

Supplementary Appendix

1. Statistical Methods

Let S denote the time when the participant is vaccinated, and T denote the time when the participant acquires SARS-CoV-2 infection (as defined by seroconversion or detectable viral RNA), with both times measured in days from the start of the clinical trial. In addition, let X denote baseline risk factors (e.g., age, occupation, race, health conditions). We specify that the hazard function of T is related to S and X through the Cox [1] regression model

$$\lambda(t|S, X) = \lambda_0(t)e^{\beta^T X + \eta(t-S)A(t)}, \quad (1)$$

where $A(t) = I(S < t)$, $I(\cdot)$ is the indicator function, $\lambda_0(\cdot)$ is an arbitrary baseline hazard function, β is a set of regression parameters representing the effects of baseline risk factors, and $\eta(\cdot)$ is the log hazard ratio characterizing the time-varying effect of vaccination. Under this formulation, the baseline hazard function varies over the calendar time, and the effect of vaccine on the risk of infection depends on the time elapsed since vaccination.

We define vaccine efficacy at day t as the proportionate reduction in the hazard rate of infection at day t for individuals who were vaccinated t days ago compared with those who have not been vaccinated, i.e., $VE_h(t) = 1 - e^{\eta(t)}$. In addition, we define t -day vaccine efficacy as the proportionate reduction in the attack rate or cumulative incidence of infection by day t for individuals who were vaccinated t days ago compared with the non-vaccinated individuals, i.e., $VE_a(t) = 1 - V(t)/t$, where $V(t) = \int_0^t e^{\eta(u)} du$ [2]. Finally, we consider vaccine efficacy in reducing the attack rate over a certain time period, say (t_1, t_2) :

$$VE_a(t_1, t_2) = 1 - \frac{V(t_2) - V(t_1)}{t_2 - t_1}.$$

Clearly, all three VE measures are simple functions of the log hazard ratio $\eta(\cdot)$.

For economic and logistical reasons, antibody (or RT-PCR) tests can only be performed infrequently for each participant. Thus, SARS-CoV-2 infection, as defined by seroconversion or detectable viral RNA, is only known to occur over a broad time interval, such that T must be treated as an interval-censored event time. Let L be the time of the last negative test, and R be the time of the first positive test, such that T is known to lie in the time interval $(L, R]$. (To be more precise, L is the last seronegative test date minus 7 days and R is the first seropositive test date minus 7 days, because it takes approximately 7 days after initial SARS-CoV-2 acquisition for the antibody test to register positive.) In addition, let E be the time when the participant enters the clinical trial. Like T and S , the time variables E , L , and R are measured from the start of the clinical trial. We assume that E , L , R , and S are independent of T conditional on X .

For a clinical trial with a total of n participants, the data consist of $(E_i, L_i, R_i, S_i, X_i)$ ($i = 1, \dots, n$). The likelihood takes the form

$$\prod_{i=1}^n \left[\exp \left\{ - \int_{E_i}^{L_i} e^{\beta^T X_i + \eta(t-S_i)A_i(t)} d\Lambda_0(t) \right\} - \exp \left\{ - \int_{E_i}^{R_i} e^{\beta^T X_i + \eta(t-S_i)A_i(t)} d\Lambda_0(t) \right\} \right],$$

where $\Lambda_0(t) = \int_0^t \lambda_0(s) ds$.

This likelihood involves two infinite-dimensional functions $\Lambda_0(\cdot)$ and $\eta(\cdot)$, which are not identifiable if both are unrestricted. We let $\Lambda_0(\cdot)$ be completely nonparametric and estimate it by a step function with non-negative jumps at the unique values of $L_i > 0$ and $R_i < \infty$ ($i = 1, \dots, n$). We approximate $\eta(\cdot)$ by a sequence of B-splines functions, denoted by $B_1(t), \dots, B_K(t)$, such that $\eta(t) = \sum_{k=1}^K \gamma_k B_k(t)$. Write $\gamma = (\gamma_1, \dots, \gamma_K)^T$, and $Z_i(t) = (B_1(t - S_i)A_i(t), \dots, B_K(t - S_i)A_i(t))^T$. Then the likelihood becomes

$$\prod_{i=1}^n \left[\exp \left\{ - \int_{E_i}^{L_i} e^{\beta^T X_i + \gamma^T Z_i(t)} d\Lambda_0(t) \right\} - \exp \left\{ - \int_{E_i}^{R_i} e^{\beta^T X_i + \gamma^T Z_i(t)} d\Lambda_0(t) \right\} \right].$$

This is the interval-censored data likelihood for the standard Cox model with time-independent covariates X and time-dependent covariates Z [3], except that the event time is measured from the start of the study rather than the participant's entry time. We compute the nonparametric maximum likelihood estimator for $(\beta, \gamma, \Lambda_0)$, denoted by $(\hat{\beta}, \hat{\gamma}, \hat{\Lambda}_0)$, through an EM algorithm based on unobserved Poisson random variables [3]. We then estimate $\eta(t)$ and $V(t)$ by $\hat{\eta}(t) = \sum_{k=1}^K \hat{\gamma}_k B_k(t)$ and $\hat{V}(t) = \int_0^t e^{\hat{\eta}(u)} du$, respectively.

By appealing to modern empirical process theory [4], we can show that $\hat{\beta}$, $\hat{\gamma}$, and $\hat{\Lambda}_0(\cdot)$ are consistent. In addition, $\hat{\beta}$ and $\hat{\gamma}$ are asymptotically normal and their covariance matrix can be consistently estimated by the Hessian matrix [3,5] of the profile log-likelihood for (β, γ) , where the log-likelihood is maximized with respect to Λ_0 for fixed β and γ via the EM algorithm. These results allow us to estimate $VE_h(t)$, $VE_a(t)$, and $VE_a(t_1, t_2)$, construct confidence intervals, and perform hypothesis testing.

Remark. In our previous work on (potentially right-censored) symptomatic disease, we approximate $\log \lambda_0(\cdot)$ by B -spline functions while letting $\eta(\cdot)$ be completely nonparametric [2]. With interval-censored data, a completely nonparametric function cannot be estimated at the parametric rate, making it difficult to construct confidence intervals. Thus, we let $\lambda_0(\cdot)$ be completely nonparametric and approximate $\eta(\cdot)$, which is the parameter of main interest, by B -spline functions. One benefit of this approach is that it provides a unified framework to study constant versus time-varying VE (by choosing appropriate B -spline functions). This framework also unifies the analysis of symptomatic disease and asymptomatic infection because potentially right-censored data can be treated as a special case of interval-censored data. For potentially right-censored data, we can adopt very flexible B -spline functions for $\eta(\cdot)$; for truly interval-censored data, we have to be more rigid unless the sample size is very large, the infection rate is high, or antibody/RT-PCR tests are performed frequently.

2. Simulation Studies

2.1. Antibody Tests

We designed the first series of simulation studies to mimic the BNT162b2 phase 3 trial [6]. We assumed that 40,000 participants entered the study at a constant rate over four months, i.e., $E \sim \text{Uniform}(0, 4)$ months. (In the actual trial, the number of participants was slightly below 40,000, after exclusion of those who were seropositive at baseline.) We created a composite baseline risk score X , which takes values 1, 2, 3, 4, and 5 with equal probability. We randomly assigned half of the participants at study entry to vaccine and half to placebo. We generated the infection time T from model (1) with $\beta = 0.2$ and

$$\log \lambda_0(t) = -5.5 + 0.1t - 0.3(t - 7)_+,$$

where $t_+ = t$ if $t > 0$ and 0 otherwise. We assumed that VE_h starts at 0 at $t = 0$, increases to some maximum value at $t = t_m$, and then stays constant or decreases gradually over time. Specifically, we set $\eta(t) = b_1 t$ ($0 \leq t < t_m$) and $\eta(t) = b_1 t_m + b_2(t - t_m)$ ($t \geq t_m$), where $t_m = 4$ weeks, and b_1 and b_2 were chosen such that $VE_h(t_m) = 0.8$ and $VE_h(1 \text{ year}) = 0.8$ or 0.

In the BNT162b2 phase 3 trial [6], serum samples were scheduled to be drawn on Day 1, Day 22, Day 52, and Day 209 (as measured from the participant’s entry time) [6]. To allow for small random departures from the schedule, we used Day 1, Day 22 + Uniform(-1,3) days, Day 52 + Uniform(-2, 8) days, and Day 209 + Uniform(-5,10) days.

Serum samples were also drawn at the crossover visits. We considered:

Priority-dependent crossover: Crossover occurs at month $(11 - X + G)$ of the study, where G follows the exponential distribution with mean of 0.5 month.

Priority-independent crossover: Crossover occurs at month $6 + G$ of the study, where G follows the exponential distribution with mean of 0.5 month.

We assumed that the analysis is performed at 10.5 months after the start of the study, such that only the blood samples that were drawn before the 10.5 month mark can be included. We considered both blinded and unblinded crossover designs. Under blinded crossover, participants receive the opposite of their original assignments, and we used all the data that are collected before the time of analysis. At the point of unblinded crossover, participants are notified of their original assignments, and placebo participants receive the vaccine; we disregarded any data collected after unblinded crossover in order to avoid bias due to behavioral confounding.

We applied the proposed methods to each simulated dataset by setting $\eta(t)$ in model (1) to be piecewise linear with a change point placed at t_m and with the slope after t_m fixed at 0 or estimated from data. For comparison, we performed maximum partial likelihood estimation of the same model with the same data by treating the time of the first positive antibody test as a potentially right-censored event time and by using Efron’s method of handling tied event times.

2.2. RT-PCR Tests

We conducted a second series of simulation studies to mimic the Prevent COVID U study. In our simulation, a total of 12,000 participants enter the study at a constant rate over one month; half of them receive the Moderna vaccine at enrollment and the other half 4 months later. We generated the infection time T from model (1) without X and with $\log \lambda_0(t) = -4.0 - 0.2t$. We assumed the same VE patterns as in the first series of simulation studies, but the change point was set at 6 weeks instead of 4 weeks.

We investigated various swabbing/RT-PCR testing schedules, ranging from every day to every 2 weeks. All participants are followed for 4 months, and the study ends at month 5. In addition, we considered a scenario in which placebo participants may receive vaccines outside of the study 1 month after enrollment. We assumed that the time to outside vaccination follows the Weibull distribution with shape parameter of 3 and scale parameter of 4, such that the cumulative probability of outside vaccination is approximately 50%.

For each simulated dataset, we analyzed the data in the same way as in the first series of simulation studies. Specifically, we implemented both the proposed method and its right-censored data counterpart, to be referred to as the naive method. We discarded the data collected on placebo participants after they received outside vaccines.

References

1. Cox DR (1972). Regression models and life-tables. *J Roy Stat Soc B*, 34: 187-202.
2. Lin DY, Zeng D, Gilbert PB (2021). Evaluating the long-term efficacy of COVID-19 vaccines. *Clin Infect Dis*, ciab226, <https://doi.org/10.1093/cid/ciab226>
3. Zeng D, Mao L, Lin DY (2016). Maximum likelihood estimation for semiparametric transformation models with interval-censored data. *Biometrika* 103: 253–271.
4. van der Vaart AW, Wellner JA (1996). *Weak Convergence and Empirical Processes*. Springer, New York, NY.
5. Zeng D, Gao F, Lin DY (2017). Maximum likelihood estimation for semiparametric regression models with multivariate interval-censored data. *Biometrika* 104: 505–525.
6. Polack FP, Thomas SJ, Kitchin N, et al (2021). Safety and efficacy of the BNT162b2 mRNA Covid-19 vaccine. *New Eng J Med*, 383: 2603-15.