

554 **Supplementary Materials**

555 Additional file 1 — Stan code for the Bayesian multivariate hierarchical model

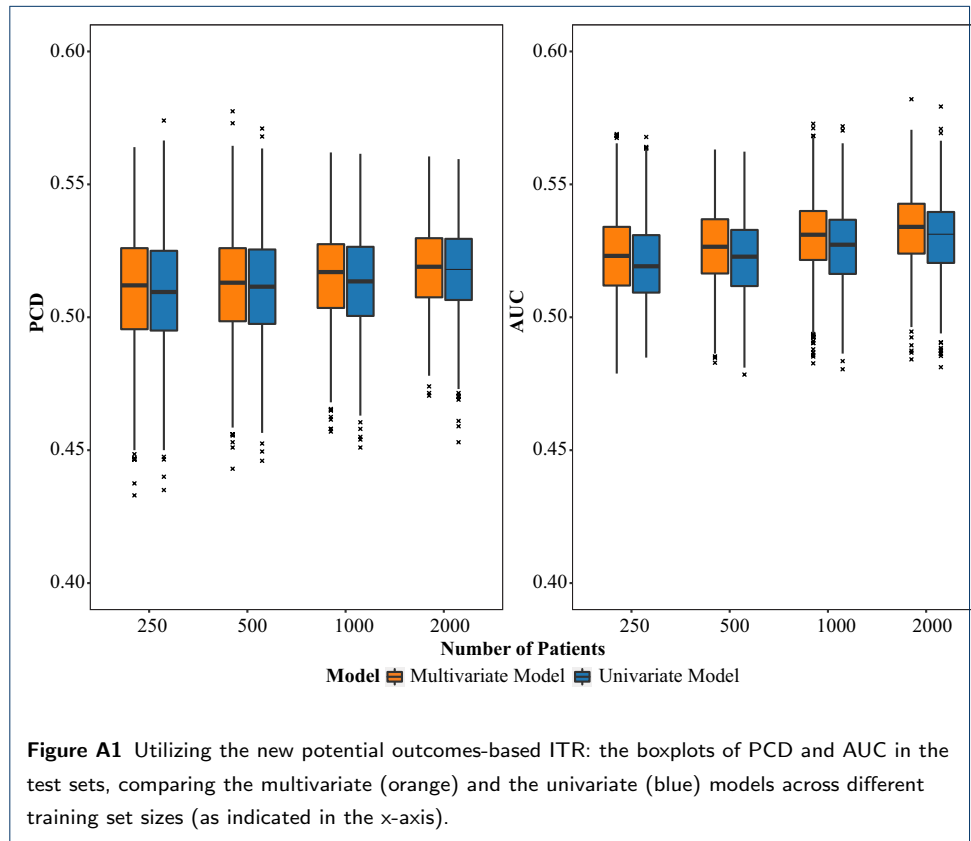
```

556 1 data {
557 2   int<lower=0> N;           // number of observations
558 3   int<lower=0> D;           // number of of binary outcomes
559 4   int<lower=2> L;           // number of WHO categories
560 5   int<lower=0> P_main;      // number of pre-treatment
561   characteristics in the main effects term
562 6   int<lower=0> P_inter;     // number of pre-treatment
563   characteristics in the interaction effects term
564 7   int<lower=1,upper=L> y_ord[N]; // vector of ordinal outcomes
565 8   int<lower=0,upper=1> y_b[N,D]; // matrix of D binary outcomes (N x D
566   matrix)
567 9   int<lower=0,upper=1> A[N]; // treatment or control
568 0   row_vector[P_main] x_main[N]; // pre-treatment characteristics in
569   the main effects term (N x P_main matrix)
570 1   row_vector[P_inter] x_inter[N]; // pre-treatment characteristics in
571   the interaction effects term (N x P_inter matrix)
572 2 }
573 3
574 4 parameters {
575 5   ordered[L-1] tau;         // cut-points for cumulative odds
576   model
577 6   vector<lower=0>[(P_inter + 1)] sigma_beta; // sd of outcome-specific
578   treatment main effect and interaction effect
579 7   vector[D] beta_0;        // outcome-specific intercepts for D
580   binary outcomes
581 8   matrix[P_main,(D + 1)] beta_1; // covariates main effect for D
582   binary outcomes and 1 ordinal outcome (P_main x (D + 1) matrix)
583 9   vector[(P_inter + 1)] beta_star; // pooled treatment main effect and
584   pooled interaction effect across all outcomes
585 0
586 1   // non-central parameterization
587 2
588 3   matrix[(P_inter + 1),(D + 1)] z_beta_int;
589 4 }
590 5
591 6 transformed parameters {
592 7   matrix[(P_inter + 1),(D + 1)] beta_int; // outcome-specific
593   treatment main effect and interaction effect ((P_inter + 1) x (D +
594   1) matrix)
595 8   vector[D] yhat_b[N];
596 9   real yhat_ord[N];
597 0
598 1   for (j in 1:(P_inter + 1))
599 2     for (k in 1:(D + 1)){
600 3       beta_int[j,k] = beta_star[j] + sigma_beta[j] * z_beta_int[j,k];
601 4     }
602 5
603 6   for (i in 1:N){

```

```
6047   for (k in 1:D){
6058     yhat_b[i,k] = beta_0[k] + x_main[i] * beta_1[,k] + (
606     append_col(1, x_inter[i]) * beta_int[,k]) * A[i];
6079   }
6080   yhat_ord[i] = x_main[i] * beta_1[,D+1] + (append_col(1, x_inter[
609   i]) * beta_int[,D+1]) * A[i];
6101 }
6112 }
6123
6134
6145 model {
6156
6167   // priors
6178
6189   sigma_beta ~ exponential(1);
6190   beta_star ~ normal(0,2.5);
6201   to_vector(beta_1) ~ normal(0,2.5);
6212   to_vector(z_beta_int) ~ std_normal();
6223
6234   for (l in 1:(L-1)){
6245     tau[l] ~ student_t(3,0,8);
6256   }
6267
6278   for (k in 1:D){
6289     beta_0[k] ~ student_t(3,0,8);
6290   }
6301
6312   // outcome model
6323
6334   for (i in 1:N){
6345     y_ord[i] ~ ordered_logistic(yhat_ord[i], tau);
6356     for (k in 1:D){
6367       y_b[i,k] ~ bernoulli_logit(yhat_b[i,k]);
6378     }
6389   }
6390 }
```

640 Additional file 2 — Main analysis: comparing the performance of the Bayesian multivariate and univariate models
 641 when the true ITR is determined by potential outcomes
 642 To implement this potential outcomes-based ITR, we first consider the patient characteristics \tilde{x}_i along with the true
 643 values of parameters from the data generation process. Next, we use the *simstudy* package [43] to generate
 644 potential primary ordinal outcomes for subjects receiving the control treatment ($y_{A=0}^{(1)}$) and the experimental
 645 treatment ($y_{A=1}^{(1)}$). The optimal ITR is derived from the indicator function $I(y_{A=1}^{(1)} < y_{A=0}^{(1)})$, which evaluates
 646 whether the experimental treatment outcome is better than the control treatment outcome.
 647 Utilizing this new potential outcomes-based ITR, the subsequent plot (Figure A1) illustrates the comparison of PCD
 and AUC values between the Bayesian multivariate and univariate models across varying training set sizes. In



648
 649 comparison to Figure 1, the improvement in prediction using the multivariate model is less remarkable. This can be
 650 attributed to the fact that generating potential outcomes based on probability inherently involves more randomness.
 651 The gain in estimation is relatively small compared to the magnitude of this randomness. Consequently, when
 652 considering prediction error, the improvement becomes less noticeable as it is overshadowed by the noise introduced
 653 by the randomness.

654 Additional file 3 — Sensitivity analysis: comparing the performance of the Bayesian multivariate and univariate
 655 models when the true ITR is determined by potential outcomes

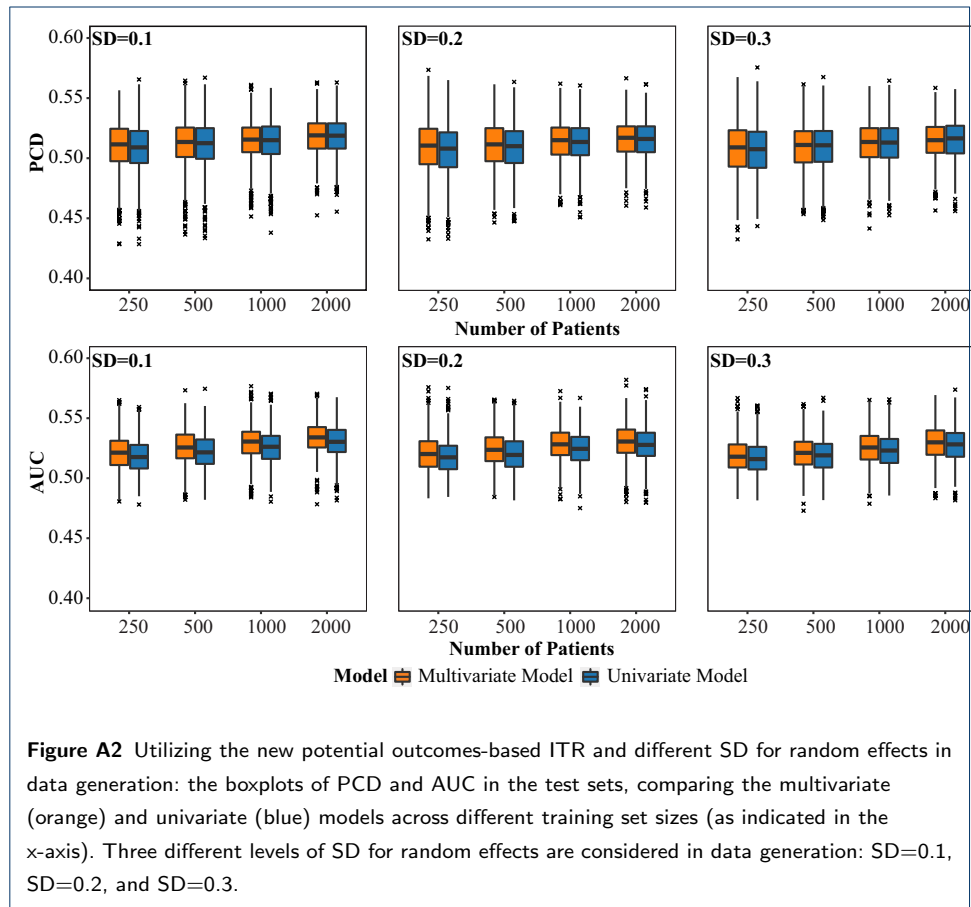


Figure A2 Utilizing the new potential outcomes-based ITR and different SD for random effects in data generation: the boxplots of PCD and AUC in the test sets, comparing the multivariate (orange) and univariate (blue) models across different training set sizes (as indicated in the x-axis). Three different levels of SD for random effects are considered in data generation: SD=0.1, SD=0.2, and SD=0.3.

656 Compared to Figure 2, the utilization of this new potential outcomes-based ITR yields less remarkable improvement
 657 in the multivariate model's performance. This could be due to the probabilistic nature of generating potential
 658 outcomes, which inherently involves more randomness. Despite the gain in estimation, the magnitude of this
 659 randomness is relatively large, resulting in a small improvement that is overshadowed by the introduced noise when
 660 considering prediction error.

661 Additional file 4 — The WHO 11-point COVID-19 clinical status scale.

| | |
|-----|--|
| 0: | Uninfected, no viral RNA detected |
| 1: | Asymptomatic, viral RNA detected |
| 2: | Symptomatic, independent |
| 3: | Symptomatic, assistance needed |
| 4: | Hospitalized, no oxygen therapy |
| 5: | Hospitalized, oxygen by mask or nasal prongs |
| 6: | Hospitalized, oxygen by non-invasive ventilation or high flow |
| 7: | Intubation & mechanical ventilation, $pO_2/FiO_2 \geq 150$ (or $SpO_2/FiO_2 \geq 200$) ^a |
| 8: | Mechanical ventilation, $pO_2/FiO_2 < 150$ (or $SpO_2/FiO_2 < 200$) or vasopressors |
| 9: | Mechanical ventilation, $pO_2/FiO_2 < 150$ and vasopressors, dialysis, or ECMO ^b |
| 10: | Dead |

^a pO_2 : partial pressure of oxygen, FiO_2 : fraction of inspired oxygen, SpO_2 : oxygen saturation.

^bECMO: extracorporeal membrane oxygenation.

Table A1 The WHO 11-point COVID-19 scale definition[39].

| | Control (n = 1097) | CCP (n = 1190) |
|----------|--------------------|----------------|
| WHO = 0 | 114 | 150 |
| WHO = 1 | 151 | 168 |
| WHO = 2 | 365 | 386 |
| WHO = 3 | 134 | 142 |
| WHO = 4 | 45 | 45 |
| WHO = 5 | 86 | 84 |
| WHO = 6 | 29 | 52 |
| WHO = 7 | 23 | 30 |
| WHO = 8 | 23 | 36 |
| WHO = 9 | 33 | 22 |
| WHO = 10 | 94 | 75 |

Table A2 The number of patients at different clinical stages of COVID-19 measured on the WHO 11-point scale at day 14 by treatment group.

662 Additional file 5 — Goodness-of-fit using posterior predictive checking
 663 In evaluating the suitability of a statistical model, it is crucial to determine whether the model provides an accurate
 664 representation of the observed data. This is particularly important for models like the *co* model, which is built upon
 665 a strong assumption of proportional cumulative odds. Posterior predictive checking serves as an effective way of
 666 assessing a model's *goodness-of-fit* [47, 48]. This method operates on the premise that a well-fitted model should
 667 enable the generation of replicated data (D^{rep}) that resembles the observed data ($D^{original}$) [49].
 668 The lack of fit can be measured by the Bayesian p-value, which represents the probability of the test statistic (e.g.,
 669 $P(Y \leq y), y = 0, \dots, 9$) for D^{rep} being equal to or exceeding the test statistic for $D^{original}$. A Bayesian p-value
 670 approaching zero or one signifies a potential issue with the model's fit, whereas a value near 0.5 suggests that the
 671 model captures the data well [42, 57]. We employed the procedure outlined in [41] to examine the *co* model's fit to
 672 the observed data and to compute the Bayesian p-value.
 673 Table A3 provides the results of posterior predictive checking based on ten test statistics (along with their 95% CIs)
 674 for both the multivariate model (6) and univariate model (7). Our analysis confirmed the satisfactory fit of both
 675 models to the data.

| Treatment | Control | | | CCP | | |
|---------------------------|-------------------|---------------------------|------------------|-------------------|---------------------------|------------------|
| Test quantity: % subjects | $T(D^{original})$ | 95% int. for $T(D^{rep})$ | Bayesian P value | $T(D^{original})$ | 95% int. for $T(D^{rep})$ | Bayesian P value |
| Multivariate model | | | | | | |
| WHO ≤ 0 | 10.39 | [8.39, 12.94] | 0.55 | 12.61 | [10.08, 14.71] | 0.39 |
| WHO ≤ 1 | 24.16 | [20.51, 26.80] | 0.34 | 26.72 | [23.28, 29.58] | 0.41 |
| WHO ≤ 2 | 57.43 | [52.42, 59.43] | 0.19 | 59.16 | [55.71, 62.52] | 0.48 |
| WHO ≤ 3 | 69.64 | [65.36, 71.93] | 0.29 | 71.09 | [68.24, 74.45] | 0.55 |
| WHO ≤ 4 | 73.75 | [69.92, 76.12] | 0.32 | 74.87 | [72.44, 78.32] | 0.63 |
| WHO ≤ 5 | 81.59 | [78.21, 83.68] | 0.34 | 81.93 | [80.17, 85.29] | 0.74 |
| WHO ≤ 6 | 84.23 | [82.22, 87.15] | 0.64 | 86.30 | [83.78, 88.40] | 0.44 |
| WHO ≤ 7 | 86.33 | [84.78, 89.43] | 0.74 | 88.82 | [86.13, 90.50] | 0.33 |
| WHO ≤ 8 | 88.42 | [87.69, 91.89] | 0.90 | 91.85 | [88.82, 92.69] | 0.14 |
| WHO ≤ 9 | 91.43 | [90.43, 94.07] | 0.81 | 93.70 | [91.26, 94.71] | 0.22 |
| Univariate model | | | | | | |
| WHO ≤ 0 | 10.39 | [8.57, 13.13] | 0.62 | 12.61 | [9.83, 14.54] | 0.32 |
| WHO ≤ 1 | 24.16 | [20.78, 27.07] | 0.41 | 26.72 | [23.03, 29.41] | 0.36 |
| WHO ≤ 2 | 57.43 | [52.42, 59.80] | 0.23 | 59.16 | [55.55, 62.44] | 0.46 |
| WHO ≤ 3 | 69.64 | [65.45, 72.20] | 0.31 | 71.09 | [68.15, 74.45] | 0.55 |
| WHO ≤ 4 | 73.75 | [69.83, 76.30] | 0.34 | 74.87 | [72.44, 78.40] | 0.62 |
| WHO ≤ 5 | 81.59 | [78.12, 83.87] | 0.35 | 81.93 | [80.17, 85.29] | 0.74 |
| WHO ≤ 6 | 84.23 | [82.04, 87.24] | 0.64 | 86.30 | [83.78, 88.49] | 0.45 |
| WHO ≤ 7 | 86.33 | [84.69, 89.52] | 0.74 | 88.82 | [86.13, 90.50] | 0.34 |
| WHO ≤ 8 | 88.42 | [87.60, 91.89] | 0.89 | 91.85 | [88.82, 92.69] | 0.14 |
| WHO ≤ 9 | 91.43 | [90.34, 94.17] | 0.81 | 93.70 | [91.34, 94.79] | 0.23 |

Table A3 Summary of posterior predictive checking based on the ten test statistics.