

## Supplementary Material

### **Developing a COVID-19 mortality risk prediction model when individual level data is not available**

Noam Barda<sup>1,2,3</sup>, Dan Riesel<sup>1</sup>, Amichay Akriv<sup>1</sup>, Joseph Levy<sup>1</sup>, Uriah Finkel<sup>1</sup>, Gal Yona<sup>4</sup>, Daniel Greenfeld<sup>5</sup>, Shimon Sheiba<sup>5</sup>, Jonathan Somer<sup>5</sup>, Eitan Bachmat<sup>3</sup>, Guy N. Rothblum<sup>4</sup>, Uri Shalit<sup>5</sup>, Doron Netzer<sup>6</sup>, Ran Balicer<sup>1,2</sup>, Noa Dagan<sup>1,2,3</sup>

### **Contents**

Supplementary Table 1: Population table for baseline model .....	2
Supplementary Table 2: Corrections table for the "multi-calibration" recalibration procedure .....	4
Supplementary Table 3: Sensitivity analysis results for the discrimination performance considering the composite outcome (severe cases and mortality) ....	5
Supplementary Table 4: Variable and outcome definitions .....	6
Supplementary Figure 1: Odds ratio as a function of variable values for the baseline predictor .....	8
Supplementary Figure 2: Cumulative Distribution Function of time-to-death among CHS COVID-19 patients .....	9
Supplementary Methods 1: Transformation of SHAP values to odds ratios.....	10
Supplementary Methods 2: Pseudo-code of the "multi-calibration" algorithm.....	11
Supplementary Methods 3: Technical description of adjustment procedure.....	12
References .....	14

**Supplementary Table 1: Population table for baseline model**

Variable	Train/Validation	Test	Missing (%)
<b>Overall</b>	735,000	315,000	
<b>Outcome, n (%)</b>			0.0
No	726,749 (98.9)	311,533 (98.9)	
Yes	8,251 (1.1)	3,467 (1.1)	
<b>Age, Mean (SD)</b>	41.00 (21.26)	40.99 (21.23)	0.0
<b>Sex, n (%)</b>			0.0
Female	378,617 (51.5)	162,144 (51.5)	
Male	356,383 (48.5)	152,856 (48.5)	
<b>Pack years, Mean (SD)</b>	23.65 (24.07)	23.52 (23.31)	75.5
<b>Chronic Respiratory Disease, n (%)</b>			0.0
No	669,347 (91.1)	286,961 (91.1)	
Yes	65,653 (8.9)	28,039 (8.9)	
<b>Cardiovascular Disease, n (%)</b>			0.0
No	657,985 (89.5)	281,953 (89.5)	
Yes	77,015 (10.5)	33,047 (10.5)	
<b>Diabetes, n (%)</b>			0.0
No	655,383 (89.2)	280,804 (89.1)	
Yes	79,617 (10.8)	34,196 (10.9)	
<b>Malignancy, n (%)</b>			0.0
No	694,706 (94.5)	297,718 (94.5)	
Yes	40,294 (5.5)	17,282 (5.5)	
<b>Hypertension, n (%)</b>			0.0
No	612,797 (83.4)	262,802 (83.4)	
Yes	122,203 (16.6)	52,198 (16.6)	
<b>COPD, n (%)</b>			0.0
No	721252 (98.1)	309082 (98.1)	
Yes	13748 (1.9)	5918 (1.9)	
<b>Wheezing/Dyspnea diagnosis, Mean (SD)</b>	0.05 (0.47)	0.05 (0.49)	0.0
<b>Albumin, Mean (SD)</b>	4.24 (0.38)	4.24 (0.38)	64.1
<b>Red cell distribution width, Mean (SD)</b>	13.69 (1.32)	13.69 (1.31)	49.9
<b>C-Reactive Peptide, Mean (SD)</b>	1.16 (2.80)	1.14 (2.73)	86.6
<b>Urea, Mean (SD)</b>	31.49 (14.34)	31.38 (14.37)	50.4
<b>Lymphocyte%, Mean (SD)</b>	31.12 (9.07)	31.16 (9.09)	43.7
<b>Chloride, Mean (SD)</b>	103.94 (3.45)	103.95 (3.46)	94.8

<b>Creatinine, Mean (SD)</b>	0.79 (0.41)	0.79 (0.41)	46.6
<b>High Density Lipoprotein, Mean (SD)</b>	49.08 (13.07)	49.04 (13.04)	51.7
<b>Duration of hospitalizations, Mean (SD)</b>	0.55 (5.72)	0.54 (5.86)	0.0
<b>Count of hospitalizations, Mean (SD)</b>	0.10 (0.45)	0.10 (0.45)	0.0
<b>Count of ambulance rides, Mean (SD)</b>	0.03 (0.26)	0.03 (0.27)	0.0
<b>Count of Sulfonamide dispenses, Mean (SD)</b>	0.15 (1.22)	0.15 (1.21)	0.0
<b>Count of Anti-cholinergic dispenses, Mean (SD)</b>	0.09 (0.91)	0.09 (0.91)	0.0
<b>Count of Glucocorticoid dispenses, Mean (SD)</b>	0.16 (0.94)	0.16 (0.93)	0.0

Population characteristics for the three population sets used to build and validate the baseline model.

Abbreviations: SD, Standard Deviation;

**Supplementary Table 2: Corrections table for the "multi-calibration" recalibration procedure**

<b>Total Corrections</b>	23	
	<b>Sex</b>	
<b>Age Group</b>	<b>Male</b>	<b>Female</b>
10-19	+0.003	-0.011214
20-29	+0.00368	-0.010539
30-39	+0.00086	-0.013356
40-49	+0.00723	-0.006987
50-59	+0.01217	-0.002049
60-69	+0.02944	+0.01523
70-79	+0.05912	+0.0449
80+	+0.08484	+0.07063

Total number of corrections performed and cumulative magnitude of the correction for each subgroup of age and sex.

**Supplementary Table 3: Sensitivity analysis results for the discrimination performance considering the composite outcome (severe cases and mortality)**

<b>Model</b>	<b>Metric</b>	<b>Value</b>
COVID-19 Model	AUROC	0.901, 95% CI: 0.881-0.920
COVID-10 Model 10% risk cut-off	Percent Positive	8%, 95% CI: 7%-9%
	Sensitivity	53%, 95% CI: 47%-60%
	PPV	37%, 95% CI: 32%-43%
COVID-10 Model 5% risk cut-off	Percent Positive	15%, 95% CI: 14-16%
	Sensitivity	73%, 95% CI: 67%-79%
	PPV	27%, 95% CI: 24%-31%

Abbreviations: AUROC, Area Under the Receiver Operating Characteristic curve; PPV, Positive Predictive Value

**Supplementary Table 4: Variable and outcome definitions**

Variable Definitions	Category	Units	Time Frame (prior to index date)	Details	Mechanism of feature selection	Potential for missing values
Outcome	Diagnoses	No / Yes	NA	ICD9 codes 48[0-8]*, 46[0-6]*, 490*, 510.9*, 511.0, 511.89, 518.82, 785.52 or positive influenza PCR	NA	Always available (cannot be missing)
Age	Demographics	Years	Current	Age in full years	From the top-features list	Always available (cannot be missing)
Sex	Demographics	Male / Female	Current		Added manually	Always available (cannot be missing)
Pack years	Clinical Covariates	Count	Last recorded value in last 1 year	#years smoked X average number of packs per year	From the top-features list	Cannot be missing (if nonsmoker then 0)
COPD	Diagnoses	No / Yes	Ever	As per Clalit's chronic disease registry	From the top-features list	Cannot be missing (if no documentation is available then the variable is set to "No")
Number of wheezing / dyspnea diagnoses	Diagnoses	Count	Cumulative count over 1 year	ICD9: 786.0	From the top-features list	Cannot be missing (if there is no documentation of a relevant diagnosis the variable is set to 0)
Albumin	Labs	g/dl	Last recorded value in last 1 year		From the top-features list	Missing if the lab test was not done in the last year.
Red cell distribution width	Labs	%	Last recorded value in last 1 year		From the top-features list	Missing if the lab test was not done in the last year.
C-Reactive Peptide	Labs	mcg/ml	Last recorded value in last 1 year		From the top-features list	Missing if the lab test was not done in the last year.
Urea	Labs	mg/dl	Last recorded value in last 1 year		From the top-features list	Missing if the lab test was not done in the last year.
Lymphocyte	Labs	%	Last recorded value in last 1 year		From the top-features list	Missing if the lab test was not done in the last year.
Chloride	Labs	mEq/L	Last recorded value in last 1 year		From the top-features list	Missing if the lab test was not done in the last year.
Creatinine	Labs	mg/dl	Last recorded value in last 1 year		From the top-features list	Missing if the lab test was not done in the last year.
High Density Lipoprotein	Labs	mg/dl	Last recorded value in last 1 year		From the top-features list	Missing if the lab test was not done in the last year.

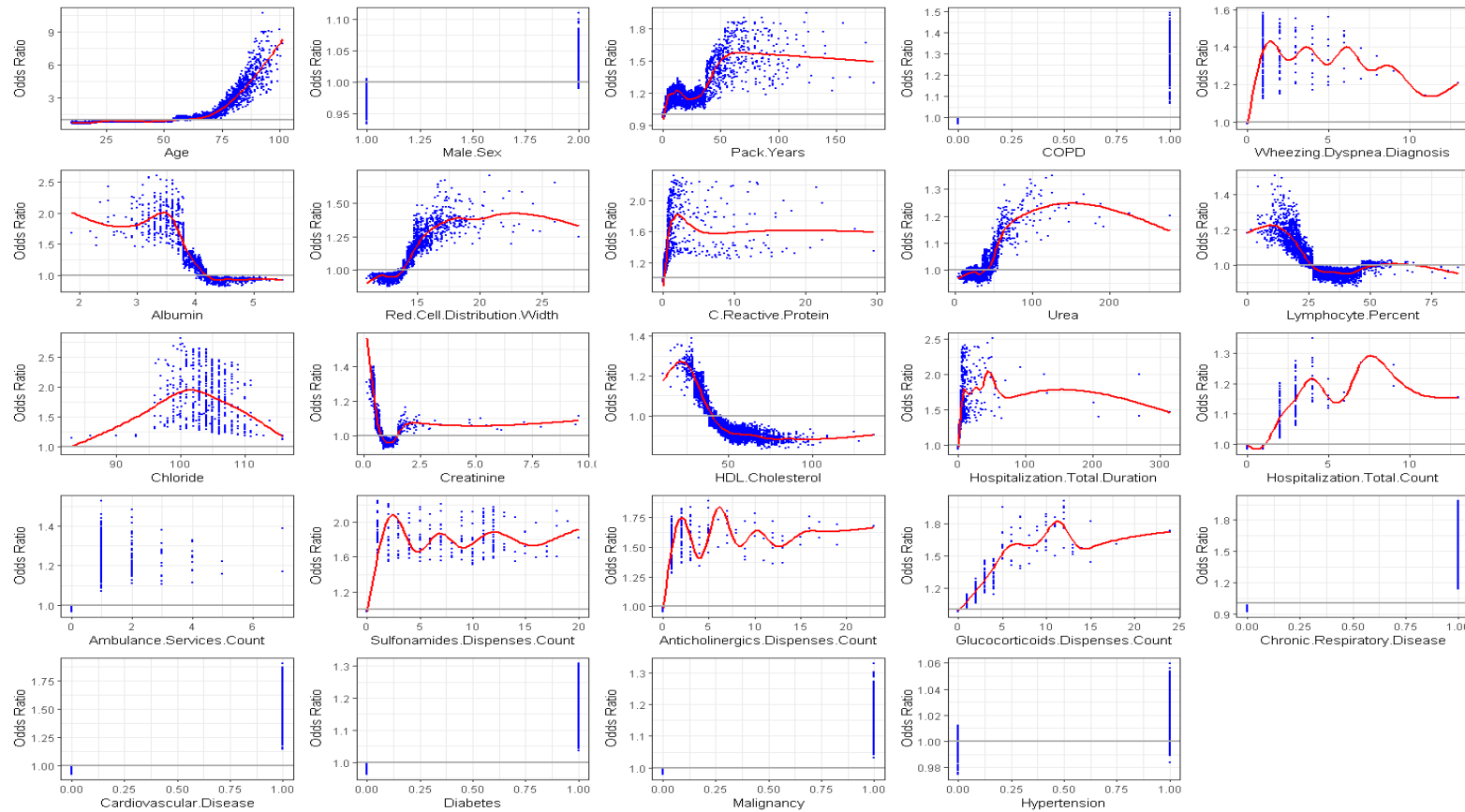
Duration of hospitalizations	Healthcare Utilization	Days	Cumulative count over 1 year		From the top-features list	Cannot be missing (if there were no hospitalizations then the variable is set to 0)
Count of hospitalizations	Healthcare Utilization	Count	Cumulative count over 1 year		From the top-features list	Cannot be missing (if there were no hospitalizations then the variable is set to 0)
Count of ambulance rides	Healthcare Utilization	Count	Cumulative count over 1 year		From the top-features list	Cannot be missing (if there were no ambulance rides then the variable is set to 0)
Count of Sulfonamide dispenses	Drugs	Count	Cumulative count over 1 year	ATC4: C03CA	From the top-features list	Cannot be missing (if there were no drug dispenses then the variable is set to 0)
Count of Anti-cholinergic dispenses	Drugs	Count	Cumulative count over 1 year	ATC4: R03BB	From the top-features list	Cannot be missing (if there were no drug dispenses then the variable is set to 0)
Count of Glucocorticoid dispenses	Drugs	Count	Cumulative count over 1 year	ATC4: H02AB	From the top-features list	Cannot be missing (if there were no drug dispenses then the variable is set to 0)
Chronic Respiratory Disease	Diagnoses	No / Yes	Ever	As per Clalit's chronic disease registry	Added manually	Cannot be missing (if no documentation is available then the variable is set to "No")
Cardiovascular Disease	Diagnoses	No / Yes	Ever	As per Clalit's chronic disease registry	Added manually	Cannot be missing (if no documentation is available then the variable is set to "No")
Diabetes	Diagnoses	No / Yes	Ever	As per Clalit's chronic disease registry	Added manually	Cannot be missing (if no documentation is available then the variable is set to "No")
Malignancy	Diagnoses	No / Yes	Ever	As per Clalit's chronic disease registry	Added manually	Cannot be missing (if no documentation is available then the variable is set to "No")
Hypertension	Diagnoses	No / Yes	Ever	As per Clalit's chronic disease registry	Added manually	Cannot be missing (if no documentation is available then the variable is set to "No")

Caption: A list of the variables used in the model, including their type, units, time frame of extraction, definitions, how they were selected and potential for missingness.

All variables were extracted at or before February 1st, 2020, prior to the onset of the COVID-19 pandemic in Israel.

Abbreviations: COPD, Chronic Obstructive Pulmonary Disease; g, gram; mcg, microgram; mEq, mili-equivalent; ml, milliliter; dl, deciliter; L, liter; PCR, Polymerase Chain Reaction; ICD9, International Classification of Disease, 9<sup>th</sup> revision; ATC4, Anatomical Therapeutic Chemical (ATC) Classification System, 4<sup>th</sup> level;

## Supplementary Figure 1: Odds ratio as a function of variable values for the baseline predictor



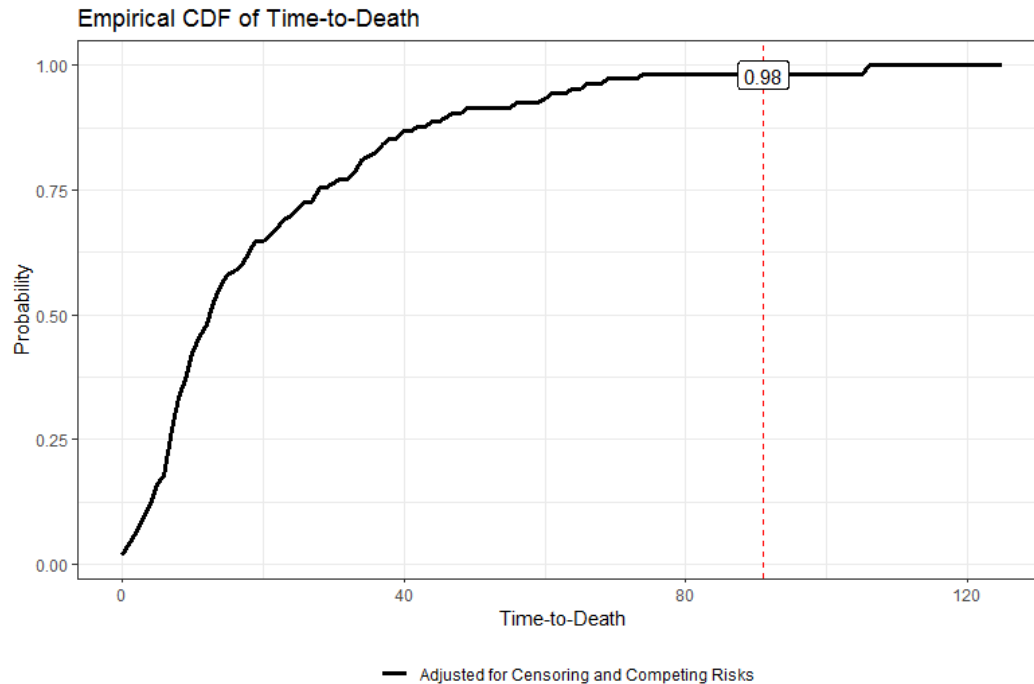
A plot of the odds ratios for different values of the predictors in the baseline model, using SHAP values<sup>1</sup>. A smoothed red line is fit to the curves and a horizontal grey line is drawn at odds ratio = 1.

All figures use a random sample of patients from this same population, n=10,000 unique patients.

Abbreviations: SHAP, SHapley Additive exPlanations; COPD, Chronic Obstructive Pulmonary Disease; HDL, High Density Lipoprotein;



## Supplementary Figure 2: Cumulative Distribution Function of time-to-death among CHS COVID-19 patients



Empirical cumulative distribution function for time-to-death of all COVID-19 patients in CHS' database. The black line shows the cumulative incidence as calculated after accounting for censoring and the "competing risk" of cure, derived using the Aalen-Johansen estimator<sup>2</sup>. The vertical dashed line denotes the follow-up time used in the study, 91 days.

The figure uses all COVID-19 patients among the Clalit Health Services insured population until the extraction date (July 16<sup>th</sup>, 2020), n=16,049 unique patients.

Abbreviations: CDF, Cumulative Distribution Function; CHS, Clalit Health Services; COVID-19, Corona Virus Disease 2019;

## Supplementary Methods 1: Transformation of SHAP values to odds ratios

In Shapley analysis for binary outcomes, the log-odds of the outcome probability is mathematically expressed as a summation of the SHAP values of each of the covariates. The exponent of a single SHAP value (similar to logistic regression), is thus the ratio of the odds of the outcome given that SHAP value divided by the baseline odds of the outcome (using only the intercept), i.e. the odds ratio.

For example, in a model with a single predictor, with  $x_i$  marking the SHAP value of that predictor:

- $\begin{cases} \log\left(\frac{p_0}{1-p_0}\right) = x_0 \\ \log\left(\frac{p_1}{1-p_1}\right) = x_0 + x_1 \end{cases} \Rightarrow$
- $\log\left(\frac{\frac{p_1}{1-p_1}}{\frac{p_0}{1-p_0}}\right) = \log\left(\frac{p_1}{1-p_1}\right) - \log\left(\frac{p_0}{1-p_0}\right) = x_0 + x_1 - x_0 = x_1 \Rightarrow$
- $\frac{\frac{p_1}{1-p_1}}{\frac{p_0}{1-p_0}} = e^{x_1}$

## Supplementary Methods 2: Pseudo-code of the "multi-calibration" algorithm

1. **Input:** An initial predictor  $\mathbf{p}'$ , a collection of subpopulations  $\mathbf{C}$ , a training set  $\mathbf{D} = \{(x_i, y_i)_{i=1}^m\}$  and a violation parameters  $\alpha > 0$
2. **Output:** a post-processed predictor  $\mathbf{p}$  satisfying  $(\mathbf{C}, \alpha)$ -multiaccuracy on  $\mathbf{D}$
3.  $\mathbf{p} \leftarrow \mathbf{p}'$
4.  $done \leftarrow False$
5. **while**  $\neg done$  **do**
  - 5.1.  $done \leftarrow True$
  - 5.2. **foreach**  $S \in \mathbf{C}$  **do** // iterate over subpopulations
    - 5.2.1.  $\Delta_S = \frac{1}{|S \cap \mathbf{D}|} (\sum_{x_i \in S} y_i - \sum_{x_i \in S} p_i)$  // magnitude of violation on  $S$
    - 5.2.2. **if**  $|\Delta_S| > \alpha$  **then**
      - 5.2.2.1.1.  $\mathbf{p} \leftarrow \mathbf{p} + \Delta_S \mathbf{1}_S$  // update predictor by "nudging"  $S$
      - 5.2.2.1.2.  $done \leftarrow False$
  - 5.2.3. **end if**
  - 5.3. **end foreach**
6. **end while**
7. **return**  $\mathbf{p}$

Detailed pseudo-code of the multi-calibration algorithm, as described in Hebert-Johnson et al.<sup>3</sup>

In this study, the algorithm accepts as input the predictions ( $\mathbf{p}'$ ) from the baseline model and the collection of subpopulations ( $\mathbf{C}$ ) defined by intersections of age and sex groups. It returns predictions recalibrated to COVID-19 outcome rates, which replace the (usually empirically derived) outcomes  $y_i$ .

## Supplementary Methods 3: Technical description of adjustment procedure

Purpose: to adjust predictions to external one-way conditional probabilities.

Inputs:

- Baseline predictions for another outcome on the entire population of interest,  $P(Y_{old}|\mathbf{X})$ , conditioned on as many independent variables as wanted.
- Published one-way conditional probabilities for the outcome of interest, from an external population, conditional on  $n$  variables.  
 $P_{external}(Y_{new}|X_i), i = 1..n$
- A co-occurrence matrix  $\mathbf{M}_{n \times n}$ , where  $\mathbf{M}_{i,j} = P_{external}(X_j|X_i)$  in the population on which the one-way conditional probabilities were published
- Rates of the different population characteristics in the local population,  $P_{local}(X_i), i = 1..n$

Process:

- By solving a system of linear equations,  $P_{external}(Y|\cdot) = \mathbf{M}\mathbf{w}$ , the  $n$  coefficients for each variable ( $W_i$ ) are calculated.
- These coefficients are used to derive  $P_{local}(Y|\cdot)$  by multiplying them with the local population characteristics.

Steps:

### 1. Solving for the coefficients of the linear probability model

This is done by deriving the coefficients for each independent variable in a linear probability model.

That is, the model is  $P(Y) = W_1P(X_1) + \dots + W_nP(X_n)$ , and the coefficients are solved for by solving a system of  $i = 1..n$  equations:

$$P_{external}(Y|X_i) = W_iP_{external}(X_1|X_i) + \dots + W_nP_{external}(X_n|X_i)$$

Where  $P(X_i|X_j)$  is the  $i, j$  entry of the  $\mathbf{M}$  matrix.

### 2. Adjust external one-way conditionals to the local population

With the coefficients in hand, rates for the local population are calculated using the local probabilities of the independent variables, for  $i = 1 \dots n$

$$P_{local}(Y|X_i) = W_i P_{local}(X_1|X_i) + \dots + W_n P_{local}(X_n|X_i)$$

3. Recalibrate the predictions for the baseline outcome to the outcome rates of the new outcome

With the calculated local outcome rates in hand, the baseline predictions are adjusted using the multi-accuracy algorithm depicted in Supplementary Methods 2.

When the process terminates, we have  $P(Y_{new}|\mathbf{X})$  for the entire population of interest.

## References

1. Lundberg, S.M. & Lee, S.-I. A unified approach to interpreting model predictions. in *Advances in neural information processing systems* 4765-4774 (2017).
2. Borgan, Ø. Aalen–Johansen Estimator. *Wiley StatsRef: Statistics Reference Online*, 1-13 (2014).
3. Hébert-Johnson, U., Kim, M.P., Reingold, O. & Rothblum, G.N. Calibration for the (computationally-identifiable) masses. *arXiv preprint arXiv:1711.08513* (2017).