

**Sieve Analysis: Statistical methods for
assessing differential vaccine protection
against HIV types**

Biostat 578A Lecture 7

Research Goal: Develop statistical methods for assessing from vaccine efficacy trial data how vaccine protection may depend on characteristics of the various circulating HIV-1 strains

Outline

I. Introduction to Sieve Analysis

II. Models for Sieve Analysis, Binary Endpoint
(HIV infection, Yes or No)

A. Discrete HIV types

B. Continuous HIV distance

III. Models for Sieve Analysis, Failure Time
Endpoint (Time to HIV infection Diagnosis)

A. Discrete HIV types

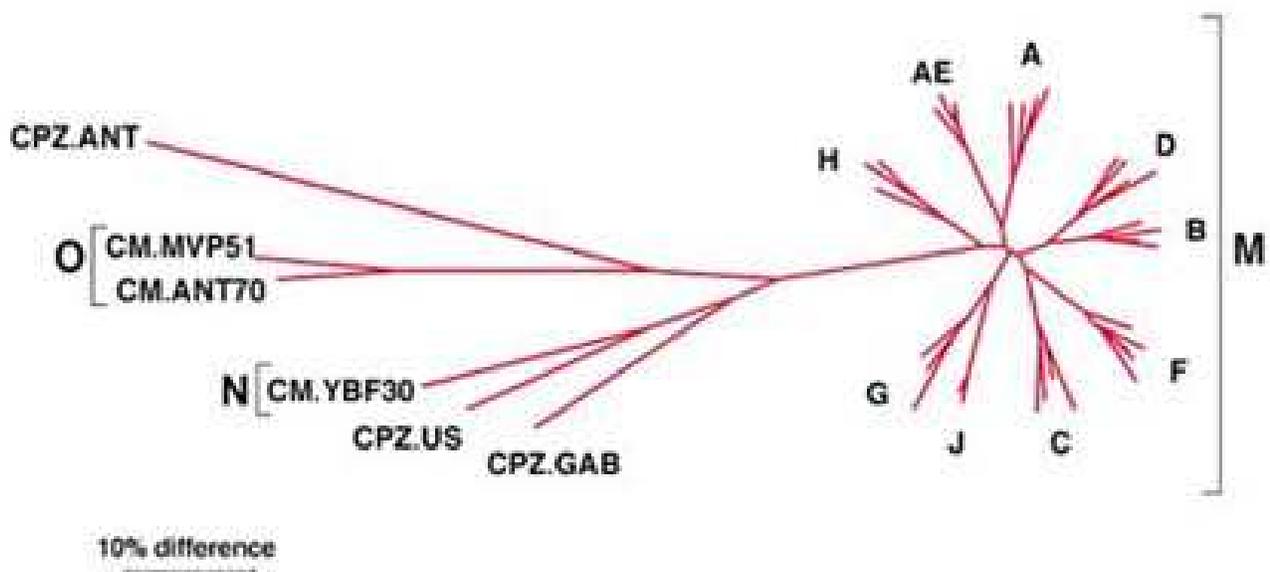
B. Continuous HIV distance: Lecture 8

Introduction to “Sieve” Analysis

- HIV-1 extremely diverse
- How broadly does a candidate vaccine protect?
- Vaccine protection depends on which characteristics of challenge virus? How so?

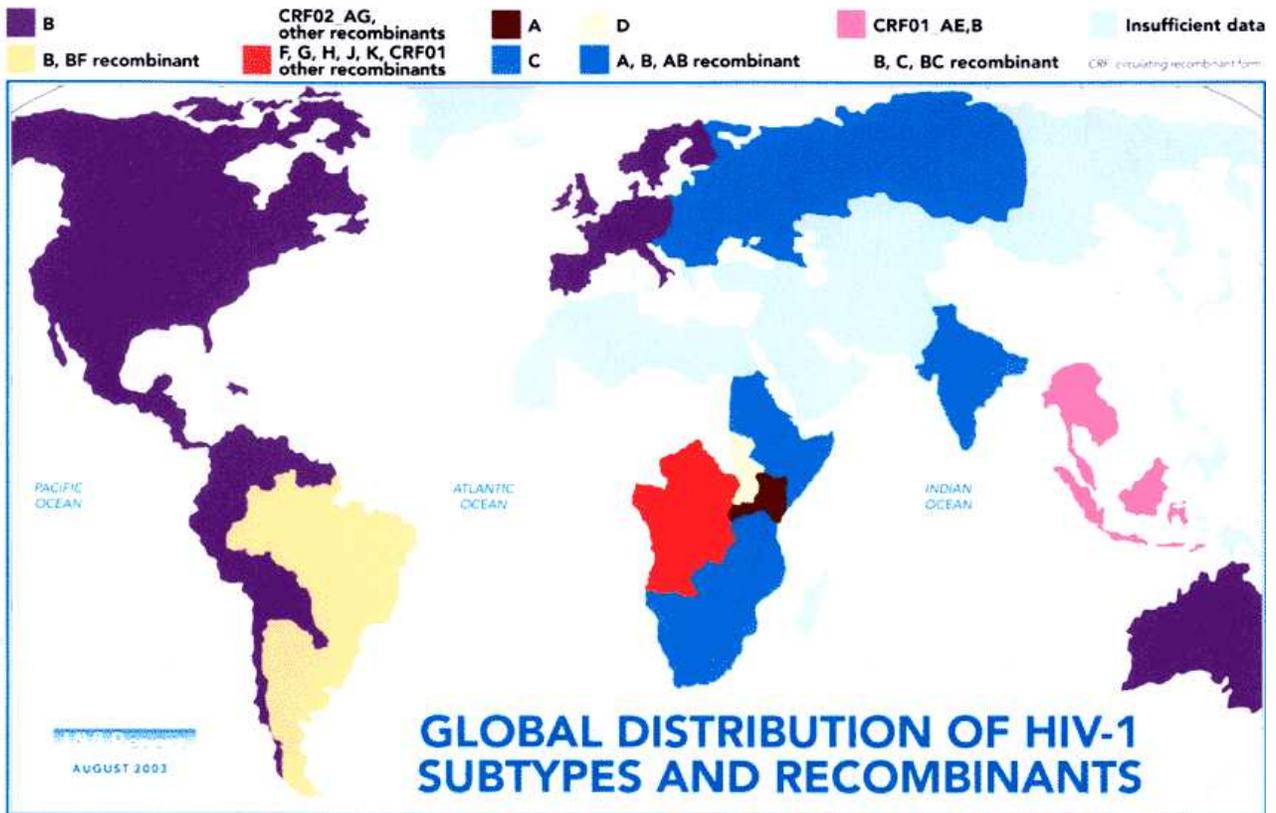
Phylogenetic Tree of HIV-1 Subtypes

Genetic Subtypes of HIV-1

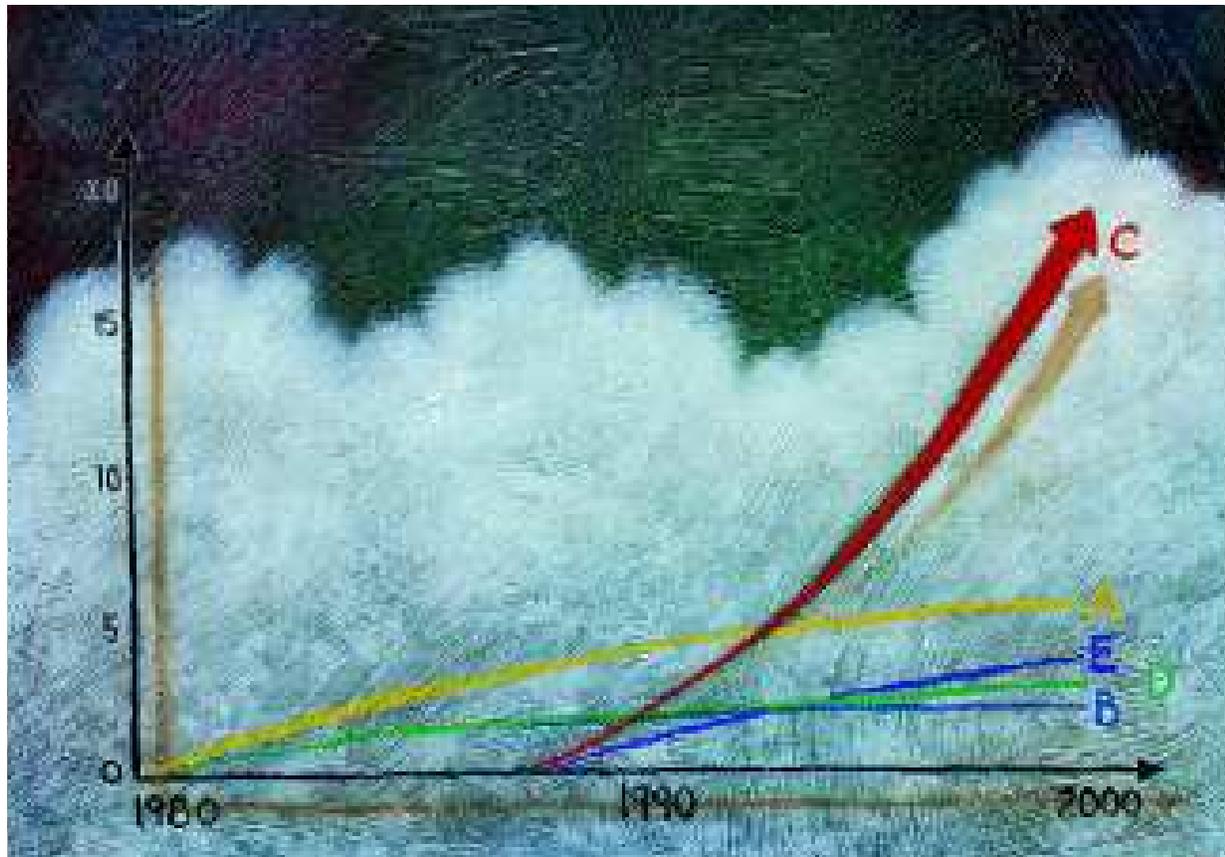


Source: Korber, B et al. Los Alamos National Laboratory

2003 Global Map of HIV-1 Subtypes



Global Dist'n of HIV-1 Subtypes 1980-1999



Introduction

- Human trials of preventive vaccines against heterogeneous pathogens
 - **hepatitis** Szmuness et al. 1981
 - **cholera** Clemens et al. 1991
van Loon et al. 1993
 - **rotavirus** Lanata et al. 1989
Ward et al. 1992
Ukae et al. 1994
Jin et al. 1996
Rennels et al. 1996
 - **pneumococcus** Amman et al. 1977
Smit et al. 1977
John et al. 1984
 - **influenza** Govaert 1994
 - **malaria** Alonso et al. 1994
- Some of these data summarized in Gilbert et al. (2001, J Clin Epidem)

Introduction

- Often no quantitative statistical assessment of type-specific vaccine efficacy
- When there is, the interpretation and validity of the analysis is often unclear

Data

- **Randomized vaccine trial**
- **Data collection**
 - Measure virus characteristics of isolated virus from breakthrough infections
 - E.g., VaxGen trials obtained 3 sequences from each infected participant, from a blood sample drawn at infection diagnosis

HIV Sequence Data

Sliding window for
analyzing 9-mers

V3 loop amino acid sequence
of reference GNE8 strain

...TRPNNNTRRSIHIG-PGR-AFYATGEIIGDIRQ...

Vaccine group V3 loop sequences

1. ...TRPNNNTRRRIHLG-PGR-AFYATG-IIGDIRQ...

2. ...TRPNNNTRKGIHIG-PGR-AFYATGEIIGNIRQ...

.

.

.

217. ...TRPSNNTRKGIHIG-PGR-AFYATEEITGDIRQ...

Placebo group V3 loop sequences

1. ...TRPNNNTRTGVLHG-PGR-VWYATGDIIGDIRQ...

2. ...TRPNNNTRRSIHIQ-PGR-AFYAT-DIIGDIRK...

.

.

.

119. ...TRPNNNTISKIRIR-PGRGSFYATNNIIGDIRQ...

Viral Variation Structure

0 = vaccine prototype strain

1. Nominal categorical:

K+1 distinct strains in circulation

$0, \dots, K$

2. Ordered categorical:

K+1 distinct strains in circulation

$0, \dots, K$

- ordered from prototype strain 0

3. Continuous:

Each strain is a continuous distance from prototype strain 0

- A vast number of meaningful ways to structure pathogen variation

The Problem with Sieve Analysis

Each viral isolate is genetically unique (under close examination) so that $2 \times K + 1$ table is too sparse and unstructured for meaningful analysis.

Solution: Add structure to the table

- a) Categorize infecting strains into nominal groups putatively related to strain-specific VE

E.g., categorize by

- subtype/clade
- phenotype (e.g., R5 vs X4)

- b) Order infecting strains by putative correlate of strain-specific VE

E.g., Order by some measure of similarity to strain used in vaccine construction

- nucleotide or amino acid sequence
- protein character

- c) Multidimensional viral feature

Categorical Model for Sieve Analysis

- **Counts data**

	Infecting strains				
	1	2	K
Placebo					
Vaccine					

- Assume $K + 1$ viral strains in circulation
- Nominal or ordinal response

Multinomial Logistic Regression Model (Cox, 1970; Anderson, 1972)

$$\Pr(Y = s|x) = \frac{\exp\{\alpha_s + \beta_s^T x\}}{1 + \sum_{l=1}^K \exp\{\alpha_l + \beta_l^T x\}}$$

$$s = 0, \dots, K; \quad \alpha_0 = \beta_0 \equiv 0$$

- Generalized linear logit model

$$\log \left\{ \frac{\Pr(Y = s|x)}{\Pr(Y = 0|x)} \right\} = \alpha_s + \beta_s^T x$$

- Interpretation of regression parameters:

$$\begin{aligned} \beta_s &= \log \left\{ \frac{\Pr(Y=s|\text{vacc})}{\Pr(Y=0|\text{vacc})} / \frac{\Pr(Y=s|\text{plac})}{\Pr(Y=0|\text{plac})} \right\} \\ &= \log \{OR(s)\} \end{aligned}$$

Model Properties

- Minimal assumptions
- Estimation by maximum likelihood
- Exact methods an option

Hirji (1992, JASA, **87**)

*Computing Exact Distributions for
Polytomous Response Data*

Strain-Specific Vaccine Efficacy

- Define “per strain-specific contact” vaccine efficacy by $VE^{pc}(s) = 1 - RR^{pc}(s)$

where

$$RR^{pc}(s) = \frac{\Pr(\text{Inf}|\text{Expos. to strain } s, \text{Vaccine})}{\Pr(\text{Inf}|\text{Expos. to strain } s, \text{Placebo})}$$

- $RR^{pc}(s)$ has an interpretation in terms of biological vaccine efficacy

Prospective Interpretation of Regression Model Parameters

Assumptions

1. Infection is possible from at most one strain during follow-up
2. The relative prevalence of strains is constant over time
3. Equal exposure of vaccine and control groups
4. $\Pr(\text{Infection} | \text{Exp to strain } s, V) = \exp \{ \alpha_{0s} + \gamma_s V \}$

$$\longrightarrow \boxed{\beta_s = \gamma_s}$$

(Proof in Gilbert, Self, and Ashby, 1998, Biometrics)

- $OR(s) = \frac{RR^{pc}(s)}{RR^{pc}(0)}$
- $\beta_s = \log \left\{ \frac{RR^{pc}(s)}{RR^{pc}(0)} \right\}$
- $\beta_s - \beta_t = \log \left\{ \frac{RR^{pc}(s)}{RR^{pc}(t)} \right\}$

Alternative Ordinal Categorical Model

- **Cumulative strain categories model**
(McCullagh 1980)

$$\begin{aligned} \exp \{ \beta_s \} &= \frac{\Pr(Y > s | v) / \Pr(Y > s | u)}{\Pr(Y \leq s | v) / \Pr(Y \leq s | u)} \\ &= OR(> s), \quad s = 1, \dots, K - 1 \end{aligned}$$

- **Scored ordinal models**

- replace β_s with $s \times \textit{beta}$
- Scored models have **increased precision**, but stronger modeling assumptions

Nonparametric Tests for Differential VE

Null hypothesis: all $OR(s) = 1$

- **Nominal categorical:** Likelihood ratio
Chi-squared test (Armitage 1971)
- **Ordinal categorical:** test for trend in strain-specific odds ratios (Breslow and Day 1980)
- **Multiple vaccine dose groups**
 - Kruskal-Wallis test
 - Linear-by-linear association test
- **Exact tests:** StatXact software

Parametric Tests for Differential VE

- **MLR or cumulative categories**

Null hypothesis: all $\beta_s = 0$

- likelihood ratio Chi-squared test
- Zelen's test (1991)

Finer null hypothesis: a subset of $\beta'_s = 0$

- **Categorical scored models**

Null hypothesis: $\beta = 0$

- **Continuous model**

Null hypothesis: $\beta = 0$

- likelihood ratio, Wald, and score test

Example

- **Hepatitis B vaccine trial in New York**
(Szmuness et al., 1981)

	Hepatitis			
	B	A	non-A,B	
Placebo	63	27	11	101
Vaccine	7	21	16	44

- **Likelihood ratio statistic:** $\chi_2^2 = 30.2$, $p < 0.0001$

$$\frac{RR^{pc}(A)}{RR^{pc}(B)} = 7.0 \quad (2.7, 18.4) \text{ 95\% CI}$$

$$\frac{RR^{pc}(non-A,B)}{RR^{pc}(B)} = 13.1 \quad (4.3, 39.3) \text{ 95\% CI}$$

Example

- **Ordered categorical viral feature:**

Number sub/del to the prototype
hexapeptide tip sequence of V3 loop

E.g., VaxGen's MN/GNE8 rgp120 vaccine:

GPGRAF

Estimate

$$\frac{RR^{pc}(1 \text{ sub/del})}{RR^{pc}(GPGRAF)}$$

$$\frac{RR^{pc}(2 \text{ sub/del})}{RR^{pc}(GPGRAF)}$$

$$\frac{RR^{pc}(3^+ \text{ sub/del})}{RR^{pc}(GPGRAF)}$$

Generalized Logistic Regression Model (Gilbert et al., 1999, 2000)

- **Continuous analog of MLR model**

- parameterized MLR model: $\beta_s = g(s)\theta$

- g a deterministic function

- **Generalized Logistic Regression (GLR) model:**

$$Pr(Y = y|vaccine) = \frac{\exp\{g(y)\theta\}f(y)}{\int_0^\infty \exp\{g(z)\theta\}dF(z)}$$

$$f(y) \equiv Pr(Y = y|placebo)$$

- **Parametric component:**

- regression relationship

- **Nonparametric component:**

- baseline strain distribution F

Interpretation of GLR Model

- $g(y)\theta = \log \{OR(y)\}$

•
assumptions

•
 $= \log \left\{ \frac{RR^{pc}(y)}{RR^{pc}(0)} \right\}$

- $[g(y_1) - g(y_2)]\theta = \log \left\{ \frac{RR^{pc}(y_1)}{RR^{pc}(y_2)} \right\}$

- e.g. $g(y) = y$:

$$RR^{pc}(y + 1) = \exp \{ \theta \} RR^{pc}(y)$$

Summary of Sieve Model Parameters

1. MLR

$$e^{\beta_2} = OR(2), \dots, e^{\beta_K} = OR(K)$$

2. Scored MLR

$$e^{2\beta} = OR(2), \dots, e^{K\beta} = OR(K)$$

3. Cumulative categories

$$e^{\beta_2} = OR(> 1), \dots, e^{\beta_K} = OR(> K - 1)$$

4. Scored cumulative categories

$$e^{\beta} = OR(> 1), \dots, e^{\beta} = OR(> K - 1)$$

5. GLR

$$e^{g(y)\beta} = OR(y)$$

- In all cases assumptions as in Gilbert et al. (1998) imply the ORs equal ratios of strain-spec. RR^{pc} 's

Multidimensional Pathogen Variation

- The MLR and GLR models can accommodate pathogen variation described by multiple features
- **Examples**
 1. **cholera:** biotype, serotype, disease severity
 2. **rotavirus:** serotype, disease severity
 3. **HIV-1:** vast possibilities

Multivariate GLR Model

- $Y = (Y_1, \dots, Y_K) \in [0, \infty)^K$

- e.g. $K=2$:

$$Pr(Y = (y_1, y_2) | vaccine) =$$

$$\frac{\exp\{g_1(y_1)\theta_1 + g_2(y_2)\theta_2 + g_1(y_1)g_2(y_2)\theta_3\}}{\int_0^\infty \int_0^\infty \exp\{g_1(u_1)\theta_1 + g_2(u_2)\theta_2 + g_1(u_1)g_2(u_2)\theta_3\} dF(u_1, u_2)}$$

- Can investigate dependency of VE on marginal distance, adjusting for other distances

- e.g. estimate $\frac{RR^{pc}(y_1)}{RR^{pc}(y'_1)}$ adjusted for Y_2

- Can investigate interactions

- does $VE(Y_1, Y_2) = VE(Y_1)VE(Y_2)$?

Example

- Merck's Adenovirus-5 vaccine vector
 - includes core proteins coded by *gag*, *pol*, and *nef*
- $Y = (Y_{gag}, Y_{pol}, Y_{nef})$
 - Y_{gag} a metric based on *gag*
 - Y_{pol} a metric based on *pol*
 - Y_{nef} a metric based on *nef*
- Investigate how vaccine protection depends on heterogeneity in *gag*, *pol*, and *nef*

Example

- **Question:** What are the roles of CD4+ cellular responses and CD8+ cellular responses in conferring homologous and heterologous protection?
- $Y = (Y_{CD4+}, Y_{CD8+})$

Y_{CD4+} = strength of CD4+ T helper response against the vaccine strain- a **T help metric**

Y_{CD8+} = strength of CD8+ T cell response against the vaccine strain- a **CTL metric**

Example

- **Six-variate GLR model:** Investigate correlation of Merck's Ad 5 vaccine protection with MHC Class I and Class II T cell responses against *gag*, *pol*, and *nef*

$$(gag, pol, nef) \times (CD4+, CD8+)$$

$$Y = (Y_{CD4+,gag}, Y_{CD4+,pol}, Y_{CD4+,nef}, \\ Y_{CD8+,gag}, Y_{CD8+,pol}, Y_{CD8+,nef})$$

s-sample GLR Model

- s distinct covariate groups

$$x_1, \dots, x_s$$

$$Pr(Y = y|x_i) = \frac{\exp \left\{ \sum_{k=1}^d g_{ik}(y) \theta_k \right\} f(y)}{\int_0^\infty \exp \left\{ \sum_{k=1}^d g_{ik}(u) \theta_k \right\} dF(u)}$$

- multiple vaccine dose groups
- stratify by covariates
- adjust for low-dimensional covariates

Estimation in GLR Model

- Maximum likelihood estimation
- s -sample GLR model is a *semiparametric biased sampling model*:

$$Pr(Y = y|i) = \frac{w_i(y, \theta) f(y)}{\int_0^{\infty} w_i(u, \theta) dF(u)}, \quad i = 1, \dots, s$$

-e.g. two-sample GLR model:

$$w_1(y, \theta) \equiv 1, \quad w_2(y, \theta) = \exp\{g(y)\theta\}$$

Partial Likelihood Estimation

- Partial likelihood

$$L_{n1}(\theta, V) = \prod_{i=1}^s \prod_{j=1}^{n_i} \left[\frac{\lambda_{ni} w_i(y_{ij}, \theta) V_i}{\sum_{k=1}^s \lambda_{nk} w_k(y_{kj}, \theta) V_k} \right]$$

with $V_i \equiv \frac{1}{w_i(F, \theta)}$

Maximization Algorithm

1. Maximize L_{n1} over θ and V , $V > 0$
2. Compute Vardi's (1985) NPMLE \hat{F}_n
using 'known' weight functions $w_i(\cdot, \hat{\theta}_n)$

Properties of MLE in GLR Model

- Described in Gilbert et al. (1999, Biometrika; 2000, Ann Stat)
- GLR model identifiable
- GLR model uniquely estimable- log profile partial likelihood strictly concave
- $(\hat{\theta}_n, \hat{F}_n)$ uniformly consistent, asymptotically normal, and asymptotically efficient
- Confidence intervals and variance estimation
 1. sample estimator of generalized Fisher information
 2. bootstrap
- Satisfactory finite-sample properties
- **Comparable to MLE in Cox model**

Simulation Study of Gilbert et al. (1999)

- Study performance of $(\hat{\theta}_n, \hat{F}_n)$
- Investigate bias and estimation of variance via observed inverse generalised Fisher information and via the bootstrap
- Investigate power of likelihood ratio, Wald and score tests of $H_0 : \theta = 0$ (no differential vaccine protection), and coverage accuracy of corresponding CIs

Simulation Study of Gilbert et al. (1999)

- Y = percent amino acid difference between an infecting virus and the global subtype B consensus in the V3 loop
- Specify true log RR ratio as

$$\log \{RR^{pc}(y)/RR^{pc}(0)\} = \frac{y}{35}\theta$$

-Set $\theta = 0, 2,$ and 4

$$\theta = 0 : RR^{pc}(35) = 1.0RR^{pc}(0)$$

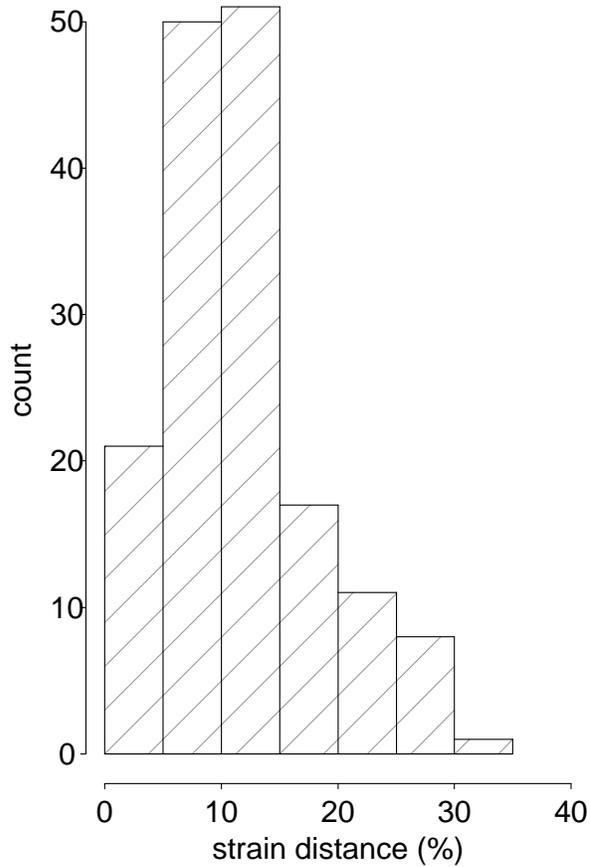
$$\theta = 2 : RR^{pc}(35) = 7.39RR^{pc}(0)$$

$$\theta = 4 : RR^{pc}(35) = 54.6RR^{pc}(0)$$

Simulation Study of Gilbert et al. (1999)

- Consider 4 baseline distribution functions F
 - Unif(0, 35)
 - Normal(0.1157, 0.710²)
 - Expon(0.1157/2)
 - Thai empirical (based on 94 sequences)
 - 0.1157 and 0.710 based on 159 subtype B V3 loop sequences in the LANL database
- 4 sample sizes (numbers of infections):
 $(n_p, n_v) = (100, 100), (100, 50), (50, 50), (50, 25)$
- Simulations based on 1000 trials

(a) U.S. distance distribution



(b) Thai distance distribution

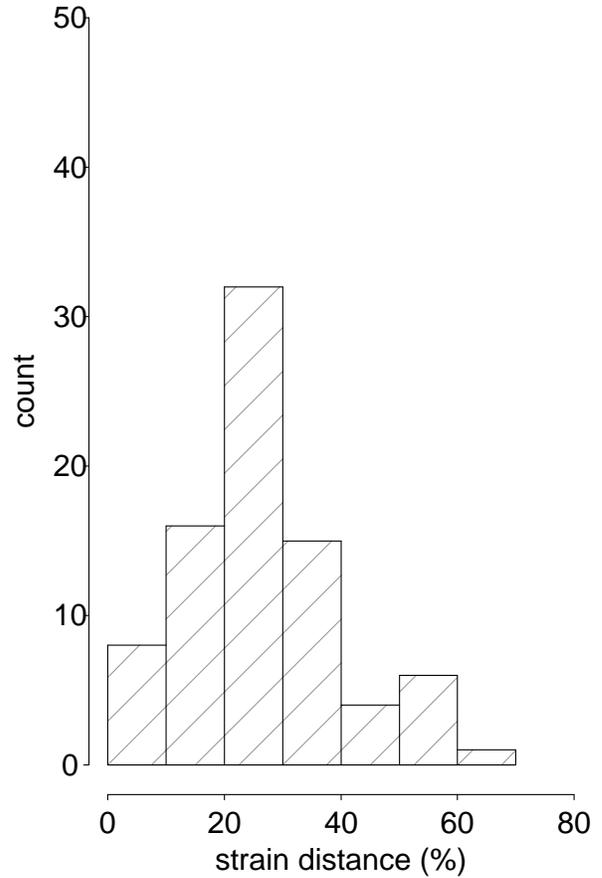


Fig. 1(a) shows the distribution of the V3 loop amino acid distance between 159 U.S. subtype B sequences and the global subtype B consensus sequence. (b) shows the distribution of 94 V3 loop amino acid distances of infecting strains in Thailand.

Table 1. *Bias and variance of the maximum likelihood estimator $\hat{\theta}$; finite-sample variance s^2 , observed generalised Fisher information variance estimate var_F , bootstrap variance estimate var_B*

			Uniform				Normal			
n_p	n_v	θ	<i>bias</i>	F s^2	var_F	var_B	<i>bias</i>	F s^2	var_F	var_B
100	100	0	-0.01	0.24	0.24	0.20	-0.03	0.47	0.50	0.33
50	25	0	-0.00	0.71	0.76	0.52	0.24	1.21	1.56	0.96
100	100	2	0.02	0.30	0.29	0.30	-0.01	0.52	0.54	0.57
50	25	2	0.08	0.97	0.95	1.01	0.20	1.63	1.71	1.78
100	100	4	0.08	0.49	0.46	0.48	0.01	0.66	0.67	0.72
50	25	4	0.27	1.87	1.67	1.77	0.29	2.28	2.14	2.67

			Exponential				Thai			
n_p	n_v	θ	<i>bias</i>	F s^2	var_F	var_B	<i>bias</i>	F s^2	var_F	var_B
100	100	0	0.07	0.67	0.79	0.47	0.01	0.49	0.52	0.38
50	25	0	0.38	1.48	2.49	1.66	0.21	1.26	1.63	0.97
100	100	2	0.05	0.64	0.64	0.67	0.04	0.54	0.53	0.54
50	25	2	0.10	1.68	1.91	1.92	0.12	1.80	1.65	1.55
100	100	4	0.06	0.64	0.65	0.67	0.09	0.67	0.62	0.66
50	25	4	0.23	1.94	1.87	2.27	0.25	2.19	1.96	2.23

Table 2. *Power of likelihood ratio, Wald and score tests of $H_0 : \theta = 0$ with $\alpha = 0.05$*

n_p	n_v	F θ	Uniform	Normal	LRatio	Exponential		Thai
			LRatio	LRatio		Wald	Score	LRatio
100	100	0	0.05	0.05	0.02	0.02	0.04	0.05
50	25	0	0.04	0.03	0.02	0.02	0.04	0.03
100	100	2	0.97	0.81	0.79	0.77	0.86	0.84
50	25	2	0.61	0.40	0.37	0.33	0.51	0.41
100	100	4	1.00	1.00	1.00	1.00	1.00	1.00
50	25	4	0.99	0.93	0.99	0.97	0.99	0.92

Table 3. *Score statistic confidence intervals about θ , $\alpha = 0 \cdot 05$*

n_p	n_v	F θ	Uniform	Normal	Exponential	Thai
100	100	0	-1.16,1.12	-1.55,1.58	-1.38,1.84	-1.63,1.88
50	25	0	-1.68,2.14	-1.81,3.35	-1.71,3.17	-1.60,3.07
100	100	2	0.93,3.03	0.66,3.49	0.52,3.59	0.62,3.58
50	25	2	0.28,4.01	-0.15,4.42	0.18,4.92	-0.12,4.68
100	100	4	2.77,5.50	2.58,5.53	3.01,4.78	2.53,6.11
50	25	4	1.81,6.60	1.67,7.26	2.33,6.41	1.59,6.78

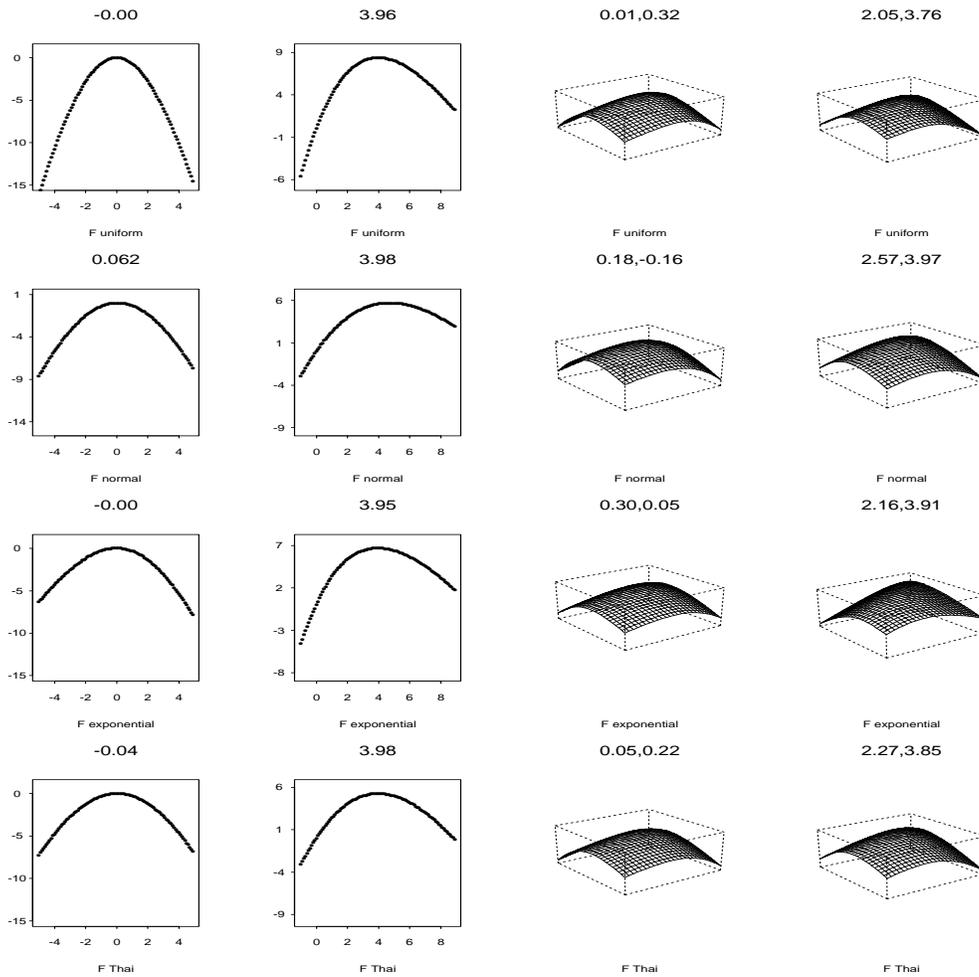


Fig. 2. shows the log profile partial likelihood versus θ for a spectrum of generated data sets. The obtained $\hat{\theta}$ is written above each plot. The first two columns are plots for the two-sample problem representing sample size $n_p = 100, n_v = 50$, and the second two columns are plots for the three-sample problem representing sample size $n_p = 100, n_{v1} = 50, n_{v2} = 25$.

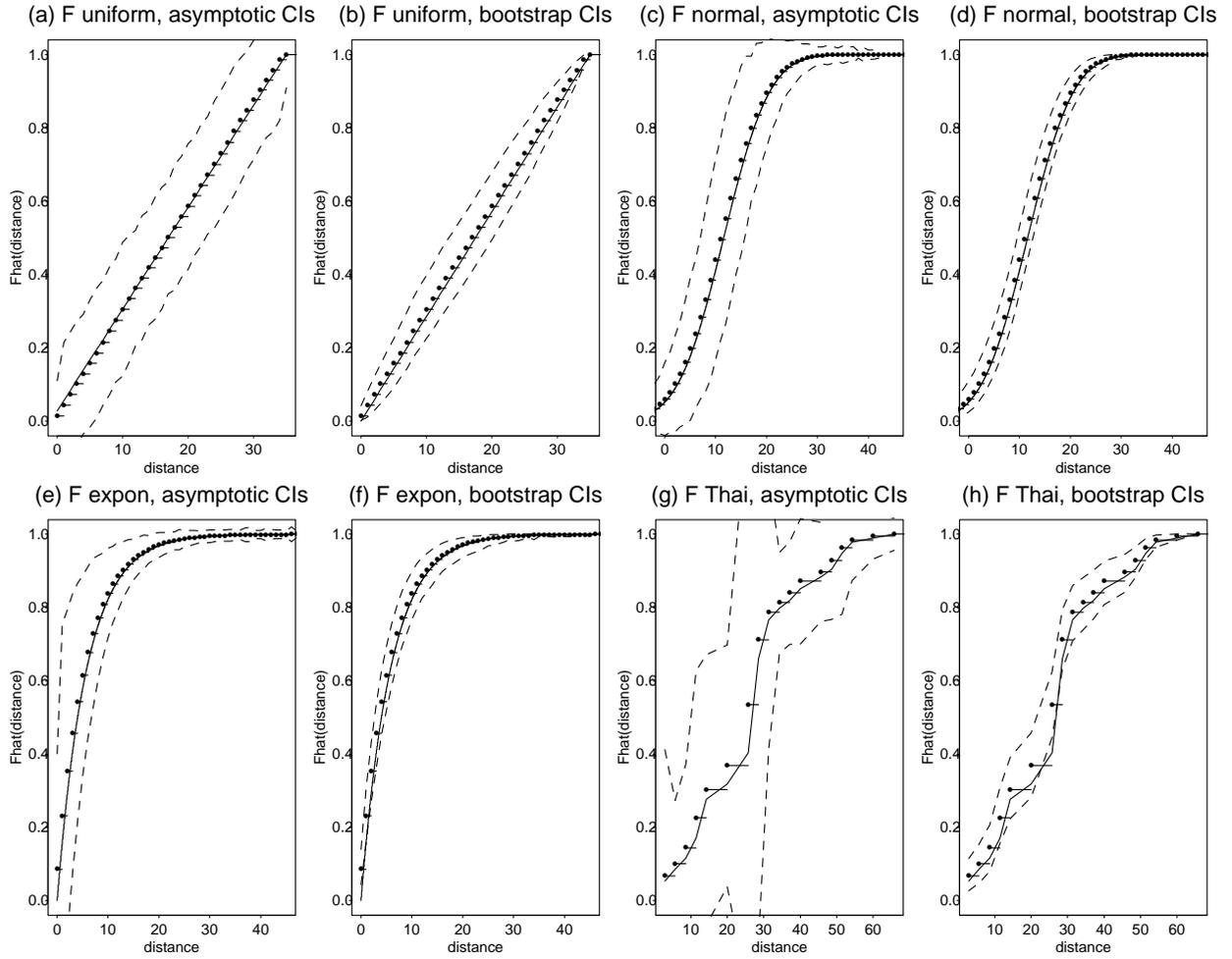


Fig. 3. portrays data sets generated from $n_p = 100, n_v = 50, \theta = 2$. (a),(c),(e) and (g) show the mean of \hat{F} across the 1000 replications, with 95% symmetric asymptotic normal approximation confidence bands. The true distribution is depicted by a solid line. (b),(d),(f) and (h) include 95% bootstrap confidence bands.

Pseudo-Example from Gilbert et al. (1999)

- Generated a single dataset using the empirical Thai strain distribution and assuming that:

$$VE^{pc}(y \leq 0.10) = 50\%$$

$$VE^{pc}(0.11 \leq y \leq 0.20) = 40\%$$

$$VE^{pc}(0.21 \leq y \leq 0.30) = 30\%$$

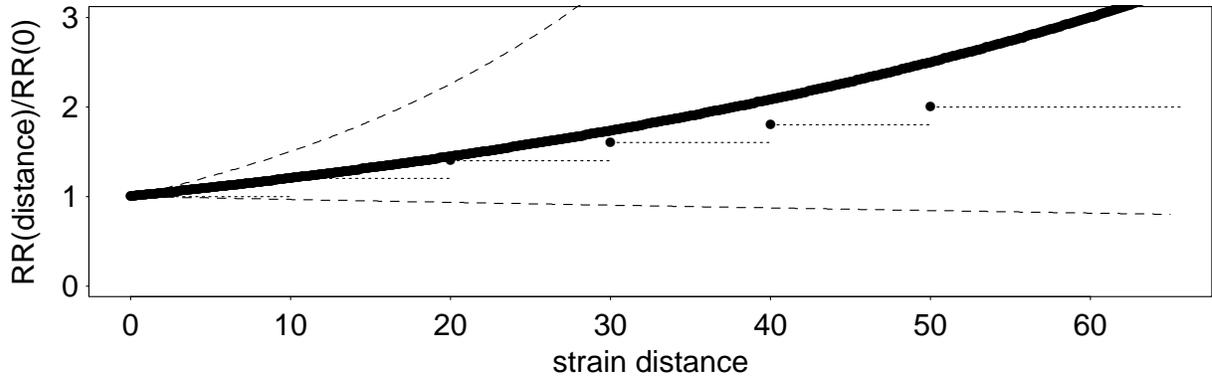
$$VE^{pc}(0.31 \leq y \leq 0.40) = 20\%$$

$$VE^{pc}(0.41 \leq y \leq 0.50) = 10\%$$

$$VE^{pc}(0.51 \leq y) = 0\%$$

- Number of infections $n_p = 100, n_v = 69$
- Fit the same GLR model studied in simulations
- LR, Wald and score tests: $p = 0.10, 0.10, 0.09$

(a) Vaccine protection versus strain distance



(b) NPMLE \widehat{F} and F versus strain distance

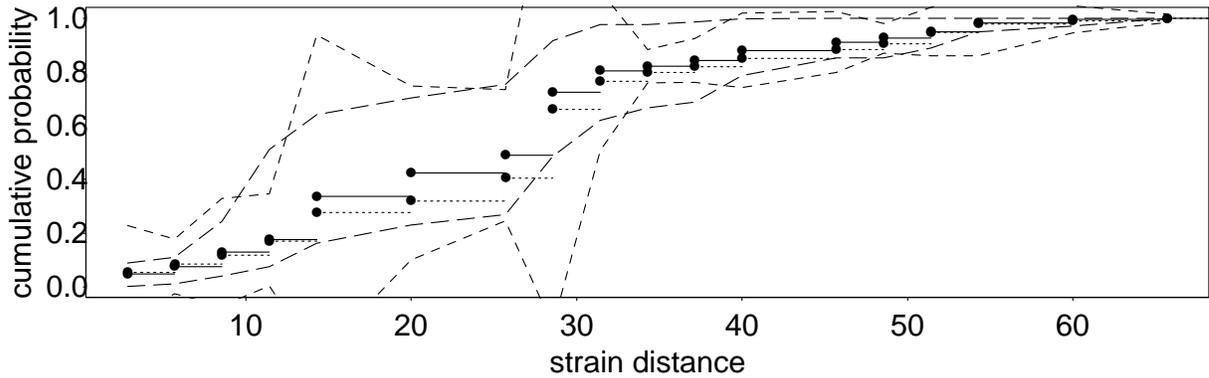


Fig. 4. (a) shows the estimated ratio of relative risks $\widehat{RR}(y)/\widehat{RR}(0)$ versus strain distance y as a solid line. The broken lines are profile likelihood-based confidence intervals, and the dotted line step function is the true relative risk ratio. (b) shows \widehat{F} as solid lines, with 95% asymptotic normal approximation confidence bands as short dashed lines and 95% bootstrap confidence bands as long dashed lines. The true F is portrayed as dotted lines.

III. Incorporating Time to Infection Diagnosis

- Cause-specific hazards approach
Prentice et al. (1978, Biometrics);
foundational paper
- Suppose K circulating strains
- Let Y_1, \dots, Y_K be *conceptual* or *latent*
infection Dx times corresponding to the K strains
- Classic competing risks data: Data are
iid observations $(T_i, \delta_i, S_i, z_i)$

$$T_i = \min(Y_1, \dots, Y_K)$$

δ_i = failure indicator (1 if infected)

S_i = infecting strain (NA if not infected)

z_i = covariate vector

Cause-specific Hazards

- Prentice et al. (1978) emphasized that all functions of cause-specific hazards λ_s are estimable from the data $(T_i, \delta_i, S_i, z_i)$

$$\lambda_s(t|z) = \frac{\lim_{\Delta t \searrow 0} \Pr(t \leq T < t + \Delta t, S = s | T \geq t, z)}{\Delta t}$$

Cause-Specific Proportional Hazards Model

- Prentice et al. (1978) proposed a *cause-specific proportional hazards model*:

$$\lambda_s(t|z) = \exp\{\beta_s^T z\} \lambda_s(t|0)$$

- arbitrary baseline hazard $\lambda_s(t|0)$
- when z is vaccination status,

$$\beta_s = \log\{\lambda_s(t|\text{vaccine})/\lambda_s(t|\text{placebo})\}$$

- $\beta_s = \log$ -relative hazard (vaccine vs placebo) of infection by strain s
 - $VE(s) \equiv 1 - \exp\{\beta_s\}$ measures strain-specific vaccine efficacy
- β_s can be estimated by the standard maximum partial likelihood estimator (MPLE), treating infection by all non- s strains as censoring

Interpretation of β_s

- λ_s has a “crude” interpretation, which is restricted to the particular vaccine trial conditions
- Additional assumptions needed for the strain s -specific vaccine efficacy estimate

$$\widehat{VE}(s) = 1 - \exp\{\widehat{\beta}_s\}$$

to have a meaningful biological interpretation

- Would like $VE(s) = VE^{pc}(s)$, where $VE^{pc}(s)$ is one minus the relative conditional probability (vaccine vs placebo) of a specified amount of exposure to strain s causing HIV infection

Interpretation of β_s

- **Assumptions:**

A1: For each strain $s \in \{1, \dots, K\}$, the probability of infection with strain s resulting from a specified amount of exposure is homogeneous and constant over time among vaccinated and placebo subjects, so that vaccination reduces the transmission probability by the same fraction $\exp\{\gamma_s\}$ for all vaccinees (i.e., “leaky” protection against each strain; Halloran, Haber, and Longini, 1992)

A2: The pattern of risk behavior and exposure to each strain $s \in \{1, \dots, K\}$ during the follow-up period $[0, \tau]$ for a trial participant is the same whether vaccine or placebo was assigned (justified by randomization and blinding)

Interpretation of β_s

- Under A1 and A2, the crude hazard ratio

$$\exp\{\beta\} = \frac{\lambda_s(t|\text{vaccine})}{\lambda_s(t|\text{placebo})}$$

equals the biologically interpretable parameter

$$\exp\{\gamma_s\} = 1 - VE^{pc}(s)$$

- Therefore, under A1 and A2 the MPLE $\hat{\beta}_s$ in the strain s -specific proportional hazards model estimates γ_s (and $\widehat{VE}(s)$ estimates $VE^{pc}(s)$)
- Based on Rhodes, Halloran, and Longini (1992, JRSS B), under randomization and blinding, $\hat{\beta}_s$ should be \approx unbiased if the strain s infection rate is low

Sketch of Proof (from Gilbert, 2000, Stat Med)

$$\lambda_s(t|z) = \lambda_{Es}(t|z) \times$$
$$\Pr(t \leq T < t + \Delta t, S = s | T \geq t, z,$$

exposed to strain s in $[t, t + \Delta t)$),

- $\lambda_{Es}(t|z)$ is the Markov intensity of the counting process counting exposures to strain s for participants with covariate z

-The second term conditions on a specified exposure during $[t, t + \Delta t)$,
e.g., on a sexual or needle contact with a strain s -infected individual

Sketch of Proof (from Gilbert, 2000, Stat Med)

- A1 implies a constant strain-specific transmission probability over time in each group
- A2 implies $\lambda_{Es}(t|\text{vacc}) = \lambda_{Es}(t|\text{plac})$ for all t
- Together these results imply $\beta_s = \gamma_s$
- Therefore the MPLE $\hat{\beta}_s$ in the strain s -specific proportional hazards model estimates γ_s

Assessing Differential Protection

- Since each γ_s is estimated from a separate model fit, the strain-specific proportional hazards models do not permit direct comparisons of vaccine efficacy across strains
- Lunn and McNeil (1995, Biometrics) showed how to reparametrize the strain s Cox model so that $\exp\{\beta_s\}$ ($s \geq 2$) equals

$$\frac{\lambda_s(t|v)}{\lambda_s(t|u)} / \frac{\lambda_1(t|v)}{\lambda_1(t|u)}$$

- β_s measures relative vaccine efficacy against strain s compared to the reference strain 1

- Therefore standard Cox model software (e.g., in Splus/R) can be used to estimate β_s with a confidence interval

Tests for Vaccine Efficacy Against Strain s

- Standard Cox model software provides tests of

$$H_{0\text{haz}} : \lambda_s(t|\text{vaccine}) = \lambda_s(t|\text{placebo})$$

- Through creative coding also provides tests of

$$H_0 : VE(1) = VE(2)$$

- An alternative to a hazards-based approach would apply Gray's (1988, Ann Stat) method to test different *cumulative incidence functions*,
 $H_{0\text{ci}} : F_{vs} = F_{ps}$, where

$$F_{vs}(t) = \Pr(T \leq t, S = s|\text{vaccine})$$

$$F_{ps}(t) = \Pr(T \leq t, S = s|\text{placebo})$$

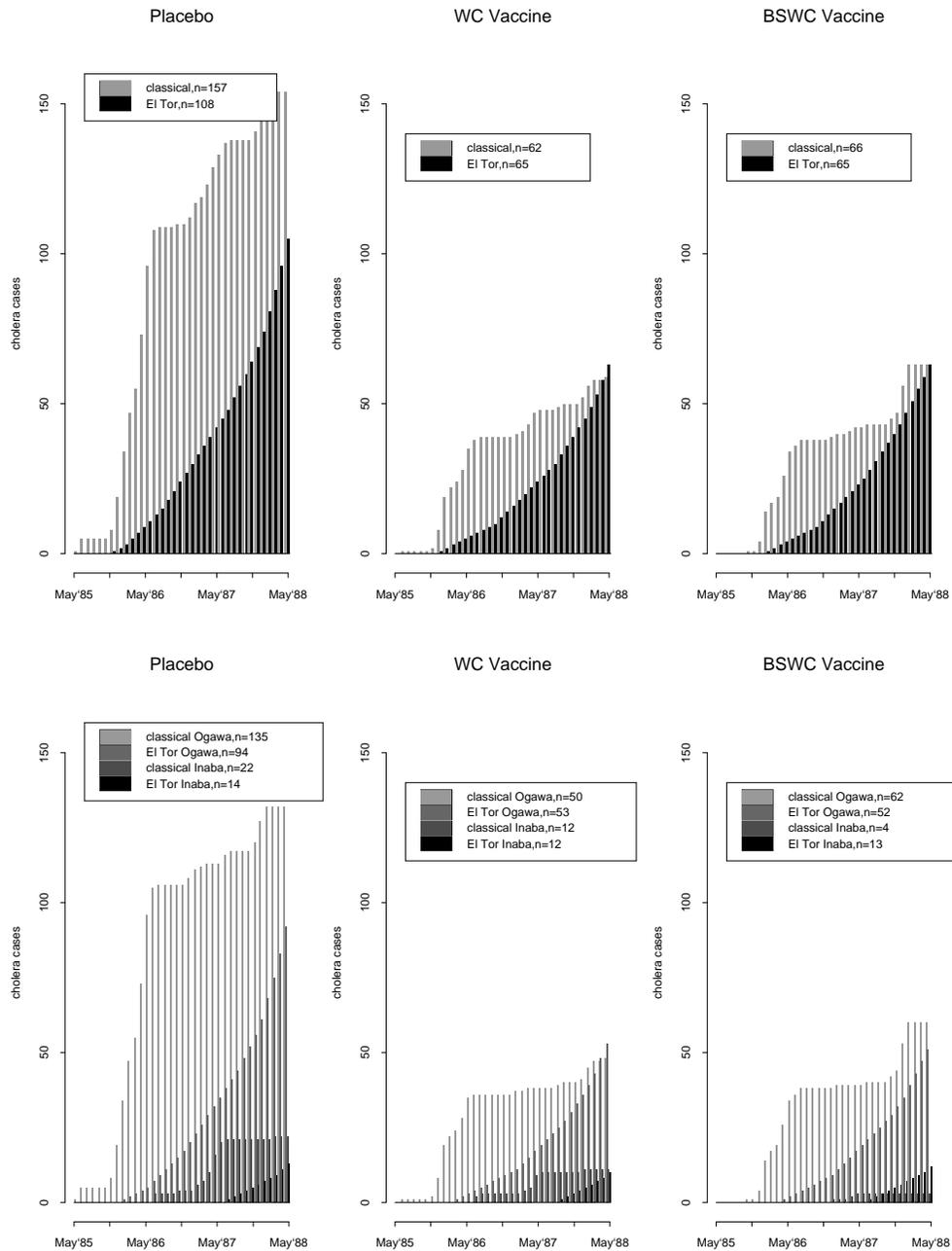
Example: Oral Cholera Vaccine Trial

- A randomized double-blind field trial was conducted in rural Bangladesh among children over 2 years and adult women (1985-1992)
- Assessed the efficacy of B subunit killed whole cell (BSWC) and killed whole-cell-only (WC) oral cholera vaccines
- Case endpoint: First diarrheal episode in which *Vibrio cholerae 01* was isolated
- 2 cholera biotypes (classical, El Tor) and 2 cholera serotypes (Inaba, Ogawa) circulated during the trial
 - The causal infecting biotype and serotype was measured for each case

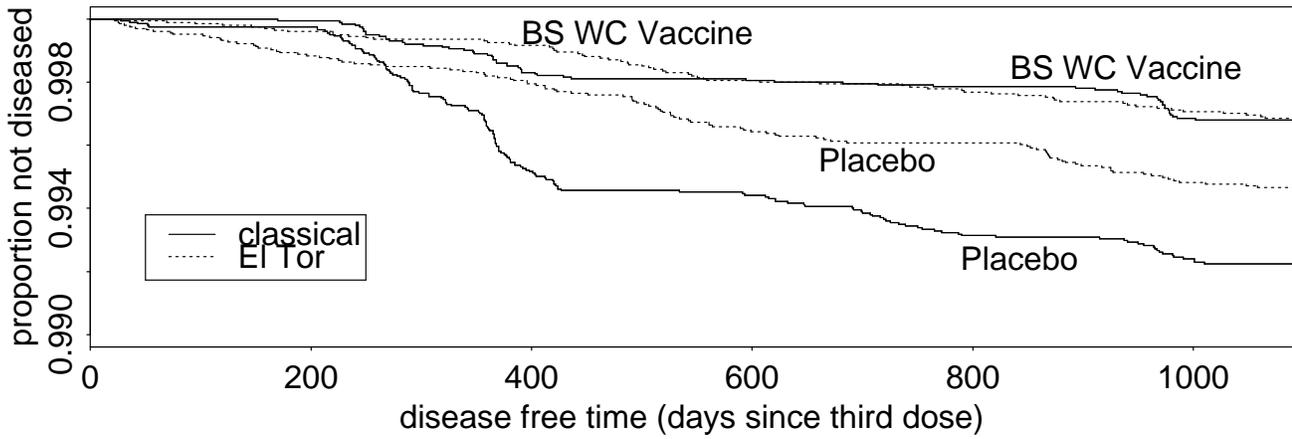
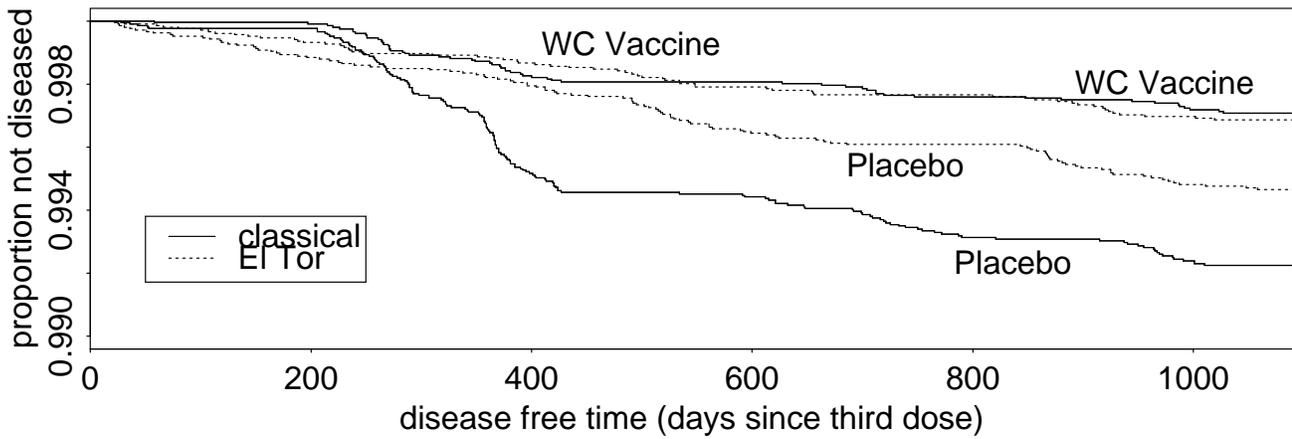
Example: Oral Cholera Vaccine Trial

- Analysis of the study population of 62,285 children and women who received three doses the BSWC vaccine (20,705), the WC vaccine (20,743), or the *Escherichia coli* K12 strain placebo (20,837)
- Overall the two vaccines performed similarly
 - Each vaccine had about 50% efficacy sustained for 2 or 3 years, waning to nil at 5 years

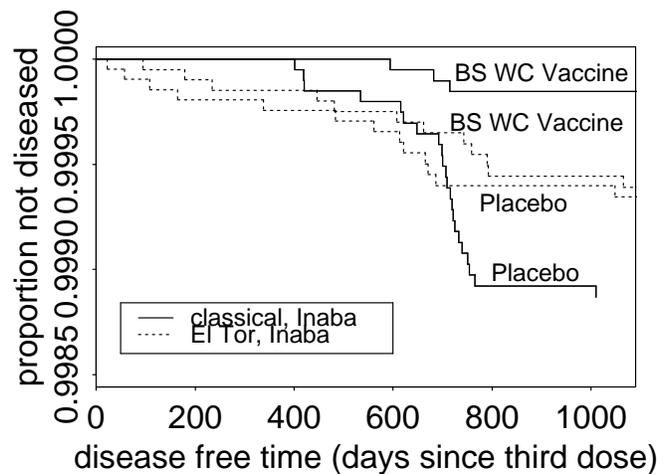
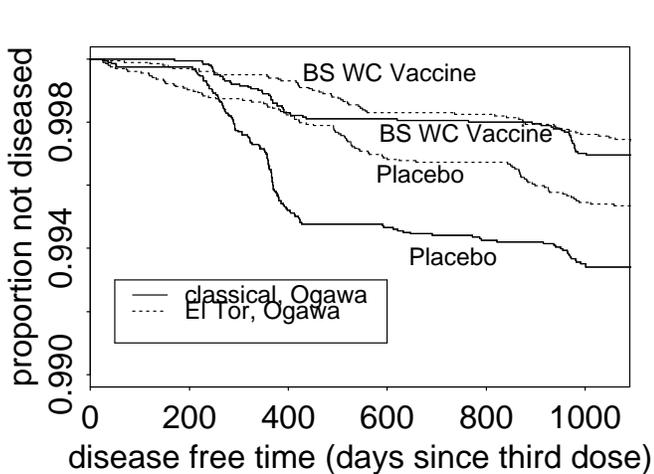
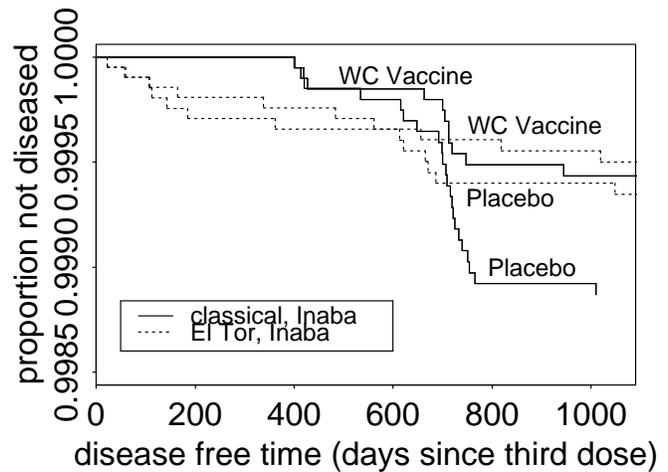
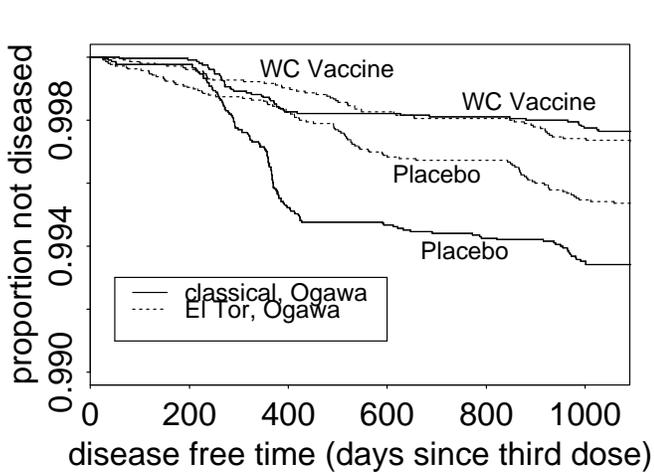
Cum. Inc. of Strain-Specific Cholera Cases



Biotype-specific 1-Cum. Inc. Curves



Biotype/Serotype-specific 1-Cum. Inc. Curves



Sieve Analyses to Assess Differential Protection

- Conducted sieve analysis to compare $VE(1)$ and $VE(2)$ (where 1 and 2 indicate different biotypes or serotypes) using 4 methods:
 1. MLR model
 2. MLR model stratified by the 3 years of follow-up (account for temporal trends in shifting biotype/serotype prevalence)
 3. Cause-specific Cox model with Lunn and McNeil recoding
 4. Cause-specific Cox model with Lunn and McNeil recoding and with a proportional baseline risks assumption $\lambda_1(t|0) = \lambda_2(t|0)$
- For cause-specific Cox model, results obtained using standard coxph function in Splus/R

Fit of Sieve Models to Cholera Data

- Compare VE for El Tor vs Classical

Vaccine	Model	$\hat{\beta}^a$	$SE(\hat{\beta})$	Robust $SE(\hat{\beta})$	$\exp\{\hat{\beta}\} = \frac{\exp\{\gamma_{El}\}}{\exp\{\gamma_{Cl}\}}$	95% CI ^b	P-value
WC	MLR	0.421	0.217		1.524	(0.946,2.332)	0.052
WC	Stratified MLR	0.389	0.219		1.475	(0.961,2.265)	0.076
WC	PH ^c	0.433	0.217	0.218	1.541	(1.006,2.360)	0.047
WC	PH, PBR ^d	0.428	0.217	0.217	1.534	(1.002,2.347)	0.049
BS WC	MLR	0.359	0.215		1.432	(0.940,2.181)	0.095
BS WC	Stratified MLR	0.318	0.221		1.375	(0.891,2.122)	0.150
BS WC	PH	0.369	0.215	0.212	1.446	(0.949,2.204)	0.085
BS WC	PH, PBR	0.365	0.215	0.212	1.440	(0.946,2.194)	0.089

^a $\beta = \gamma_{El\ Tor} - \gamma_{classical}$

^b 95% CIs derived from a normality approximation and the information matrix

^c PH model fit by duplication Method B of Lunn and McNeil

^d PH model fit under the proportional baseline risks assumption by duplication Method A of Lunn and McNeil

Summary of Results

- For each vaccine, the 4 methods perform similarly
 - Result explained by the very low failure rate
- Results suggest that both vaccines protect $\approx 50\%$ better against classical than El Tor cholera
- A possible explanation is that the vaccine contains 3 times as many Classical as El Tor antigens

Summary of Utility of Sieve Analysis Methods

1. Statistical inference of differential vaccine protection according to a pathogen variation structure chosen *a priori*
2. Exploratory tools for identifying which pathogen features are potentially correlated with vaccine protection