

# COVID-19 Baseline Risk Score Analysis Report

## MockENSEMBLE Study

USG COVID-19 Response Biostatistics Team

July 23, 2021



# Contents

<b>1</b>	<b>Baseline Risk Score (Proxy for SARS-CoV-2 Exposure)</b>	<b>9</b>
<b>2</b>	<b>Appendix</b>	<b>19</b>



# 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



# 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 $\geq 1$ 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 $\geq 18$ , NA=missing). Note that the randomization strata included Age 18-59 vs. Age $\geq 60$ .	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

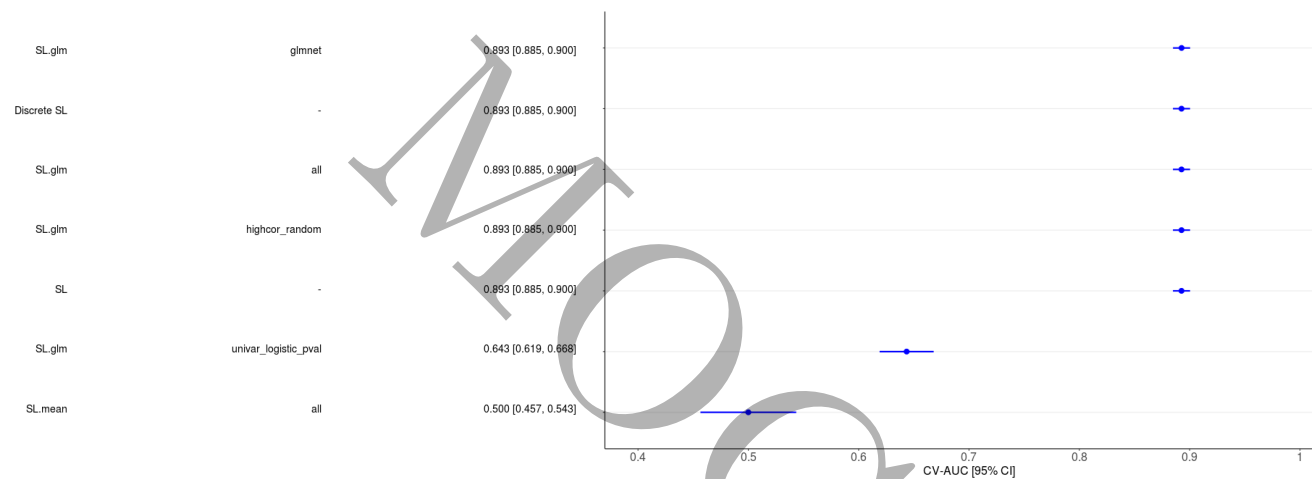


Figure 1.1: Cross-validated AUC (95% CI) of algorithms for predicting COVID-19 disease status starting 7 days after Day 29.

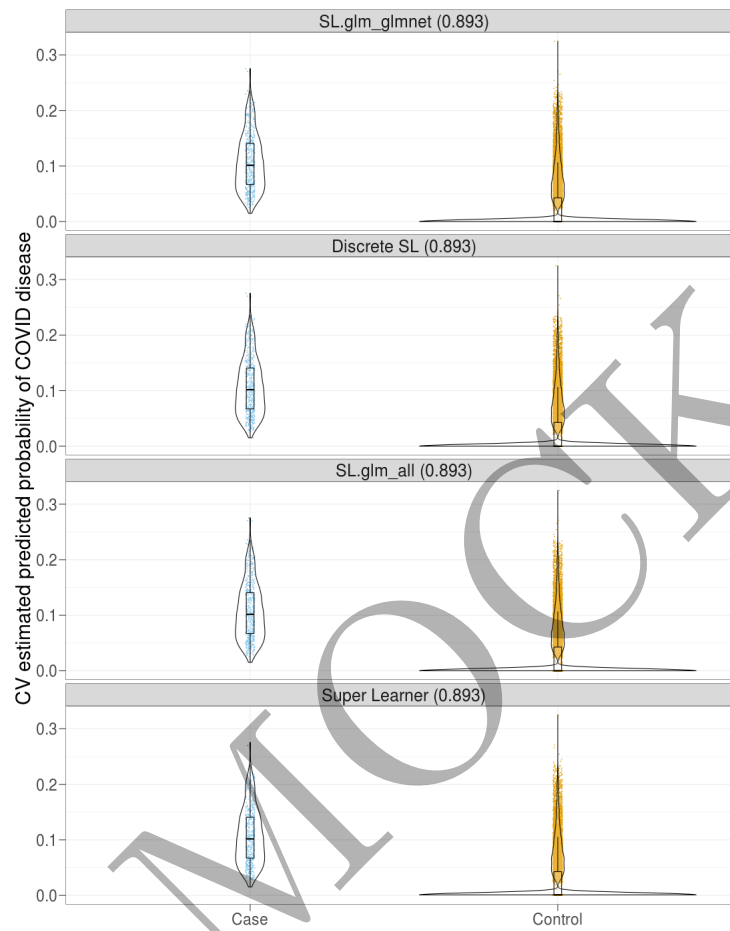


Figure 1.2: CV-estimated predicted probabilities of COVID-19 disease 7 days after Day 29 by case/control status for top 2 learners, SuperLearner and Discrete SL.

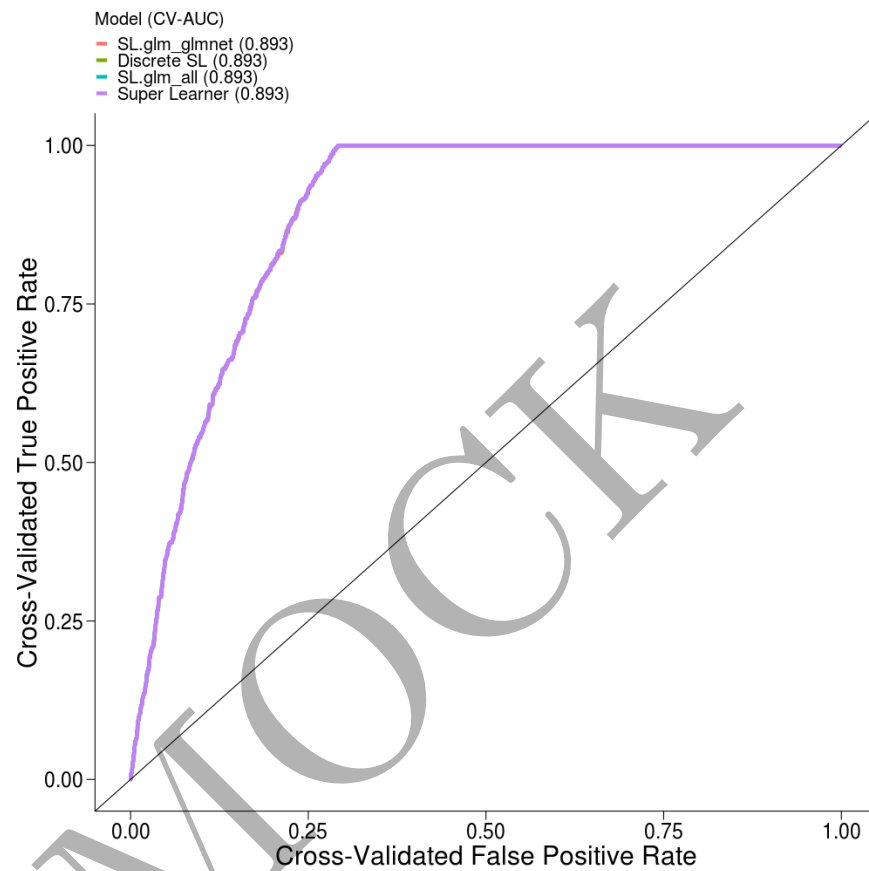


Figure 1.3: ROC curves based off CV-estimated predicted probabilities for the top 2 learners, Superlearner and Discrete SL.

Table 1.3: Weights assigned by Superlearner.

<b>Learner</b>	<b>Screen</b>	<b>Weight</b>
SL.glm	screen_glmnet	1
SL.mean	screen_all	0
SL.glm	screen_all	0
SL.glm	screen_univariate_logistic_pval	0
SL.glm	screen_highcor_random	0

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

Learner	Screen	Weight	Predictors	Coefficient	Odds.Ratio
SL.glm	screen_glmnet	1	(Intercept)	-19.348	0
SL.glm	screen_glmnet	1	EthnicityHispanic	0.078	1.081
SL.glm	screen_glmnet	1	EthnicityNotreported	0.013	1.014
SL.glm	screen_glmnet	1	EthnicityUnknown	0.051	1.052
SL.glm	screen_glmnet	1	Black	-0.062	0.94
SL.glm	screen_glmnet	1	Asian	-0.058	0.944
SL.glm	screen_glmnet	1	NatAmer	-0.033	0.967
SL.glm	screen_glmnet	1	Multiracial	-0.028	0.973
SL.glm	screen_glmnet	1	Notreported	-1.409	0.244
SL.glm	screen_glmnet	1	Unknown	-0.07	0.932
SL.glm	screen_glmnet	1	URMforsubcohortsampling	-0.066	0.936
SL.glm	screen_glmnet	1	HighRiskInd	0.383	1.467
SL.glm	screen_glmnet	1	Sex	12.043	169843.858
SL.glm	screen_glmnet	1	Age	0.355	1.427
SL.glm	screen_glmnet	1	BMI	-0.045	0.956
SL.glm	screen_glmnet	1	Country	-0.013	0.987
SL.glm	screen_glmnet	1	HIVinfection	0.034	1.034
SL.glm	screen_glmnet	1	CalendarDateEnrollment	0.023	1.023
SL.glm	screen_glmnet	1	(Intercept)	-19.348	0
SL.glm	screen_glmnet	1	EthnicityHispanic	0.078	1.081
SL.glm	screen_glmnet	1	EthnicityNotreported	0.013	1.014
SL.glm	screen_glmnet	1	EthnicityUnknown	0.051	1.052
SL.glm	screen_glmnet	1	Black	-0.062	0.94
SL.glm	screen_glmnet	1	Asian	-0.058	0.944
SL.glm	screen_glmnet	1	NatAmer	-0.033	0.967
SL.glm	screen_glmnet	1	Multiracial	-0.028	0.973
SL.glm	screen_glmnet	1	Notreported	-1.409	0.244
SL.glm	screen_glmnet	1	Unknown	-0.07	0.932
SL.glm	screen_glmnet	1	URMforsubcohortsampling	-0.066	0.936
SL.glm	screen_glmnet	1	HighRiskInd	0.383	1.467
SL.glm	screen_glmnet	1	Sex	12.043	169843.858
SL.glm	screen_glmnet	1	Age	0.355	1.427
SL.glm	screen_glmnet	1	BMI	-0.045	0.956
SL.glm	screen_glmnet	1	Country	-0.013	0.987
SL.glm	screen_glmnet	1	HIVinfection	0.034	1.034
SL.glm	screen_glmnet	1	CalendarDateEnrollment	0.023	1.023



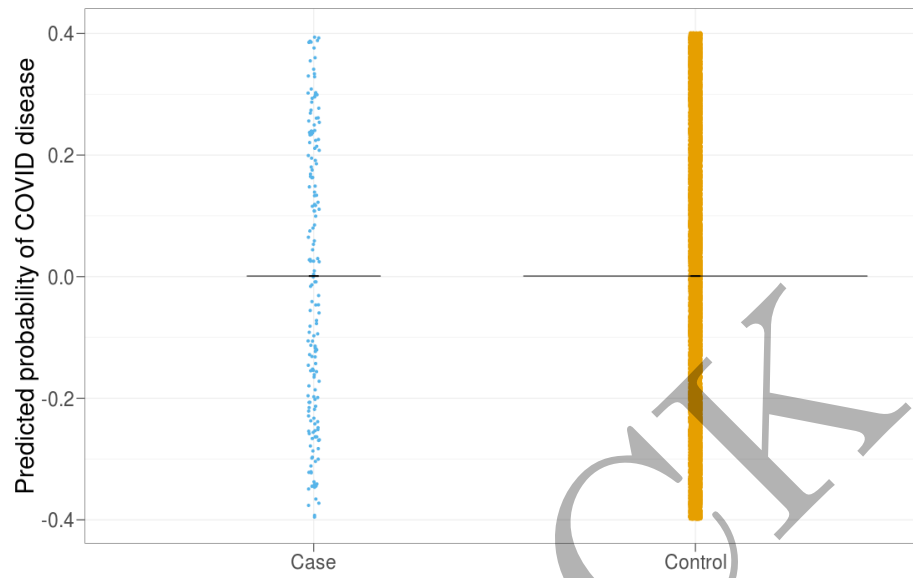


Figure 1.4: Superlearner predicted probabilities of COVID-19 disease in vaccinees 7 days after Day 29 by case/control status.

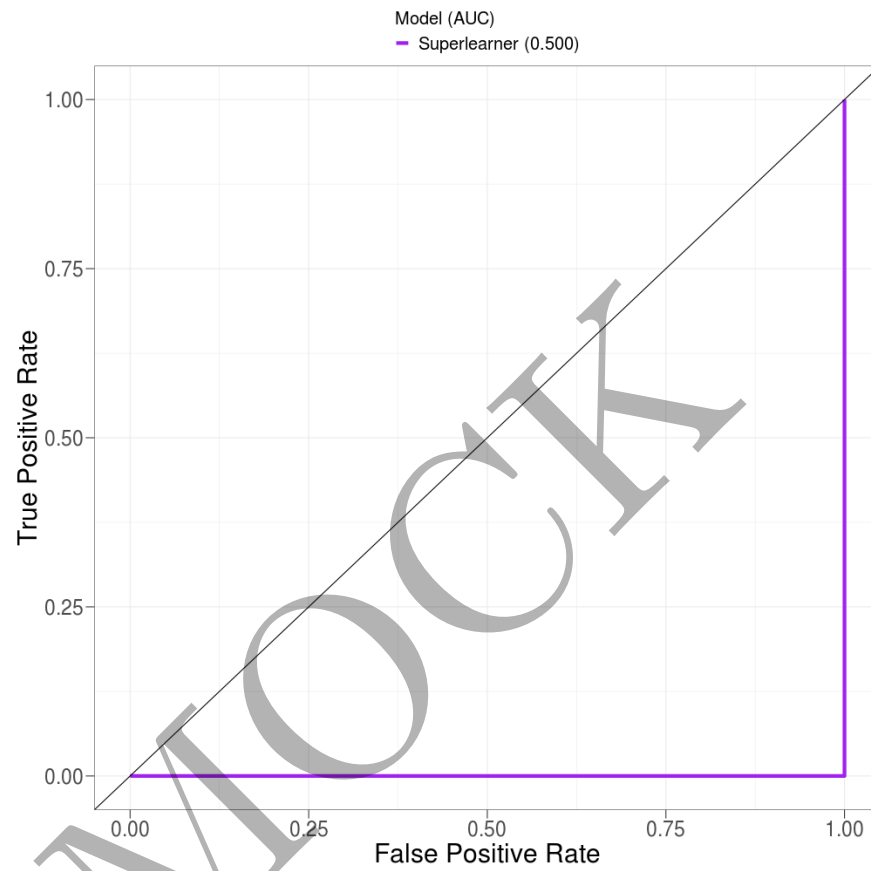


Figure 1.5: ROC curve based off Superlearner predicted probabilities in vaccinees.

## Chapter 2

# Appendix

- This report was built from the [CoVPN/correlates\\_reporting](https://github.com/CoVPN/correlates_reporting) repository with commit hash 5d679cb6d32c711f1bdb0e5950b332ad5d1395ce. A diff of the changes introduced by that commit may be viewed at [https://github.com/CoVPN/correlates\\_reporting/commit/5d679cb6d32c711f1bdb0e5950b332ad5d1395ce](https://github.com/CoVPN/correlates_reporting/commit/5d679cb6d32c711f1bdb0e5950b332ad5d1395ce)
- The sha256 hash sum of the raw input file, “COVID\_ENSEMBLE\_practicedata.csv”:  
e89657612e73e565302a8ce218bce5c13d956b89894135193efc82caf98bb26a
- The sha256 hash sum of the processed file, “janssen\_pooled\_mock\_data\_processed.csv”:  
19074a59ddd1179104f418b88b2da389c4ce7e91976e15eb4290f902b89a9315