

COVID-19 Baseline Risk Score Analysis Report

MockENSEMBLE Study

USG COVID-19 Response Biostatistics Team

August 09, 2021

Contents

1	Baseline Risk Score (Proxy for SARS-CoV-2 Exposure)	9
2	Appendix	21

List of Tables

1.1	Variables considered for risk score analysis.	9
1.2	All learner-screen combinations (14 in total) used as input to the Superlearner.	11
1.3	Weights assigned by Superlearner.	15
1.4	Predictors in learners assigned weight > 0.0 by Superlearner. . .	16
1.5	Cases per treatment arm prior to risk score analysis.	19
1.6	Cases per treatment arm post risk score analysis.	19

List of Figures

1.1	Cross-validated AUC (95% CI) of algorithms for predicting COVID-19 disease status starting 7 days after Day 29.	12
1.2	CV-estimated predicted probabilities of COVID-19 disease 7 days after Day 29 by case/control status for top 2 learners, Super-Learner and Discrete SL.	13
1.3	ROC curves based off CV-estimated predicted probabilities for the top 2 learners, Superlearner and Discrete SL.	14
1.4	Superlearner predicted probabilities of COVID-19 disease in vaccinees 7 days after Day 29 by case/control status.	17
1.5	ROC curve based off Superlearner predicted probabilities in vaccinees.	18

MOCK

Chapter 1

Baseline Risk Score (Proxy for SARS-CoV-2 Exposure)

Table 1.1: Variables considered for risk score analysis.

Variable.Name	Definition	Total.missing.values	Comments
EthnicityHispanic	Indicator ethnicity = Hispanic (1 = Hispanic, 0 = complement)	0/19503 (0.0%)	NA
EthnicityNotreported	Indicator ethnicity = Not reported (1 = Not reported, 0 = complement)	0/19503 (0.0%)	NA
EthnicityUnknown	Indicator ethnicity = Unknown (1 = Unknown, 0 = complement)	0/19503 (0.0%)	NA
Black	Indicator race = Black (1=Black, 0=complement)	0/19503 (0.0%)	NA
Asian	Indicator race = Asian (1=Asian, 0=complement)	0/19503 (0.0%)	NA
NatAmer	Indicator race = American Indian or Alaska Native (1=NatAmer, 0=complement)	0/19503 (0.0%)	NA
Multiracial	Indicator race = Multiracial (1=Multiracial, 0=complement)	0/19503 (0.0%)	NA
Notreported	Indicator race = Not reported (1=Notreported, 0=complement)	0/19503 (0.0%)	NA
Unknown	Indicator race = unknown (1=Unknown, 0=complement)	0/19503 (0.0%)	NA
URMforsubcohortsampling	Indicator of under-represented minority (1=Yes, 0=No)	0/19503 (0.0%)	NA
HighRiskInd	Baseline covariate indicating ≥ 1 Co-existing conditions (1=yes, 0=no, NA=missing)	0/19503 (0.0%)	NA
Sex	Sex assigned at birth (1=female, 0=male/undifferentiated/unknown)	0/19503 (0.0%)	NA
Age	Age at enrollment in years (integer ≥ 18 , NA=missing). Note that the randomization strata included Age 18-59 vs. Age ≥ 60 .	0/19503 (0.0%)	NA

Table 1.1: Variables considered for risk score analysis. (*continued*)

Variable.Name	Definition	Total.missing.values
BMI	BMI at enrollment (Ordered categorical 1,2, 3, 4, NA=missing); 1 = Underweight BMI < 18.5; 2 = Normal BMI 18.5 to < 25; 3 = Overweight BMI 25 to < 30; 4 = Obese BMI >= 30	0/19503 (0.0%)
Country	Country of the study site of enrollment (0=United States, 1=Argentina,2=Brazil, 3=Chile,4=Columbia, 5=Mexico, 6=Peru, 7=South Africa)	0/19503 (0.0%)
HIVinfection	Indicator HIV infected at enrollment (1=infected, 0=not infected)	0/19503 (0.0%)
CalendarDateEnrollment	Date variable (used to control for calendar time trends in COVID incidence). Coded as number of days since first person enrolled until the ppt is enrolled.	0/19503 (0.0%)

Note:

Variables with more than 5% missing values were dropped from analysis; missing values for other variables were dropped. Indicator variables not meeting the threshold, such that under the null of not a risk factor there were less than 5% in the subgroup with value 1 (or 0), were dropped from analysis.

Table 1.2: All learner-screen combinations (14 in total) used as input to the Superlearner.

Learner	Screen*
SL.mean	all
SL.glm	all glmnet univar_logistic_pval highcor_random

Note:

*Screen details:

all: includes all variables

glmnet: includes variables with non-zero coefficients in the standard implementation of SL.glmnet that optimizes the lasso tuning parameter via cross-validation

univar_logistic_pval: Wald test 2-sided p-value in a logistic regression model < 0.10

highcor_random: if pairs of quantitative variables with Spearman rank correlation > 0.90 , select one of the variables at random

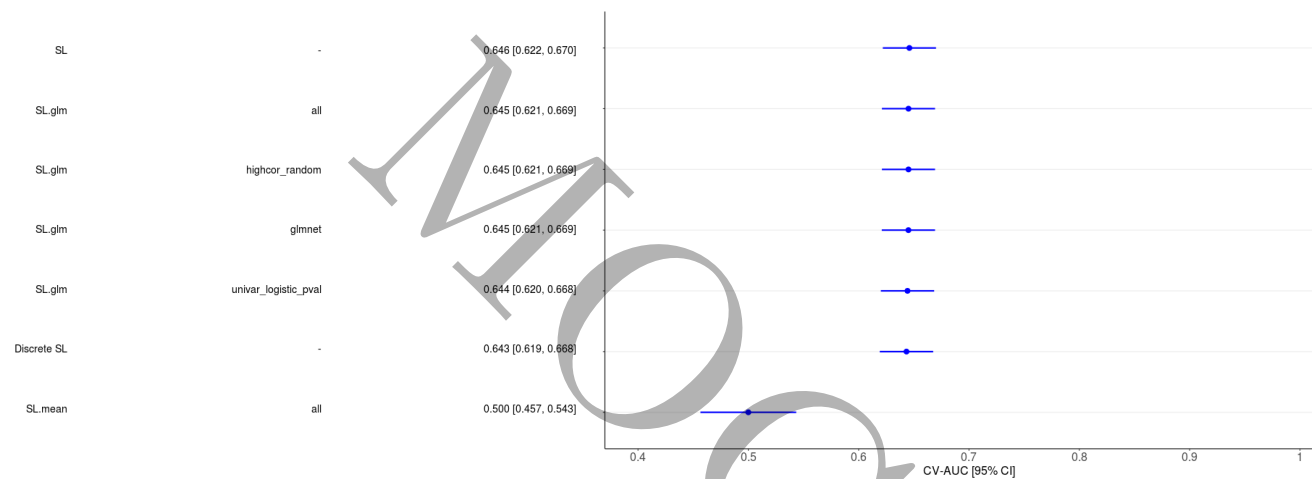


Figure 1.1: Cross-validated AUC (95% CI) of algorithms for predicting COVID-19 disease status starting 7 days after Day 29.

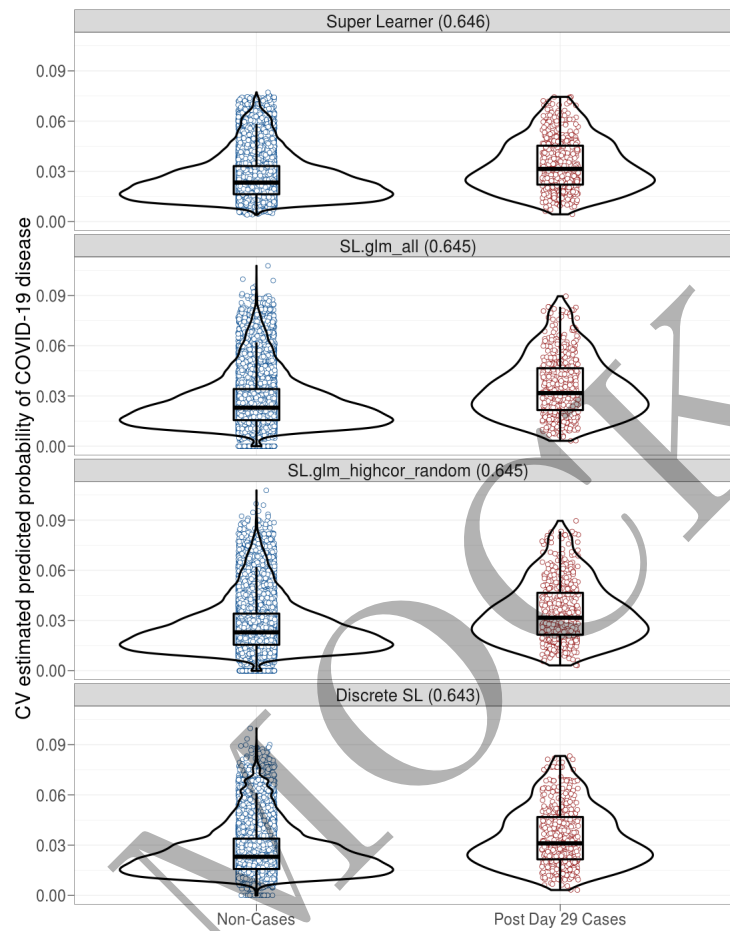


Figure 1.2: CV-estimated predicted probabilities of COVID-19 disease 7 days after Day 29 by case/control status for top 2 learners, SuperLearner and Discrete SL.

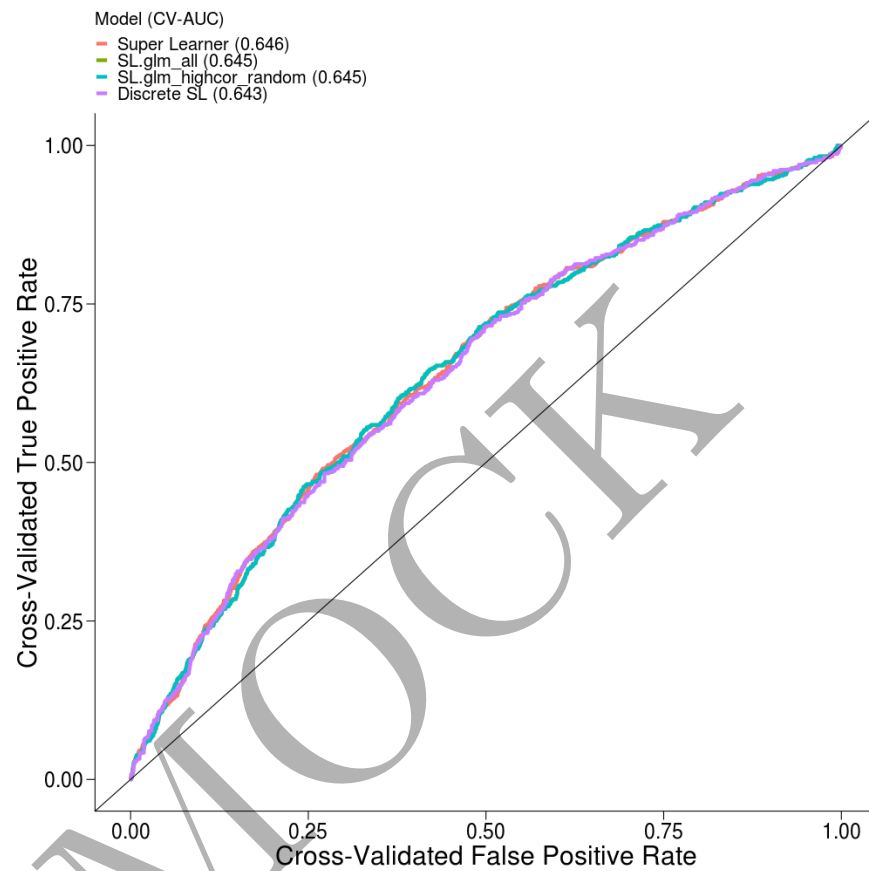


Figure 1.3: ROC curves based off CV-estimated predicted probabilities for the top 2 learners, Superlearner and Discrete SL.

Table 1.3: Weights assigned by Superlearner.

Learner	Screen	Weight
SL.glm	screen_univariate_logistic_pval	0.612
SL.glm	screen_all	0.338
SL.mean	screen_all	0.050
SL.glm	screen_glmnet	0.000
SL.glm	screen_highcor_random	0.000

Table 1.4: Predictors in learners assigned weight > 0.0 by Superlearner.

Learner	Screen	Weight	Predictors	Coefficient	Odds.Ratio
SL.glm	screen_univariate_logistic_pval	0.612	(Intercept)	-3.714	0.024
SL.glm	screen_univariate_logistic_pval	0.612	Black	-0.066	0.936
SL.glm	screen_univariate_logistic_pval	0.612	URMforsubcohortsampling	-0.046	0.955
SL.glm	screen_univariate_logistic_pval	0.612	HighRiskInd	0.394	1.483
SL.glm	screen_univariate_logistic_pval	0.612	Age	0.348	1.416
SL.glm	screen_all	0.338	(Intercept)	-3.787	0.023
SL.glm	screen_all	0.338	EthnicityHispanic	0.116	1.123
SL.glm	screen_all	0.338	EthnicityNotreported	0.022	1.022
SL.glm	screen_all	0.338	EthnicityUnknown	0.026	1.026
SL.glm	screen_all	0.338	Black	-0.016	0.985
SL.glm	screen_all	0.338	Asian	-0.074	0.929
SL.glm	screen_all	0.338	NatAmer	-0.005	0.995
SL.glm	screen_all	0.338	Multiracial	-0.055	0.947
SL.glm	screen_all	0.338	Notreported	-0.965	0.381
SL.glm	screen_all	0.338	Unknown	-0.053	0.948
SL.glm	screen_all	0.338	URMforsubcohortsampling	-0.151	0.859
SL.glm	screen_all	0.338	HighRiskInd	0.394	1.484
SL.glm	screen_all	0.338	Sex	-0.032	0.969
SL.glm	screen_all	0.338	Age	0.348	1.416
SL.glm	screen_all	0.338	BMI	-0.025	0.976
SL.glm	screen_all	0.338	Country	0.024	1.024
SL.glm	screen_all	0.338	HIVinfection	0.014	1.014
SL.glm	screen_all	0.338	CalendarDateEnrollment	0.046	1.047

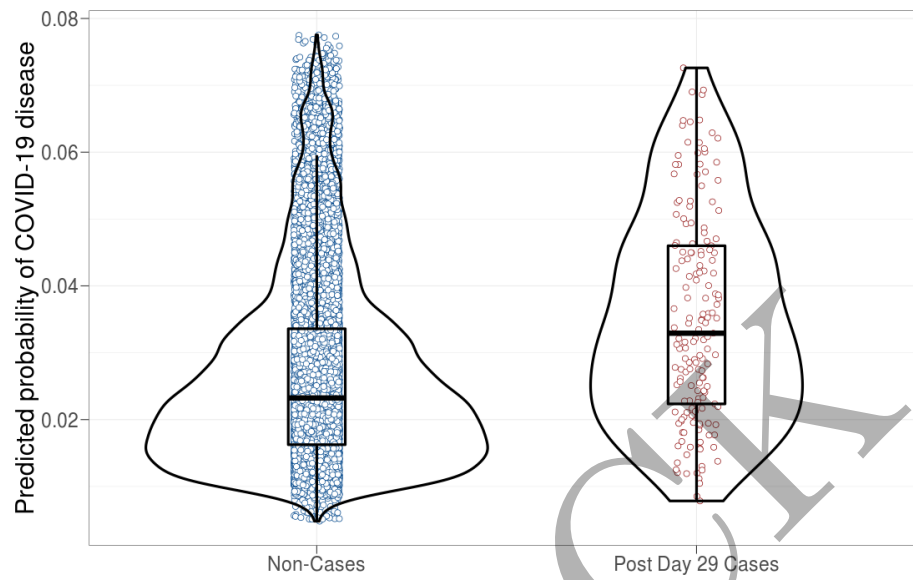


Figure 1.4: Superlearner predicted probabilities of COVID-19 disease in vaccinees 7 days after Day 29 by case/control status.

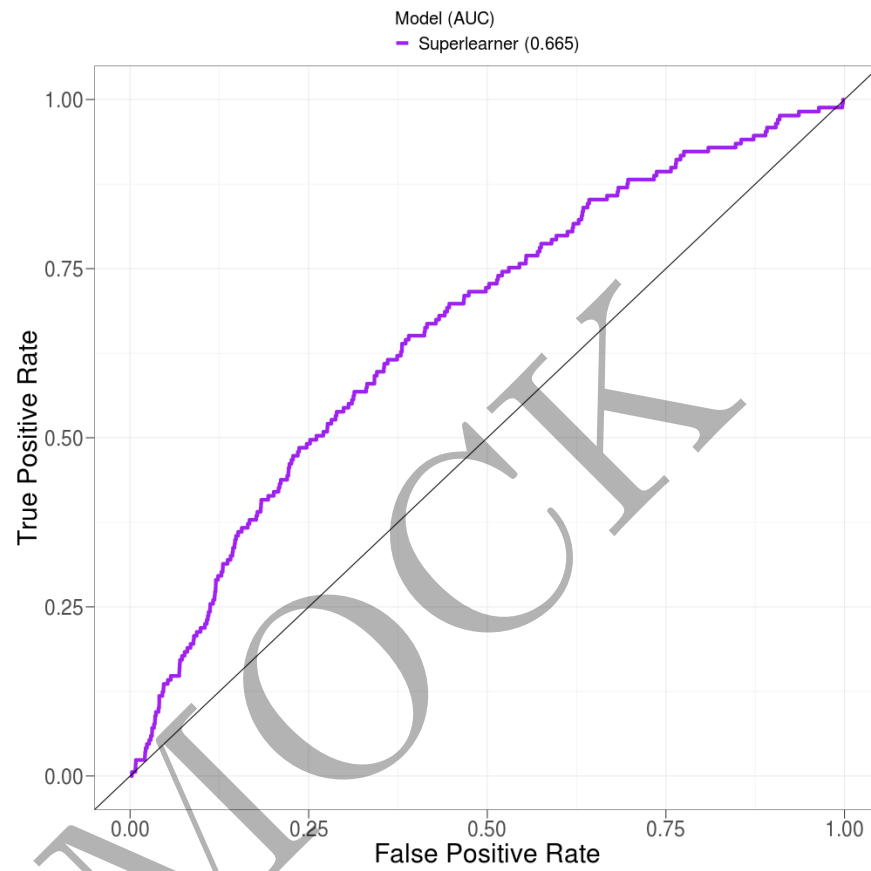


Figure 1.5: ROC curve based off Superlearner predicted probabilities in vaccinees.

Table 1.5: Cases per treatment arm prior to risk score analysis.

Study-Arm	Non-Cases	Post-Day 29-Cases
Placebo	18754	524
Vaccine	19196	169

Table 1.6: Cases per treatment arm post risk score analysis.

Study-Arm	Non-Cases	Post-Day 29-Cases
Placebo	18754	524
Vaccine	19196	169

MOCK

Chapter 2

Appendix

- This report was built from the [CoVPN/correlates_reporting](https://github.com/CoVPN/correlates_reporting) repository with commit hash b400af2c05e50f0d638fb3ba48a84b9adcecb5c. A diff of the changes introduced by that commit may be viewed at https://github.com/CoVPN/correlates_reporting/commit/b400af2c05e50f0d638fb3ba48a84b9adcecb5c
- The sha256 hash sum of the raw input file, “COVID_ENSEMBLE_practicedata.csv”:
0b430fcb0b10936460ae8fe7bfc3f78076afa07bbbca8b5e8ad9d7574d806934
- The sha256 hash sum of the processed file, “janssen_pooled_mock_data_processed.csv”:
6de79476ebd54c8c802b471c21898ae8a99edaf8fc55ba85bdd26e5b510eb056