most estimators are based on pseudoscores parallel to (1), with $\bar{Z}_F(t, \beta)$ replaced with an approximation $\bar{Z}_C(t, \beta)$

$$U_C(\beta) = \sum_{k=1}^K \sum_{i=1}^{n_k} \int_0^{\tau} \{Z_{ki}(t) - \bar{Z}_C(t, \beta)\} dN_{ki}(t)$$

- The double indices k, i reflect the stratification

- The marker sampled cohort at-risk average is defined as

$$\bar{Z}_C(t, \beta) \equiv S_C^{(1)}(t, \beta) / S_C^{(0)}(t, \beta),$$

where

$$S_C^{(1)}(t, \beta) = n^{-1} \sum_{k=1}^K \sum_{i=1}^{n_k} \rho_{ki}(t) Z_{ki}(t) \exp \left\{ \beta^T Z_{ki}(t) \right\} Y_{ki}(t)$$

$$S_C^{(0)}(t, \beta) = n^{-1} \sum_{k=1}^K \sum_{i=1}^{n_k} \rho_{ki}(t) \exp \left\{ \beta^T Z_{ki}(t) \right\} Y_{ki}(t)$$

Estimation of β_0

- $\rho_{ki}(t)$ is set to zero for subjects with incomplete data, eliminating them from the estimation
- Cases and subjects in the marker subcohort have $\rho_{ki}(t) > 0$
 - Usually $\rho_{ki}(t)$ is set as the **inverse estimated sampling probability** (Using the same idea as the weighted GEE methods of Robins, Rotnitzky, and Zhao, 1994, 1995)
- Different estimators are formed by different choices of weights $\rho_{ki}(t)$
- Two classes of estimators (case-cohort and 2-phase)

Example CoR Analysis: RV144 HIV-1 VE Trial

Haynes et al. (2012, *NEJM*) assessed in vaccine recipients the association of 6 immune response biomarkers measured at Week 26 with HIV-1 infection through 3.5 years

- **2-phase sampling design:** Measured Week 26 responses from all HIV-1 infected cases ($n = 41$) and from a stratified random sample of controls ($n = 205$ by gender \times # vaccinations \times per-protocol)

Immune Response Variable	Est. HR (95% CI)	2-Sided P-value
IgA Magnitude-Breadth to Env	1.58 (1.07–2.32)	0.02
Avidity to A244 Strain	0.90 (0.55–1.46)	0.66
ADCC to 92TH023 Strain	0.92 (0.62–1.37)	0.67
Neutralization M-B to Env	1.46 (0.87–2.47)	0.15
IgG to gp70-V1V2 Env	0.57 (0.37–0.90)	0.014
CD4 T cell Magn to 92TH023	1.17 (0.83–1.65)	0.37

Borgan et al. (2000, *Lifetime Data Analysis*) Cox model estimator II

Case-cohort Estimators (Called N-estimators in Kulich and Lin, 2004)

- The subcohort is considered a sample from all study subjects regardless of failure status
 - The whole covariate history $Z(t)$ is used for all subcohort subjects
 - For cases not in the subcohort, only $Z(T_i)$ (the covariate at the failure time) is used
- Prentice (1986, Biometrika): $\rho_i(t) = \epsilon_i/\alpha$ for $t < T_i$ and $\rho_i(T_i) = 1/\alpha$
- Self and Prentice (1988, Ann Stat): $\rho_i(t) = \epsilon_i/\alpha$ for all t

- General stratified N-estimator

- $\rho_{ki}(t) = \epsilon_i / \hat{\alpha}_k(t)$ for $t < T_{ki}$ and $\rho_{ki}(T_{ki}) = 1$
 - $\hat{\alpha}_k(t)$ is a possibly time-varying estimator of α_k
 - α_k is known by design, but nonetheless estimating α_k provides greater efficiency for estimating β_0 (Robins, Rotnitzky, Zhao, 1994)
 - A time-varying weight can be obtained by calculating the fraction of the sampled subjects among those at risk at a given time point (Barlow, 1994; Borgan et al., 2000, Estimator I)

Two-phase Sampling Estimators (Called D-estimators in Kulich and Lin, 2004)

- Weight cases by 1 throughout their entire at-risk period
- D-estimators treat cases and controls **completely separately**
 - α_k apply to controls only, so that α_k should be estimated using data only from controls
- Nested case-control estimators are the special case with one covariate sampling stratum $K = 1$

Two-phase Sampling D-estimators

- General D-estimator

$$\rho_{ki}(t) = \Delta_{ki} + (1 - \Delta_{ki})\epsilon_{ki}/\hat{\alpha}_k(t)$$

- Borgan et al. (2000, Estimator II) obtained by setting

$$\hat{\alpha}_k(t) = \frac{\sum_i^n \epsilon_{ki}(1 - \Delta_{ki})Y_{ki}(t)}{\sum_i^n (1 - \Delta_{ki})Y_{ki}(t)},$$

i.e., the proportion of the sampled controls among those who remain at risk at time t

- the `cch` package in R (by Thomas Lumley and Norm Breslow) implements the Cox model for case-cohort (N-estimators) and 2-phase sampling (D-estimators) (code for using `cch` to analyze a data set is provided at <http://faculty.washington.edu/peterg/SISMID2017.html>)

Main Distinctions Between N- and D- Estimators

- D-estimators require data on the complete covariate histories of cases
- N-estimators only require data at the failure time for cases
 - E.g., for the Vax004 HIV VE trial, the immune responses in cases were only measured at the visit prior to infection, so N-estimators are valid while D-estimators are not valid

Main Distinctions Between N- and D- Estimators

- For N-estimators, the sampling design is **specified in advance**, whereas for D-estimators, it can be **specified after the trial** (retrospectively)
 - D-estimators more flexible

Gaps of Both N- and D- Estimators

Estimator	Does Not Need Full Covariate Histories in Cases	Allows Outcome-Dependent Sampling
N (Prosp. case-cohort)	Yes	No
D (Retrospective 2-phase)	No	Yes

- For time-dependent correlates, none of the partial-likelihood based methods are flexible on both points
- All of the methods require full covariate histories in controls

- ① Traditional CoR methods: Inverse probability weighted Cox model
- ② Key issues
 - Marker sampling design
 - Marker measurement error
- ③ Improved CoR methods (Breslow et al., 2009; Rose and van der Laan, 2011)
- ④ Estimated optimal surrogate (Price, Gilbert, van der Laan, 2017)

Some Marker Sampling Questions to Consider Further

- Prospective or retrospective sampling?
- How much of the cohort to sample?
- Sampling design: Which subjects to sample?

Prospective case-cohort sampling: Select a random sample for immunogenicity measurement **at baseline**

- Advantages of prospective sampling
 - Can estimate case incidence for groups with certain immune responses
 - Can study correlations of immune response with multiple study endpoints
 - Straightforward to descriptively study the distribution of the immune responses in the whole study population at-risk when the immune responses are measured
 - **Practicality:** The lab will know what subjects to sample as early as possible, and there is one simple subcohort list

Retrospective 2-phase sampling: At or after the final analysis, select a random sample of control subjects for immunogenicity measurement

- Advantages of retrospective sampling
 - Can match controls to cases to obtain balance on important covariates
 - E.g., balanced sampling on a prognostic factor gains efficiency (balanced sampling = equal number of subjects sampled within each level of the prognostic factor for cases and controls)
 - Can flexibly adapt the sampling design in response to the results of the trial
 - E.g., Suppose the results indicate effect modification, with $VE \gg 0$ in a subgroup and $VE \approx 0\%$ in other subgroups. Could over-sample controls in the 'interesting' subgroup.

Prospective or Retrospective Sampling?

- For cases where there is one primary endpoint and it is not of major interest to estimate absolute case incidence, retrospective sampling may be typically referred

How Many Controls to Sample?

- In prevention trials, for which the clinical event rate is low, it is very expensive and unnecessary to sample all of the controls
 - Vax004 trial vaccine recipients: 225 HIV infected cases; ≈ 3000 controls
 - RV144 trial vaccine recipients: 41 HIV infected cases; ≈ 7000 controls
 - **Rule of thumb:** Under the null hypothesis, a $K : 1$ Control:Case ratio achieves relative efficiency of $1 - \frac{1}{1+K}$ compared to complete sampling

K	Relative Efficiency
1	0.50
2	0.67
3	0.75
4	0.80
5	0.83
10	0.91

- Simulations useful for studying the trade-offs of different K under alternative CoR hypotheses

Which Controls to Sample?

Two-Phase Sampling

- **Phase I:** All N trial participants are classified into K strata on the basis of information known for everyone: N_k in stratum k ;
$$N = \sum_{k=1}^K N_k$$
- **Phase II:** For each k , $n_k \leq N_k$ subjects are sampled at random, and the 'expensive' immune response biomarkers Z are measured for the resulting $n = \sum_{k=1}^K n_k$ subjects

Which Controls to Sample?

Principle: Well-powered CoR evaluation requires broad variability in the biomarker and in the risk of the clinical endpoint

- Can improve efficiency by over-sampling the “most informative” subjects
 - Disease cases (usually sampled at 100%)
 - Rare or unusual immune responses; or rare covariate patterns believed to affect immune response (e.g., HLA subgroups)
- Auxiliary Phase I variables measured in everyone are most valuable when they predict the missing data (i.e., the biomarker of interest)
- In general, optimal sampling obtained with sampling probabilities proportional to the cost-adjusted square-root variance of the efficient influence function (Gilbert, Yu, Rotnitzky, 2014, *Stat Med*)

- ① Traditional CoR methods: Inverse probability weighted Cox model
- ② Key issues
 - Marker sampling design
 - Marker measurement error
- ③ Improved CoR methods (Breslow et al., 2009; Rose and van der Laan, 2011)
- ④ Estimated optimal surrogate (Price, Gilbert, van der Laan, 2017)

Illustrative Example

- 'True' CoR $S^* \sim N(0, 1)$
- 'Measured CoR' $S = S^* + \epsilon, \epsilon \sim N(0, \sigma^2)$
- Infection status Y generated from $\Phi(\alpha + \beta S^*)$

with α set to give $P(Y = 1|S^* = 0) = 0.20$ and β set to give $P(Y = 1|S^* = 1) = 0.15$

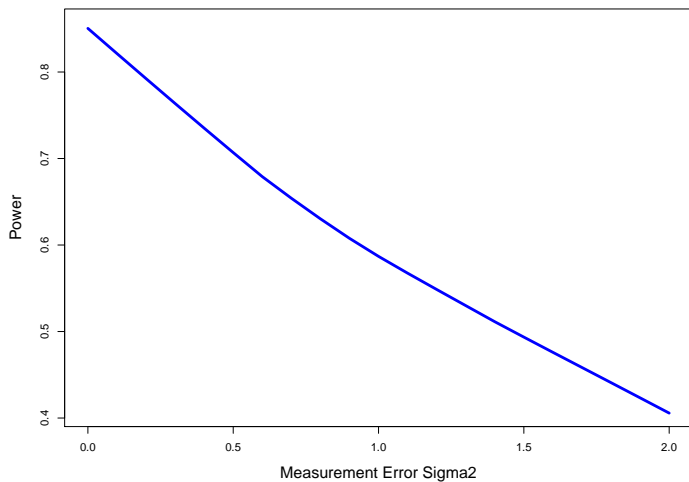
σ^2 ranges from 0 to 2 (no-to-large measurement error)

Simple Simulation Study

- Consider a study with $n = 500$ participants
- Consider power of a logistic regression model to detect an association between S and Y

Measurement Error Reduces Power to Detect a CoR

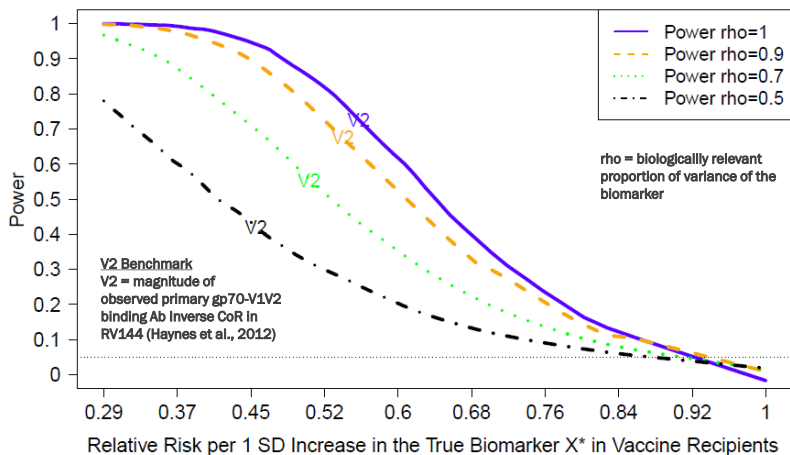
Deterioration of Power to Detect a CoR with Increasing Measurement Error



Power Calculations for Assessing CoRs

- Ideally, the power/sample size calculations should explicitly account for measurement error in the assay
 - E.g., Gilbert, Janes, Huang (2016, *Stat Med*), implemented in the R package *CoRpower* posted at <http://faculty.washington.edu/peterg/programs.html>
 - E.g., specify $\rho \equiv \sigma^2 / \sigma_{obs}^2$, the proportion of inter-vaccinee variability of the biomarker that is biologically relevant
 - **Rule of thumb:** ρ = relative efficiency for estimating a CoR odds ratio for the underlying perfect biomarker compared to the observed biomarker (McKeown-Eyssen, Tibshirani, 1994, *AJE*)
 - 'Noise' components of σ_{obs}^2 may be estimated, especially from laboratory assay validation studies
 - Within-vaccinee variability of replicates
 - Between-vaccinee variability due to variability in the time from the last immunization to marker sampling

Power to Detect a CoR of HIV Infection in Vaccinees in HVTN 505 ($\alpha = 0.05$)



Method: 2-phase logistic regression (Holubkov and Breslow, 1997)

- ① Traditional CoR methods: Inverse probability weighted Cox model
- ② Key issues
 - Marker sampling design
 - Marker measurement error
- ③ Improved CoR methods (Breslow et al., 2009; Rose and van der Laan, 2011)
- ④ Estimated optimal surrogate (Price, Gilbert, van der Laan, 2017)

Typical Correlates Assessments are Inefficient

- Broadly in epidemiology studies, biomarker-disease associations are commonly assessed ignoring much data collected in the study
- That is, only subjects with the biomarker measured are included in the analysis
- Standard analyses use inverse probability weighting of the biomarker sampled subcohort, including all of the methods discussed so far
- These ubiquitously-used methods are implemented in the R package `cch` (Breslow and Lumley)

Typical Correlates Assessments are Inefficient

- Breslow et al.* urge statisticians/epidemiologists to consider using the whole cohort in the analysis of case-cohort/2-phase sampling data
- Baseline data on demographics and potential confounders are typically collected in all subjects (the Phase I data measured in everyone)
- These Phase I data are most valuable when they predict “missing” data

*Breslow, Lumley et al. (2009, *AJE*, *Stat Biosciences*)

How to Leverage All of the Data?

- **Question:** How can we use the Phase I data to improve the assessment of CoRs?
- **One Answer:** One approach adjusts the sampling weights used in the standard analyses described above to obtain approximately efficient estimators (e.g., Breslow et al., 2009, *AJE, Stat Biosciences*)

Some Lessons Learned from Breslow et al. (2009)

- 1 Obtain 'worthwhile' efficiency gain for the CoR assessment if baseline covariates can explain at least 40% of the variation in the immunological biomarker ($R^2 \geq 0.40$)
- 2 If interested in interactions (evaluation of whether a baseline covariate measured in everyone modifies the association of the biomarker and the clinical endpoint), can obtain worthwhile efficiency gain with a lower R^2
- 3 Even if no gain for the CoR assessment, will usually dramatically improve efficiency for assessing the associations of the Phase I covariates with outcome
- 4 Therefore it may often be the preferred method, and all practitioners should have methods accounting for all of the data in their analytic toolkit
- 5 Additional research needed to make these more-efficient methods work well for multivariate markers and for time-dependent markers

How to Leverage All of the Data?

- **Question:** How can we use the Phase I data to improve the assessment of CoRs?
- **Another Answer:** Use an efficient and double-robust method: Inverse probability of censoring weighted targeted minimum loss based estimation (IPCW-TMLE) (Rose and Van der Laan, 2011, *Int J Biostat*)

Right-Censored Data Structure for Fixed Follow-up Time t

- $V =$ Phase I information: Covariates (Z, V_0) , $\tilde{T} = \min(T, C)$, $\Delta = I(T \leq C)$, $Y^* = I(\tilde{T} \leq t)\Delta$, Phase II sampling probability ϵ
- $S = (A, W) =$ Phase II information: Immune response biomarkers measured at τ
 - Focus on the marker A of interest; $W =$ all other markers
 - Repeat the analysis taking each element of W as A

- ① Traditional CoR methods: Inverse probability weighted Cox model
- ② Key issues
 - Marker sampling design
 - Marker measurement error
- ③ Improved CoR methods (Breslow et al., 2009; Rose and van der Laan, 2011)
- ④ Estimated optimal surrogate (Price, Gilbert, van der Laan, 2017)

Introduction to an Optimal Surrogate*

- **Goal:** Develop a most-promising surrogate outcome for a clinical outcome so that future randomized studies can restrict themselves to only collecting the surrogate outcome
- Data from a clinical trial for developing a surrogate: n iid observations of $O = (W, A, S, Y)$
 - W = Baseline covariates
 - A = Treatment assignment (1=vaccine, 0=placebo)
 - S = Response variables/markers measured by an intermediate time point τ
 - Y = Outcome of interest at a final time point τ_1 after τ
- Assume A is randomized conditional on W

*Price B, Gilbert PB, van der Laan MJ. Estimation of the Optimal Surrogate Based on a Randomized Trial. *Under Review*.

Optimal Surrogate = Valid Surrogate that Optimally Predicts Y

- Define an **optimal surrogate** for the current trial as the function of (W, A, S) that satisfies the Prentice definition and that optimally predicts Y
 - A true parameter that is estimated
- **Goal:** Use the estimated optimal surrogate in **future clinical trials** for estimation and testing of a mean contrast treatment effect on Y
 - Tackles the **transportability problem** of inferring the causal treatment effect in a new trial without measuring Y
 - (also addressed by Pearl and Bareinboim, 2011, 2012)

Optimal Surrogate Framework vs. Other Frameworks

- **vs. controlled/natural effects and VE curve frameworks:**
Departs by being based on average causal effects identified from standard assumptions in randomized trials
- **vs. Prentice/valid replacement endpoint framework:** Aligns in that the optimal surrogate satisfies the **Prentice definition**
 - Partially aligns with the **Prentice criteria**
 - The best optimal surrogate will have treatment and candidate surrogate highly predictive of Y , similar to Prentice criteria 1 and 2
 - The framework posits a conditional mean version of Prentice criterion 3 for licensing correct inferences on Y in a new trial
 - It handles equally well the general case where S varies or is constant in the placebo group
- **vs. meta-analysis framework:** Aligns in its objective of inference on the clinical treatment effect in a future study without collecting Y in that study (Gail et al., 2000, *Biostatistics*)
 - Departs in being based on a single (or few) trials and different transportability assumptions

Optimal Surrogate Framework

- Departs from all previous frameworks by defining the optimal surrogate as an unknown target parameter
 - Predicted values from the estimated optimal surrogate are used as the actual surrogate endpoint
 - In large samples this resulting surrogate must satisfy the Prentice definition (under the standard assumptions of an RCT)
- New approach in treating the surrogate endpoint problem as a supervised targeted learning problem
 - Previous methods evaluate a pre-selected univariable or low-dimensional vector candidate surrogate
 - the optimal surrogate approach is efficient in allowing all collected data to potentially contribute to the optimal surrogate, through unbiased machine learning
 - The optimal surrogate approach is robust in that consistent estimates of the clinical treatment effects in the current and future trials are obtained without parametric modeling assumptions

Introduction to an Optimal Surrogate

- This approach is about the **search for promising surrogates** based on an efficacy trial(s) with (W, A, S, Y) measured
- A promising surrogate is one that satisfies the Prentice definition and is optimally predictive of Y **in this original trial**
- A **best starting point** for building a surrogate that is promising for the ultimate objective of bridging/inference on the clinical treatment effect in new settings based on (W, A, S)

Statistical Formulation of an Optimal Surrogate

- W = baseline covariates
- A = binary treatment assigned at baseline
- S = vector of intermediate outcomes measured at time τ
- Y = final univariate outcome measured at time τ_1 after τ

- Potential outcomes (S_1, S_0) and (Y_1, Y_0) under treatment assignment $A = 1$ and $A = 0$
- Treatment A is randomized conditional on W

A Nonparametric Approach

- $X = (W, S_0, S_1, Y_0, Y_1)$ = full-data structure with distribution $P_{X,0}$
- $O = (W, A, S, Y)$ = observed data with distribution P_0 determined by $P_{X,0}$ and $g_0(a | X) = g_0(a | W)$
- The statistical model \mathcal{M} for P_0 makes at most some assumptions about g_0
 - Known in a randomized trial
- \mathcal{M} puts no assumptions on the marginal distribution of W nor on the conditional distribution of (S, Y) given A, W

Candidate Surrogate Outcomes

- Any real-valued function $(W, A, S) \rightarrow \psi(W, A, S) \in \mathbb{R}$ is a **candidate surrogate**, representing a measurement one can collect by time τ
- **Question:** How to define a good surrogate in terms of the true data distribution P_0 ?
- **Starting point:** Only consider $S^\psi \equiv \psi(W, A, S)$ that are valid in the actual study, according to the Prentice definition:

$$E_0(Y_1 - Y_0) = 0 \quad \text{if and only if} \quad E_0(S_1^\psi - S_0^\psi) = 0,$$

where $S_a^\psi = \psi(W, a, S_a)$, for $a \in \{0, 1\}$

- Guarantees that an α -level test for $H_0^\psi : E_0(S_1^\psi - S_0^\psi) = 0$ is also an α -level test for $H_0 : E_0(Y_1 - Y_0) = 0$

Optimal Surrogate Outcome

- Criterion for ranking valid surrogates and defining a P_0 -optimal surrogate: full-data mean squared error

$$\psi \rightarrow MSE_{P_{X,0}}(\psi) \equiv \sum_a E_{P_{X,0}} \{g_0(a | W)(Y_a - \psi(W, a, S_a))^2\}$$

- **Goal:** Minimize the weighted mean square prediction error for predicting Y_a across $a \in \{0, 1\}$ subject to the Prentice definition constraint
- Given a class Ψ of possible surrogate functions $\psi()$, the P_0 -optimal surrogate in this class is defined as

$$\psi_0^F = \arg \min_{\psi \in \Psi} MSE_{P_{X,0}}(\psi)$$

- We focus on the nonparametric class– all functions of (W, A, S)

Optimal Surrogate Outcome

The minimizer of $\psi \rightarrow MSE_{P_{X,0}}(\psi)$ over all functions $(W, A, S) \rightarrow \psi(W, A, S)$ that satisfy the Prentice definition is:

$$\bar{S}_0 = \psi_0(W, A, S) \equiv E_0(Y | W, A, S)$$

Potential outcomes of this P_0 -optimal surrogate: $\bar{S}_{0,a} = E_0(Y_a | W, S_a)$,
 $a \in \{0, 1\}$ and

$$E_{P_0}(\bar{S}_{0,a} | W) = E_{P_0}(Y_a | W)$$

- **Implications:**

- The surrogate treatment effect has the same interpretation as the clinical treatment effect
- Under P_0 , a 95% CI for the causal effect of treatment on the P_0 -optimal surrogate is also a 95% CI for the causal effect of treatment on Y

Conditions for a New Study P Under which the P_0 -Optimal Surrogate is also the P -Optimal Surrogate

Consider a new study with iid observations $O^* = (W^*, A^*, S^*, Y^*) \sim P$, where A^* is randomized conditional on W^*

Assumptions:

- **Equal Conditional Means:**

$$E[Y^* | W^* = w, A^* = a, S^* = s] = E[Y | W = w, A = a, S = s] \text{ for all } (w, a, s) \text{ in a support of } (W^*, A^*, S^*)$$

- **Contained Support:** A support of (W^*, A^*, S^*) is contained in a support of (W, A, S)
- **Positivity:** $P(A^* = a | W^*) > 0$ a.e. for $a \in \{0, 1\}$

Result: The P_0 -optimal surrogate equals the P -optimal surrogate: for all (w, a, s) in a support of (W^*, A^*, S^*)

$$\begin{aligned} E_P(Y^* | W^* = w, A^* = a, S^* = s) &= E_{P_0}(Y | W = w, A = a, S = s) \\ &= E_P(Y_a^* | W^* = w, S_a^* = s) = E_{P_0}(Y_a | W = w, S_a = s) \end{aligned}$$

Transportability Theorem Under a Prentice Criterion 3: Application to a New Treatment $A^* \neq A$

- If the new study considers a new treatment $A^* \neq A$, then generally the transportability theorem will not apply, because
$$E[Y^* | W^* = w, A^* = a, S^* = s] \neq E[Y | W = w, A = a, S = s]$$

Transportability Theorem Under a Prentice Criterion 3: Application to a New Treatment $A^* \neq A$

- Special case where the transportability assumptions may be reasonable
- **Same three assumptions as in Theorem 2**
- **Prentice criterion 3 assumption for both settings:**

$$E[Y^* | W^*, A^*, S^*] = E[Y^* | W^*, S^*]$$

$$E[Y | W, A, S] = E[Y | W, S]$$

Result: The P -optimal surrogate equals the P_0 -optimal surrogate and

$$E_{P_0}(Y_a | W = w, S_a = s) \text{ is constant in } a$$

$$E_P(Y_a^* | W^* = w, S_a^* = s) \text{ is constant in } a$$

Estimation of the P_0 -optimal Surrogate

- Estimation of the P_0 -optimal surrogate is a standard prediction problem
- Estimate $E_0(Y | W, A, S)$ by a minimizer of the risk of a loss
 - Use MSE loss (matched to the optimality criterion for defining the optimal surrogate)
- Loss-based super-learning*: yields an optimal estimator among any given class of candidate estimators
 - Oracle inequality for the cross-validation selector: the estimator is asymptotically at least as good as any candidate in the set of candidate estimators
 - $CV-R^2 \in [0, 1]$ provides a universal measure of the strength of the estimated optimal surrogate, allowing comparisons of different candidate surrogate estimators across studies and within a study

*Leo Breiman (1984); van der Laan, Polley, and Hubbard (2007); van der Laan and Rose (2011) textbook

Targeted Estimate of the Optimal Surrogate

- Let ψ_n be the super-learner estimator of

$$\psi_0(W, A, S) = E_0(Y | W, A, S)$$

- ψ_n may be updated to be a TMLE of ψ_0 , ψ_n^{TMLE}

TMLE = targeted minimum loss-based estimation (e.g., van der Laan and Rose, 2011)

The Targeted Estimated Optimal Surrogate Provides an Efficient Estimator of $\theta_0 = E(Y_1 - Y_0)$

Use $\psi_n^{TMLE}(W, A, S)$ in place of the final outcome Y

- Based on the reduced data $(W_i, A_i, \psi_n^{TMLE}(W_i, A_i, S_i))$, $i = 1, \dots, n$, compute the TMLE θ_n^{TMLE} of the data adaptive target parameter

$$\theta_{\psi_n} = E_0 \left[\psi_n^{TMLE}(W, 1, S_1) - \psi_n^{TMLE}(W, 0, S_0) \right]$$

The Targeted Estimated Optimal Surrogate Provides an Efficient Estimator of $\theta_0 = E(Y_1 - Y_0)$

Use $\psi_n^{TMLE}(W, A, S)$ in place of the final outcome Y

- Based on the reduced data $(W_i, A_i, \psi_n^{TMLE}(W_i, A_i, S_i))$, $i = 1, \dots, n$, compute the TMLE θ_n^{TMLE} of the data adaptive target parameter

$$\theta_{\psi_n} = E_0 \left[\psi_n^{TMLE}(W, 1, S_1) - \psi_n^{TMLE}(W, 0, S_0) \right]$$

- θ_n^{TMLE} is an asymptotically efficient estimator of θ_{ψ_n} based on the reduced data
- **It is also an asymptotically efficient estimator of θ_0 based on $O = (W, A, S, Y)$ in the statistical model \mathcal{M} !**

Inference on $\theta_0 = E(Y_1 - Y_0)$ Based on θ_n^{TMLE}

- θ_n^{TMLE} based on the reduced data model is asymptotically linear with influence curve equal to that of the TMLE $\tilde{\theta}_n^{TMLE}$ of $\theta_0 = E_0(Y_1 - Y_0)$ based on the data (W_i, A_i, Y_i) , $i = 1, \dots, n$
- Thus a Wald $(1 - \alpha)\%$ CI for θ_{ψ_n} based on θ_n^{TMLE} is also a $(1 - \alpha)\%$ CI for θ_0 and is as narrow as a CI based on an efficient estimator of θ_0 using (W, A, Y)
- **Conclusion: The optimal surrogate has the perfect properties for the original study**

Inference on $\theta_P^* = E_P(Y_1^* - Y_0^*)$ based on θ_n^{TMLE}

Now suppose we have built θ_n^{TMLE} based on $(W_i, A_i, S_i, Y_i) \sim P_0$ from an original efficacy trial(s) and a second trial is done with $(W_i^*, A_i^*, S_i^*, Y_i^*) \sim P$ only measuring (W_i^*, A_i^*, S_i^*)

- 1 Calculate the $\psi_n^{TMLE}(W_i^*, A_i^*, S_i^*)$ surrogate outcome values, $i = 1, \dots, n^*$
- 2 Estimate the treatment-specific surrogate means

$$\theta_{\psi_n}^a(P) = E_P \left[E_P(\psi_n^{TMLE}(W^*, a, S^*) \mid A^* = a, W^*) \right]$$

Estimate by $\theta_{\psi_n}^{TMLE,a}(P) = \frac{1}{n^*} \sum_{i=1}^{n^*} \psi_n^{TMLE}(W_i^*, a, S_i^*)$, $a = 0, 1$

- 3 Estimate $\theta_{\psi_n}(P) = \theta_{\psi_n}^1(P) - \theta_{\psi_n}^0(P)$
- 4 Compute Wald-based CIs for $\theta_{\psi_n}^1(P)$, $\theta_{\psi_n}^0(P)$, $\theta_{\psi_n}(P)$ based on the influence functions

Inference on $\theta_P^* = E_P(Y_1^* - Y_0^*)$ based on θ_n^{TMLE}

Under Theorem 2, $\theta_{\psi_n}^{TMLE,1}(P)$, $\theta_{\psi_n}^{TMLE,0}(P)$, $\theta_{\psi_n}^{TMLE}(P)$ are consistent estimators of $E_P(Y_1^*)$, $E_P(Y_0^*)$, $\theta_P^* = E_P(Y_1^* - Y_0^*)$

- The CI for $\theta_{\psi_n}(P)$ is correct for θ_P^* for an infinite sample sized original P_0 -study $n = \infty$

Inference on $\theta_P^* = E_P(Y_1^* - Y_0^*)$ based on θ_n^{TMLE}

Under Theorem 2, $\theta_{\psi_n}^{TMLE,1}(P)$, $\theta_{\psi_n}^{TMLE,0}(P)$, $\theta_{\psi_n}^{TMLE}(P)$ are consistent estimators of $E_P(Y_1^*)$, $E_P(Y_0^*)$, $\theta_P^* = E_P(Y_1^* - Y_0^*)$

- The CI for $\theta_{\psi_n}(P)$ is correct for θ_P^* for an infinite sample sized original P_0 -study $n = \infty$
- Future work is needed to obtain correct CIs for θ_P^* for finite n
- This problem is readily solved if the surrogate means $\theta_{\psi_n}^a(P)$ were estimated using a parametric model
- However, obtaining a CI when estimating $\theta_{\psi_n}^a(P)$ nonparametrically through super-learning is much harder, because $\psi_0 = E(Y|W, A, S)$ is a function that is not estimable at root- n rate
 - E.g., the nonparametric bootstrap theoretically fails

Dengue Phase 3 Trial Example

- Two randomized, double-blinded, placebo-controlled, multicenter, Phase 3 trials of a recombinant, live, attenuated, tetravalent (4 serotypes) dengue vaccine (CYD-TDV)
 - **CYD14:** Asia-Pacific region (Capeding, et al., 2014, *The Lancet*)
 - **CYD15:** Latin America (Villar et al, 2015, *NEJM*)

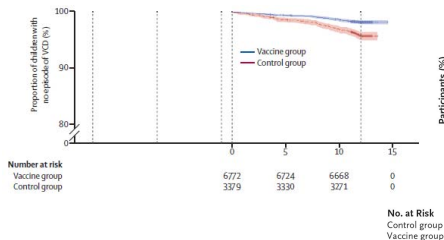
- 2:1 randomization to vaccine:placebo
- Immunizations at months 0, 6, 12
- Primary follow-up from Month 13 to Month 25 (active phase of follow-up)
- Primary endpoint: Symptomatic, virologically confirmed dengue (VCD)

Results on Vaccine Efficacy (Proportional Hazards Model)

CYD14: $\widehat{VE} = 56.5\%$ (95% CI 43.8–66.4)

CYD15: $\widehat{VE} = 64.7\%$ (95% CI 58.7–69.8)

CYD14 Trial (Capeding et al., 2014, *The Lancet*)



CYD15 Trial (Villar et al., 2015, *NEJM*)

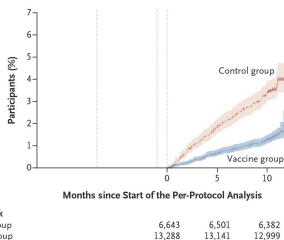


Illustration of Estimated Optimal Surrogate Approach

Analysis carried out by Brenda Price

- Based on pseudo CYD14 and CDY15 simulated data sets
- Treat CYD14 as the current trial; CYD15 as the future trial

- A = Vaccination status (1=vaccine; 0=placebo)
- Y = Disease outcome (1=VCD endpoint between Month 13 and 25; 0 = no VCD endpoint by Month 25)
- W = Baseline covariates: age, sex, baseline PRNT₅₀ neutralization titers to the 4 vaccine strains (serotypes 1–4)
- S = Month 13 PRNT₅₀ neutralization titers to the 4 vaccine strains (serotypes 1–4)

Illustration: Inference on VE_0 in CYD14

- 1 Obtain $\psi_n^{TMLE}(W, A, S)$ from the CYD14 data (W_i, A_i, S_i, Y_i) , $i = 1, \dots, n$, yielding estimates of

$$E_0(\psi_n^{TMLE}(W, 1, S_1)), \quad E_0(\psi_n^{TMLE}(W, 0, S_0)),$$
$$VE_0(\psi_n^{TMLE}) = 1 - \frac{E_0(\psi_n^{TMLE}(W, 1, S_1))}{E_0(\psi_n^{TMLE}(W, 0, S_0))}$$

- 2 Compute Wald-based CIs for the above parameters based on the influence functions
- 3 Compare these point and interval estimates to direct estimates of $E_0(Y_1)$, $E_0(Y_0)$, and VE_0 based on (W_i, A_i, Y_i) from CYD14

Illustration: Inference on VE_P^* in CYD15 Based on the Surrogate Built from CYD14

- 1 Calculate the $\psi_n^{TMLE}(W_i^*, A_i^*, S_i^*)$ surrogate values for CYD15 participants, $i = 1, \dots, n^*$
- 2 Estimate the surrogate mean parameters in CYD15

$$\theta_{\psi_n}^a(P) = E_P \left[E_P(\psi_n^{TMLE}(W^*, a, S^*) \mid A^* = a, W^*) \right]$$

Illustration: Inference on VE_P^* in CYD15 Based on the Surrogate Built from CYD14

- 1 Calculate the $\psi_n^{TMLE}(W_i^*, A_i^*, S_i^*)$ surrogate values for CYD15 participants, $i = 1, \dots, n^*$
- 2 Estimate the surrogate mean parameters in CYD15

$$\theta_{\psi_n}^a(P) = E_P \left[E_P(\psi_n^{TMLE}(W^*, a, S^*) \mid A^* = a, W^*) \right]$$
$$\theta_{\psi_n}^{TMLE,a}(P) = \frac{1}{n^*} \sum_{i=1}^{n^*} \psi_n^{TMLE}(W_i^*, a, S_i^*), \quad a = 0, 1$$
$$\theta_{\psi_n}^{TMLE}(P) = VE_{\psi_n}^{TMLE}(P) = 1 - \frac{\theta_{\psi_n}^{TMLE,1}(P)}{\theta_{\psi_n}^{TMLE,0}(P)}$$

- 3 Compute Wald-based CIs for the above parameters
- 4 Compare these point and interval estimates to direct estimates of $E_P(Y_1^*)$, $E_P(Y_0^*)$, and VE_P^* based on (W_i^*, A_i^*, Y_i^*) from CYD15

- Use the MSE loss function for the super-learner cross-validation selector (matched to the optimality criterion for a surrogate)

Table: Input Variables for the Learning Algorithms

Input Variables	
<i>W</i> :	Baseline demographics age (range 2–14 years), sex, Baseline neutralization titers to the 4 vaccine strains, average, min, max of the 4 titers, interactions with age
<i>S</i> :	Month 13 neutralization titers to the 4 vaccine strains, average, min, max of the 4 titers, interactions with age

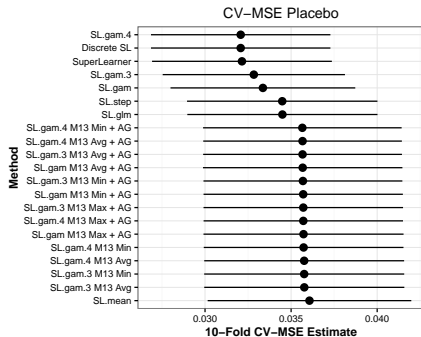
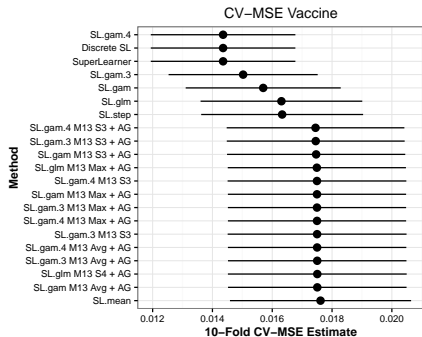
Super-learner to Estimate the Optimal Surrogate

Run super-learner separately for each treatment group $a \in \{0, 1\}$

Table: Learning Algorithms Employed

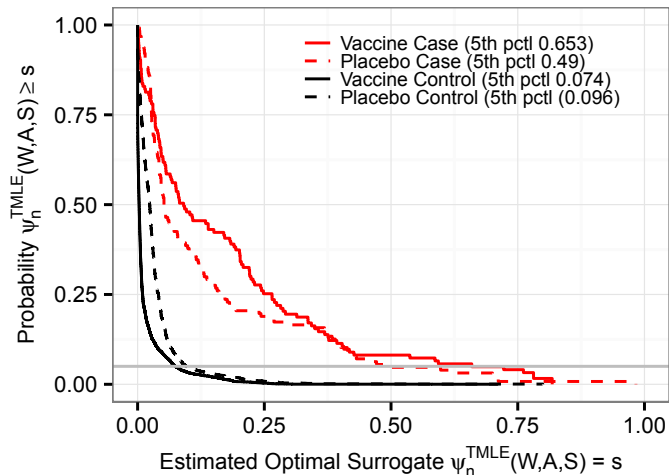
SL.mean	$E_0(Y W, A = a, S)^* = \beta_a$ for $a \in \{0, 1\}$
SL.glm	Logistic regression with all input variables
SL.step	Best logistic regression model by AIC through a step-wise search
SL.gam	gam for W & S inputs; all titer variables each with 2 df
SL.gam.3	gam for W & S inputs; all titer variables each with 3 df
SL.gam.4	gam for W & S inputs; all titer variables each with 4 df
M13 Sk	SL.glm, SL.gam, SL.gam.3, SL.gam.4 with only Month 13 serotype k titers
M13 Avg	SL.glm, SL.gam, SL.gam.3, SL.gam.4 with only Month 13 average titers
M13 Min, Max	SL.glm, SL.gam, SL.gam.3, SL.gam.4 with only Month 13 Min or Max titers
M13 Sk + AG	SL.glm, SL.gam, SL.gam.3, SL.gam.4 with Month 13 serotype k titers + (age, gender)
M13 Avg + AG	SL.glm, SL.gam, SL.gam.3, SL.gam.4 with Month 13 average titers + (age, gender)
M13 Min, Max + AG	SL.glm, SL.gam, SL.gam.3, SL.gam.4 with Month 13 Min or Max titers + (age, gender)
Discrete SL	van der Laan, Polley, and Hubbard (2007)
Super Learner (SL)	van der Laan, Polley, and Hubbard (2007)

Cross-Validated Mean-Squared Errors (CV-MSEs): CYD14



Empirical RCDFs for the Estimated Optimal Surrogate Values: CYD14

CYD14 Reverse CDFs



Estimated Optimal Surrogate (EOS) TMLEs of Target Parameters: CYD14

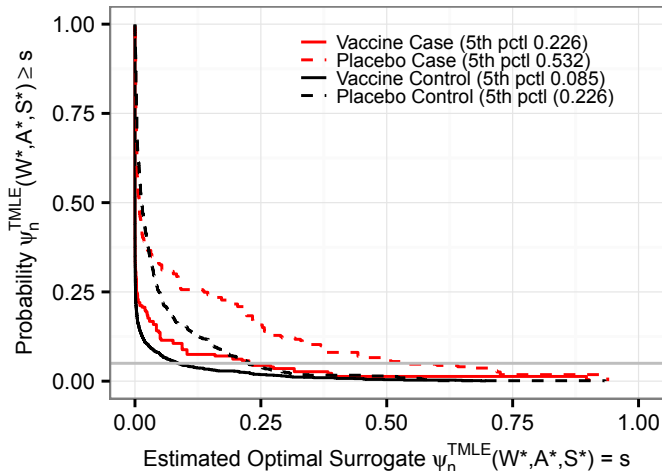
Parameter	TMLE Based on EOS	TMLE Based on (W, A, Y)
$E_0(Y_1)$	0.018 (0.014–0.023)	0.017 (0.014–0.019)
$E_0(Y_0)$	0.037 (0.023–0.060)	0.040 (0.036–0.045)
$VE_0 = 1 - \frac{E_0(Y_1)}{E_0(Y_0)}$	52% (41–66)	59% (54–65)

- The point estimate results have to be similar by construction!

Using the Estimated Optimal Surrogate (EOS) in CYD15

How well do the EOS values $\psi_n^{TMLE}(W^*, A^*, S^*)$ predict Y^* in CYD15?

CYD15 Reverse CDFs



How Well Does the Surrogate-Based Estimator Estimate VE_P^* in CYD15?

Table: Estimation in CYD15 based on the EOS built in CYD14 (not using outcome data Y^* in CYD15) vs. TMLE estimation using (W^*, A^*, Y^*) in CYD15

	TMLEs of Surrogate Parameters ¹		TMLEs of Clinical Parameters ²	
$\theta_{\psi_n}^1(P)$	0.017 (0.014–0.020)	$E_P(Y_1^*)$	0.017 (0.014–0.019)	¹ Based
$\theta_{\psi_n}^0(P)$	0.053 (0.040–0.069)	$E_P(Y_0^*)$	0.040 (0.036–0.045)	
$VE_{\psi_n}(P)$	68% (58–81)	VE_P^*	$VE_P^* = 61\%$ (54–67)	

on $(W_i, A_i, \theta_n^{TMLE}(W_i, A_i, S_i))$ and (W_i^*, A_i^*, S_i^*)

²Based on (W_i^*, A_i^*, Y_i^*) [use the actual clinical data]

Compare Predictive Ability of Input Variable Sets

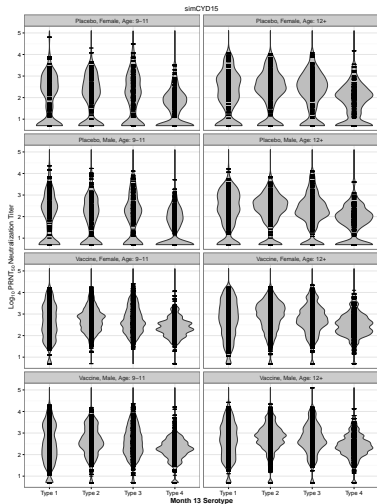
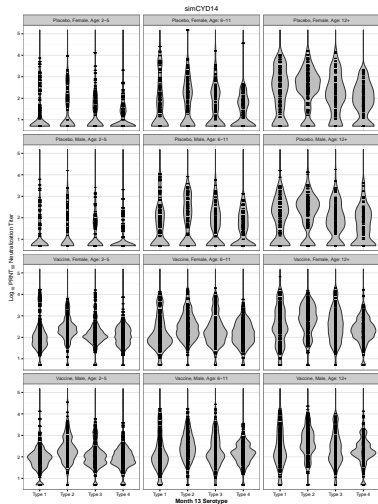
Table: Cross Validated AUCs* with 95% CIs

Input Set	CYD14 Vaccine	CYD14 Placebo	CYD15 Vaccine	CYD15 Placebo
(1) Demographics	0.61 (0.57, 0.66)	0.6 (0.55, 0.65)	0.54 (0.5, 0.58)	0.5 (0.47, 0.54)
(2) All baseline	0.89 (0.86, 0.92)	0.79 (0.76, 0.83)	0.58 (0.54, 0.61)	0.55 (0.51, 0.58)
(3) Month 13 titers	0.71 (0.67, 0.75)	0.63 (0.58, 0.68)	0.65 (0.62, 0.69)	0.57 (0.54, 0.61)
(4) All data	0.89 (0.86, 0.91)	0.76 (0.72, 0.8)	0.78 (0.76, 0.8)	0.6 (0.57, 0.64)

*Cross-validated area under the ROC-curves (Van der Laan, Hubbard, and Pajouh, 2013)

- The user can judge the tradeoff of **accuracy** and **simplicity** of the estimated optimal surrogate

Distributions of Month 13 Titers within (W, A) Strata



Checking Assumptions of the Transportability Theorem for Randomized Trials

- ① $E[Y^*|W^* = w, A^* = a, S^* = s] = E[Y|W = w, A = a, S = s]$ for all (w, a, s) in a support of (W^*, A^*, S^*)
 - ② A support of (W^*, A^*, S^*) is contained in a support of (W, A, S)
 - ③ Positivity: $P_0(A = a|W) > 0$ and $P(A^* = a|W^*) > 0$ a.e. for $a \in \{0, 1\}$
- *Condition 1* Examine by comparing estimates of $E[Y^*|W^* = w, A^* = a, S^* = s] = E[Y|W = w, A = a, S = s]$
 - *Condition 2* Examine by comparing distributions of (W, A, S) and (W^*, A^*, S^*)
 - *Condition 3* Examine by comparing distributions of W and of W^* between treatment groups

Two Simulation Studies

- **Objective of First Study:** Simple illustration that the estimated optimal surrogate will always provide unbiased estimation of $\theta_0 = E_0(Y_1 - Y_0)$ in the original trial, for any distribution of (W, A, S, Y)
- **Objective of Second Study:** Illustrate how well the estimated optimal surrogate built from one trial works for inference on $\theta_P^* = E_P(Y_1^* - Y_0^*)$ in a second trial, when Equal Conditional Means fails

Data Generating Distribution

- 10 candidate surrogates S^i , each taking values 0, 1, 2
- For each S^i :

$$P(S_1^i = 0, S_0^i = 0) = P(S_1^i = 1, S_0^i = 1) = P(S_1^i = 2, S_0^i = 2) = 0.1$$
$$P(S_1^i = 1, S_0^i = 0) = 0.5, P(S_1^i = 1, S_0^i = 2) = 0.2$$

$$Y = \sum_{i=1}^3 [0.1 * i * I(S^i = 1) + I(S^i = 2)] + \epsilon_Y, \quad \epsilon_Y \sim N(0, 0.1^2)$$

$$\theta_0 = E_0(Y_1 - Y_0) = -0.18 \quad \text{and} \quad E_0(S_1^i - S_0^i) = 0.3$$

(surrogate paradox occurs)

Comparator: Proportion of Treatment Explained Type Method

- For each S^i , estimate the Proportion of the Treatment Effect Captured (PCS)* based on a linear model

$$E[Y|S^i = s, A = a] = \beta_0 + \beta_1 * I[s = 1] + \beta_2 * I[s = 2]$$

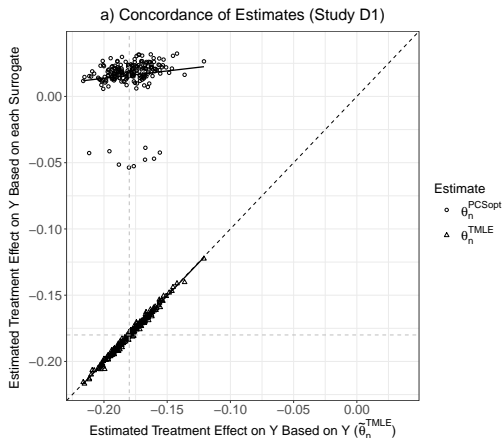
(true PCS = 0.87, 0.2, 0.002 for $i = 1, 2, 3$; PCS = 0 for $i = 4, \dots, 10$)

- Select the “best surrogate” i , $S^{\text{PCSopt}} = S^i$, as the one most frequently with greatest $\widehat{\text{PCS}}$ over 100 bootstrap data sets
- Estimate θ_0 by the difference ($a = 1$ minus $a = 0$) in average predicted Y 's in the fitted model

$$\hat{E}[Y|S^{\text{PCSopt}} = s, A = a] = \hat{\beta}_0 + \hat{\beta}_1 * I[s = 1] + \hat{\beta}_2 * I[s = 2]$$

*Kobayashi and Kuroki (2014, *Stat Med*)

Simulation 1 Results: $n = 2000$ subjects



- Surrogate paradox: $\theta_n^{\text{PCSopt}} > 0$ (vs. $\theta_0 = -0.18$)
 - Occurs in 96% of 200 generated data sets for the PCS approach (0% for SL-TMLE)

Simulation 2: Bridging to a Second Trial

- Simulate pairs of data sets (D1, D2) for the original and second trial
 - Original trial (As in Simulation 1):

$$Y = \sum_{i=1}^3 [0.1 * i * I(S^i = 1) + I(S^i = 2)] + \epsilon_Y, \quad \epsilon_Y \sim N(0, 0.1^2)$$

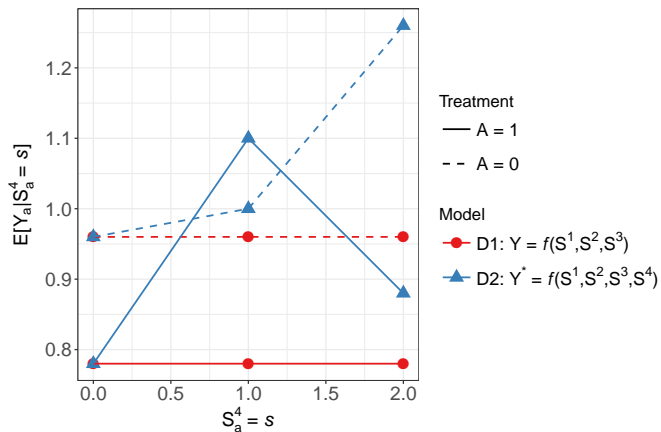
- New trial:

$$Y^* = \sum_{i=1}^4 [0.1 * i * I(S^{*i} = 1) + I(S^{*i} = 2)] + \epsilon_{Y^*}, \quad \epsilon_{Y^*} \sim N(0, 0.1^2)$$

- Equal Conditional Means fails because Y depends on (S^1, S^2, S^3) and Y^* depends on $(S^{*1}, S^{*2}, S^{*3}, S^{*4})$

Equal Conditional Means Fails

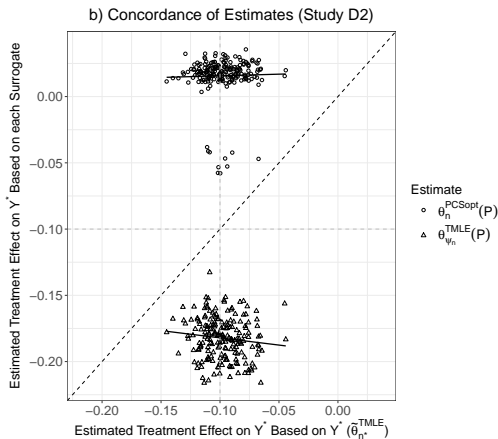
Violation of the Equal Conditional Means Assumption:
Differences in Y_a/Y_a^* by S_a^4



- The optimal surrogates $\psi_n^{TMLE}(A, S)$ and S_n^{PCSopt} are estimated from D1 as in Simulation 1
- Based on the (A^*, S^*) values in the paired data set D2, obtain surrogate-based estimates of $\theta_P^* = E_P(Y_1^* - Y_0^*)$
 - TMLE: $\theta_{\psi_n}^{TMLE}(P)$ as above
 - PCS:

$$\begin{aligned}\theta_n^{PCSopt}(P) &= \frac{1}{n_1^*} \sum_{i=1}^{n^*} A_i^* \widehat{E}[Y | S_i^{*PCSopt}, A_i^* = 1] \\ &\quad - \frac{1}{n_0^*} \sum_{i=1}^{n^*} (1 - A_i^*) \widehat{E}[Y | S_i^{*PCSopt}, A_i^* = 0]\end{aligned}$$

Simulation 2 Results: $n^* = 2000$



- Surrogate paradox: $= \theta_n^{\text{PCSopt}}(P) > 0$ (vs. $\theta_P^* = -0.18$)
 - Occurs in 95% of 200 generated data sets for the PCS approach (0% for SL-TMLE)

Conclusion from Simulation 2

- Demonstrates that the Equal Conditional Means assumption is necessary for valid inference of θ_P^* in a new setting
- When Equal Conditional Means is majorly violated, the estimated optimal surrogate can still preserve some accuracy in bridging the clinical treatment effect to a new setting

Start at the Right Place

- VanderWeele (2013, *Biometrics*) and discussants Joffe (2013) and Pearl (2013) suggest that a minimal requirement for an intermediate endpoint to be a useful surrogate endpoint is that it avoids the surrogate paradox
- VanderWeele (2013) shows that commonly used methods for surrogate endpoint evaluation generally do not guarantee avoiding this paradox
- The optimal surrogate approach starts at this minimal requirement, defining the optimal surrogate in a way guaranteed to satisfy the Prentice definition of a valid surrogate
 - Responds to Pearl's (2013) question:
"If we take the negation of the "surrogate paradox" as a criterion for "good" surrogate, why cannot we create a new, formal definition of "surrogacy" that will automatically avoid the paradox?..."

Nonparametric Supervised Learning Approach

- Using super-learner + TMLE seeks to minimize assumptions and use all of the information in the data
- Main application is when many candidate surrogates are measured, and the objective is supervised learning of most promising surrogate endpoints that may depend on baseline covariates as well as intermediate response endpoints
- This framework also applies for generating promising candidate surrogates based on observational studies, with all of the results holding under the additional assumption that all confounders W of treatment assignment are measured and included in the super-learner

Elaborations

- Missing data on $O = (W, A, S, Y)$
 - E.g., case-cohort or nested case-control sampling of S
 - Happenstance missingness
- Some participants experience Y before S is measured at τ
- Right-censoring of Y (failure time endpoint), competing risks outcomes
- Tailoring the super-learner to contextual features [sample size, event rate, dimensionality of (W, S)]
- Confidence intervals about the clinical treatment effect
 $\theta_P^* = E(Y_1^* - Y_0^*)$ in a new setting accounting for the error in estimating the optimal surrogate

Acknowledgements

- NIH NIAID support for the grant “Statistical Methods in HIV Vaccine Efficacy Trials”
- Participants and study personnel of the CYD14 and CYD15 dengue Phase 3 trials and SanofiPasteur colleagues for collaboration and sharing the data