COVID-19 Baseline Risk Score Analysis Report $_{\rm mock\ Study}$

USG COVID-19 Response Biostatistics Team

May 13, 2021

Contents

L	Baseline Risk Score (Proxy for SARS-CoV-2 Exposure)	9
2	Appendix	19

4 CONTENTS

List of Tables

1.1	Variables considered for risk score analysis	10
1.2	All learner-screen combinations (28 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 57	12
1.2	CV-estimated predicted probabilities of COVID-19 disease 7 days after Day 57 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 57 by case/control status	17
1.5	ROC curve based off Superlearner predicted probabilities in vaccinees	18

8 LIST OF FIGURES



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
MinorityInd	Baseline covariate underrepresented minority status (1=minority, 0=non-minority)	0/30000 (0.0%)	NA
EthnicityHispanic	Indicator ethnicity = Hispanic (0 = Non-Hispanic)	0/30000 (0.0%)	NA
${\bf Ethnicity Not reported}$	Indicator ethnicity = Not reported (0 = Non-Hispanic)	0/30000 (0.0%)	NA
EthnicityUnknown	Indicator ethnicity = Unknown (0 = Non-Hispanic)	0/30000 (0.0%)	NA
Black	Indicator race = Black (0 = White)	0/30000 (0.0%)	NA
Asian	Indicator race $=$ Asian $(0 = White)$	0/30000 (0.0%)	NA
NatAmer	Indicator race = American Indian or Alaska Native (0 = White)	0/30000 (0.0%)	NA
PacIsl	Indicator race = Native Hawaiian or Other Pacific Islander (0 = White)	0/30000 (0.0%)	NA
Multiracial	Indicator race = Multiracial (0 = White)	0/30000 (0.0%)	NA
Other	Indicator race = Other (0 = White)	0/30000 (0.0%)	NA
Notreported	Indicator race = Not reported (0 = White)	0/30000 (0.0%)	NA
Unknown	Indicator race = unknown (0 = White)	0/30000 (0.0%)	NA
HighRiskInd	Baseline covariate high risk pre-existing condition (1=yes, 0=no)	0/30000 (0.0%)	NA
Sex	Sex assigned at birth (1=female, 0=male)	0/30000 (0.0%)	NA
Age	Age at enrollment in years, between 18 and 85	0/30000 (0.0%)	NA
BMI	BMI at enrollment (kg/m^2)	0/30000 (0.0%)	NA

Table 1.2: All learner-screen combinations (28 in total) used as input to the Superlearner.

Learner	Screen*
SL.mean	all
SL.glm	all glmnet univar_logistic_pval highcor_random
SL.glm.interaction	glmnet univar_logistic_pval highcor_random
SL.glmnet	all
SL.gam	glmnet univar_logistic_pval highcor_random
SL.xgboost	all
SL.ranger.imp	all

Note:

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 ${\bf r}$

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

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

^{*}Screen details:



Figure 1.1: Cross-validated AUC (95% CI) of algorithms for predicting COVID-19 disease status starting 7 days after Day 57.



Figure 1.2: CV-estimated predicted probabilities of COVID-19 disease 7 days after Day 57 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_all	0.714
SL.glm.interaction	$screen_univariate_logistic_pval$	0.271
SL.mean	screen_all	0.016
SL.glmnet	screen_all	0.000
SL.xgboost	screen_all	0.000
SL.ranger.imp	screen_all	0.000
SL.glm	$screen_glmnet$	0.000
SL.glm	screen_univariate_logistic_pval_	0.000
SL.glm	screen_highcor_random	0.000
SL.glm.interaction	screen_glmnet	0.000
SL.glm.interaction	screen_highcor_random	0.000
SL.gam	screen_glmnet	0.000
SL.gam	screen_univariate_logistic_pval	0.000
SL.gam	screen_highcor_random	0.000

Weight Predictors Coefficient Odds.Ratio Learner Screen SL.glm screen all 0.714(Intercept) -3.592 0.028 SL.glm screen_all 0.714MinorityInd 0.109 1.115 EthnicityHispanic 0.016 SL.glm screen_all 0.7141.017 $_{\mathrm{SL.glm}}$ $screen_all$ 0.714 ${\bf Ethnicity Not reported}$ 0.064 1.066 0.714EthnicityUnknown 0.055 1.057 SL.glm screen_all $_{\mathrm{SL.glm}}$ screen_all 0.714Black -0.1270.881 SL.glmscreen_all 0.714 Asian -0.051 0.950SL.glm $screen_all$ 0.714NatAmer-0.022 0.978SL.glm $screen_all$ 0.714PacIsl 0.014 1.014 SL.glm $screen_all$ 0.714Multiracial -0.106 0.900Other SL.glm $screen_all$ 0.7140.049 1.051 $_{\mathrm{SL.glm}}$ 0.714Notreported 0.0641.066 $screen_all$ SL.glm 0.714 Unknown 0.056 1.058 screen all $_{\mathrm{SL.glm}}$ ${\tt screen_all}$ 0.714 HighRiskInd 0.7932.210SL.glm screen_all 0.714-0.033 0.967 Sex 0.714 0.743SL.glm screen_all Age 2.102 SL.glm $screen_all$ 0.714вми -0.001 0.9990.271(Intercept) 0.027 SL.glm.interaction $screen_univariate_logistic_pval$ -3.604SL.glm.interaction ${\tt screen_univariate_logistic_pval}$ 0.271PacIsl 0.060 1.062 SL.glm.interaction screen_univariate_logistic_pval 0.271Other 0.112 1.119 SL.glm.interaction $screen_univariate_logistic_pval$ 0.271HighRiskInd 0.806 2.239 SL.glm.interaction screen_univariate_logistic_pval 0.271Age 0.7672.154 ${\tt screen_univariate_logistic_pval}$ PacIsl:Other SL.glm.interaction 0.271NA SL.glm.interaction screen_univariate_logistic_pval 0.271PacIsl:HighRiskInd 0.020 1.020 SL.glm.interaction PacIsl:Age -0.083 $screen_univariate_logistic_pval$ 0.2710.921Other: HighRiskIndSL.glm.interaction $screen_univariate_logistic_pval$ 0.271-0.015 0.985SL.glm.interaction screen_univariate_logistic_pval 0.271Other:Age -0.050 0.952HighRiskInd:Age ${\it SL.glm.interaction}$ $screen_univariate_logistic_pval$ 0.271-0.0210.979 SL.glm.interaction 0.271(Intercept) 0.027 screen_univariate_logistic_pval -3.604 $screen_univariate_logistic_pval$ 1.062 SL.glm.interaction 0.271PacIsl 0.060 SL.glm.interaction screen_univariate_logistic_pval 0.271Other 0.1121.119 $screen_univariate_logistic_pval$ 0.271HighRiskInd 0.806 2.239 SL.glm.interaction SL.glm.interaction $screen_univariate_logistic_pval$ 0.2710.767 2.154SL.glm.interaction screen_univariate_logistic_pval 0.271PacIsl:Other NASL.glm.interaction $screen_univariate_logistic_pval$ 0.271PacIsl:HighRiskInd 0.020 1.020 SL.glm.interaction ${\tt screen_univariate_logistic_pval}$ 0.271PacIsl:Age -0.083 0.921 Other:HighRiskInd $screen_univariate_logistic_pval$ -0.015 0.985 SL.glm.interaction 0.271SL.glm.interaction $screen_univariate_logistic_pval$ 0.271Other:Age -0.050 0.952SL.glm.interaction $screen_univariate_logistic_pval$ 0.271HighRiskInd:Age -0.021 0.979

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



Figure 1.4: Superlearner predicted probabilities of COVID-19 disease in vaccinees 7 days after Day 57 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 repository with commit hash dd6777c742d1d1c8a4bc6dd4b107f7ec7254fcbd. A diff of the changes introduced by that commit may be viewed at https://github.com/CoVPN/correlates_reporting/commit/dd6777c742d1d1c8a4bc6dd4b107f7ec7254fcbd
- The sha256 hash sum of the raw input file, "COVID_VEtrial_practicedata_primarystage1.csv": 2353971c2e14399ede55ef6ba0d4e624626433dc15ec507c2482bb886210019a
- \bullet The sha256 hash sum of the processed file, "practice_data.csv": 6250066f886245b78f7aa29fefc615ba5d10118448f298c39ec2b601b2a5049f