

# S2 Statistical Methodology for Combining Rapid Antigen Testing and Syndromic Data Improves Sensitivity and Specificity in Real-World COVID-19 Detection

Fergus J Chadwick<sup>a,b</sup>, Yacob Haddou<sup>a,b</sup>, Tasnuva Chowdhury<sup>a</sup>, David Pascall<sup>c</sup>,  
Shayan Chowdhury<sup>e</sup>, Jessica Clark<sup>a,b</sup>, Joanna Andrecka<sup>f</sup>, Mikolaj  
Kundergorski<sup>d,b</sup>, Craig Wilkie<sup>d,b</sup>, Eric Brum<sup>f</sup>, Tahmina Shirin<sup>g</sup>, A S M  
Alamgir<sup>g</sup>, Mahbubur Rahman<sup>g</sup>, Ahmed Nawsher Alam<sup>g</sup>, Farzana Khan<sup>g</sup>, Janine  
Illian<sup>d,b</sup>, Ben Swallow<sup>d,b</sup>, Davina L Hill<sup>a,b</sup>, Dirk Husmeier<sup>d</sup>, Jason  
Matthiopoulos<sup>a,b</sup>, Katie Hampson<sup>a,b</sup>, Ayesha Sania<sup>h</sup>

<sup>a</sup>Institute of Biodiversity, Animal Health and Comparative Medicine, University of Glasgow

<sup>b</sup>COVID-19 in LMICs Research Group, University of Glasgow

<sup>c</sup>MRC Biostatistics Unit, University of Cambridge

<sup>d</sup>School of Mathematics and Statistics, University of Glasgow

<sup>e</sup>a2i, United Nations Development Program, ICT Ministry, Bangladesh

<sup>f</sup>UN FAO in support of the UN Interagency Support Team, Bangladesh

<sup>g</sup>Institute of Epidemiology, Disease Control and Research, Ministry of Health, Bangladesh

<sup>h</sup>Division of Developmental Neuroscience, Department of Psychiatry, Columbia University

Below we have extended the modelling description provided in the main text to  
include more technical detail. The code used to implement these tasks is available  
at [https://github.com/fergusjchadwick/COVID19\\_SyndromicRATDiagnosis](https://github.com/fergusjchadwick/COVID19_SyndromicRATDiagnosis).

## 0.1. Modelling

### 0.1.1. Structure

We examined the ability of the two imperfect identification methods, syn-  
dromic modelling and RAT, to predict the patient's COVID-19 status when used

---

\*Corresponding Author

Email addresses: [f.chadwick.1@research.gla.ac.uk](mailto:f.chadwick.1@research.gla.ac.uk) (Fergus J Chadwick),  
[yacob.haddou@glasgow.ac.uk](mailto:yacob.haddou@glasgow.ac.uk) (Yacob Haddou), [tasnuvachowdhury2004@gmail.com](mailto:tasnuvachowdhury2004@gmail.com) (Tasnuva  
Chowdhury), [david.pascall@mrc-bsu.cam.ac.uk](mailto:david.pascall@mrc-bsu.cam.ac.uk) (David Pascall),  
[shayan.chowdhury@a2i.gov.bd](mailto:shayan.chowdhury@a2i.gov.bd) (Shayan Chowdhury), [Jessica.Clark@glasgow.ac.uk](mailto:Jessica.Clark@glasgow.ac.uk)  
(Jessica Clark), [aandrecka@gmail.com](mailto:aandrecka@gmail.com) (Joanna Andrecka), [mikolaj.kundergorski@gmail.com](mailto:mikolaj.kundergorski@gmail.com)  
(Mikolaj Kundergorski), [craig.wilkie@glasgow.ac.uk](mailto:craig.wilkie@glasgow.ac.uk) (Craig Wilkie), [eric.brum@fao.org](mailto:eric.brum@fao.org)  
(Eric Brum), [tahmina.shirin14@gmail.com](mailto:tahmina.shirin14@gmail.com) (Tahmina Shirin), [aalamgir@gmail.com](mailto:aalamgir@gmail.com) (A S M  
Alamgir), [dr\\_mahbub@yahoo.com](mailto:dr_mahbub@yahoo.com) (Mahbubur Rahman), [anawsher@yahoo.com](mailto:anawsher@yahoo.com) (Ahmed  
Nawsher Alam), [farzanakhan\\_25@yahoo.com](mailto:farzanakhan_25@yahoo.com) (Farzana Khan), [janine.illian@glasgow.ac.uk](mailto:janine.illian@glasgow.ac.uk)  
(Janine Illian), [ben.swallow@glasgow.ac.uk](mailto:ben.swallow@glasgow.ac.uk) (Ben Swallow), [davina.hill@glasgow.ac.uk](mailto:davina.hill@glasgow.ac.uk)  
(Davina L Hill), [dirk.husmeier@glasgow.ac.uk](mailto:dirk.husmeier@glasgow.ac.uk) (Dirk Husmeier),  
[jason.matthiopoulos@glasgow.ac.uk](mailto:jason.matthiopoulos@glasgow.ac.uk) (Jason Matthiopoulos), [katie.hampson@glasgow.ac.uk](mailto:katie.hampson@glasgow.ac.uk)  
(Katie Hampson), [ays328@mail.harvard.edu](mailto:ays328@mail.harvard.edu) (Ayesha Sania)

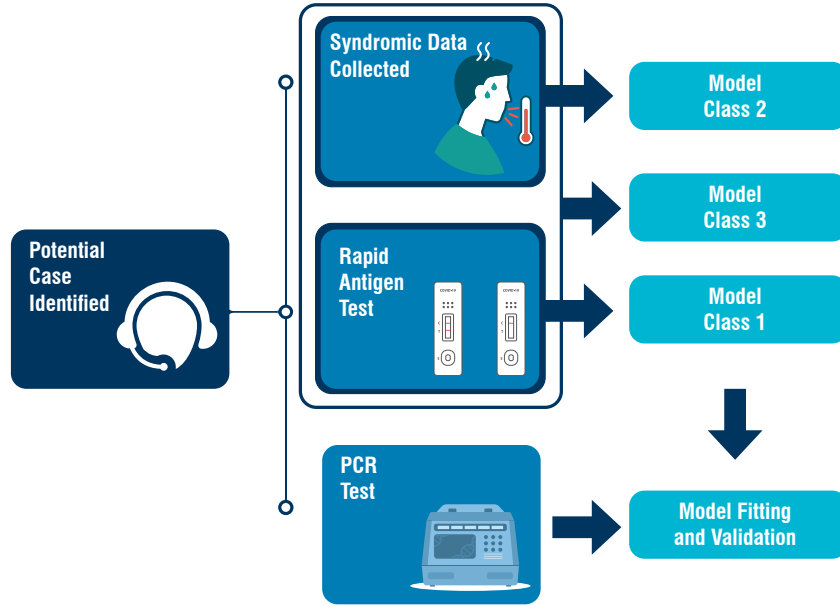


Figure 1: Schematic description of identification of likely COVID-19 patients by community support teams (CSTs), swab collection and model definitions. The teams collected syndromic data (age, gender and presence/absence of 14 predetermined symptoms), and two sets of naso-pharyngeal swabs (one each for Rapid Antigen Testing and RT-PCR). We then used rapid antigen testing (RAT) and syndromic data, two imperfect but inexpensive diagnostics, to generate three model classes: RAT result only in Model Class 1, syndromic data only in Model Class 2, and both RAT result and syndromic data in Model Class 3. The RT-PCR test result is used to train and test each model using temporal cross-validation.

separately and together. These combinations define three model classes (Figure 1).

Model Class 1 uses only the RAT result. It equates being RAT-positive with the patient being PCR-positive for COVID-19 (hereafter, PCR-positive), and being RAT-negative with PCR-negativity. We can write this formally for the  $i$ th individual as:

$$\begin{aligned} \text{PCR}_i &= \text{RAT}_i \\ \text{PCR} &\in \{0, 1\} \\ \text{RAT} &\in \{0, 1\} \end{aligned} \tag{1}$$

Model Class 2 uses only the syndromic data. For this model, we used a Bayesian multivariate probit model [1]. The multivariate probit structures the outcomes of the PCR test and symptoms presence/absence as a  $D$ -dimensional vector of binary outcomes ( $\mathbf{y}_i = (y_{i1}, y_{i2}, \dots, y_{id}), y_{ij} \in \{0, 1\}$ ). These outcomes are determined by an indicator function which takes a  $D$ -dimensional vector of *continuous latent* variables ( $\mathbf{z}_i = (z_{i1}, z_{i2}, \dots, z_{id}), z_{ij} \in \mathbb{R}$ ). These latent continuous variables then covary as realisations of a  $D$ -dimensional multivariate normal, with the mean of the error structure informed by a linear predictor,  $\sum_{j=1}^J x_{ij}\beta_{jd} + \epsilon_{id}$ , and a covariance between dimensions,  $\Sigma$ . The linear predictor allows us to condition the outcomes on risk factor variables (here, age and gender) while the covariance structure allows us to account for the correlated nature of the symptoms with each other and the outcome. As Albert and Chib [1] identify in their paper, the covariance matrix formulation is not identifiable, with the variance,  $\text{diag}(\Sigma)$  and means of the latent variables,  $\mathbf{z}_i$  trading off against each other. For this reason, we use a correlation matrix,  $\Omega$ , formulation with the variance set to 1. A correlation based framework also makes communication with clinicians and other practitioners smoother as correlations are more familiar.

$$\begin{aligned} y_{id} &= \mathbb{I}(z_{id} > 0) \\ \mathbf{z}_i &= \mathbf{x}_i\boldsymbol{\beta} + \boldsymbol{\epsilon}_i \\ z_{id} &= \sum_{j=1}^J x_{ij}\beta_{jd} + \epsilon_{id} \\ \boldsymbol{\epsilon}_i &\sim N(\mathbf{0}, \boldsymbol{\Omega}) \\ \text{diag}(\boldsymbol{\Omega}) &= 1 \\ \beta &\sim N(0, 1) \\ \boldsymbol{\Omega} &\sim \text{LKJ}(1) \end{aligned} \tag{2}$$

Model Class 3 combines the two data sources. We utilise the specificity of RAT by treating RAT-positive patients as PCR-positive patients. The RAT-negative patients are modelled using the sensitive syndromic approach using Model Class 2 to capture PCR-positive patients that are missed by the RAT. This approach leverages the potential different syndromic profiles of PCR-positive patients who are RAT-positive and -negative, allowing the model to adapt solely

Find more ways to cite Figure reffig:data-flowchart2 in this text.

55 to the latter. Structurally, the model combines Equations (??) and (2), with  
 56 RAT-positive patients being modelled using Equation (??), and Equation (2).

$$\begin{aligned}
 y_{i1} &= \begin{cases} 1, & \text{if RAT}_i = 1 \\ \text{otherwise} \end{cases} \\
 y_{id} &= \mathbb{I}(z_{id} > 0) \\
 \mathbf{z}_i &= \mathbf{x}_i \boldsymbol{\beta} + \boldsymbol{\epsilon}_i \\
 z_{id} &= \sum_{j=1}^J x_{ij} \beta_{jd} + \epsilon_{id} \\
 \boldsymbol{\epsilon}_i &\sim N(\mathbf{0}, \boldsymbol{\Omega}) \\
 \text{diag}(\boldsymbol{\Omega}) &= 1\beta \quad \quad \quad \sim N(0, 1) \\
 \boldsymbol{\Omega} &\sim \text{LKJ}(1)
 \end{aligned} \tag{3}$$

57  
 58 By using a Bayesian formulation, we generate full posteriors for our parameter  
 59 estimates, allowing natural quantification of uncertainty. Bayesian methods also  
 60 facilitate the use of more informative priors. While we used minimally informative  
 61 priors here (standard normals in the probit scale for betas and an LKJ correlation  
 62 prior with minimal shrinkage,  $\eta = 1$  [2]), more informative priors that incorporate  
 63 spatio-temporal effects, for instance, would be natural extensions. The models  
 64 were fitted to the data using Bayesian inference techniques based on Hamiltonian  
 65 Monte Carlo in the Stan programming language [3]. The models all converged  
 66 with zero divergent transitions and large effective sample sizes.

This isn't quite right - need to demonstrate subsetting of i's, otherwise it implies that z is calculated for all data points

#### 67 0.1.2. Model Selection

68 We conducted backwards model selection (starting with the most complex,  
 69 biologically plausible model) to identify a subset of models with the highest  
 70 predictive power under temporal cross-validation (Figure 2). For the cross-  
 71 validation, we divided the data into 5 folds of equal sizes in time order (i.e. the  
 72 first fold is formed of the chronologically first  $\frac{N}{K}$  patients, where  $N$  is the number  
 73 of patients and  $K$  is the number of folds, the second fold by the next  $\frac{N}{K}$  etc.) To  
 74 test the sensitivity of this cross-validation structure, we also did a strict temporal  
 75 division (i.e. the first  $\frac{T}{K}$  days where  $T$  is the number of days samples were taken  
 76 on). The results did not change qualitatively between these approaches.

77 The coarse round of model selection (Figure 2) selected candidate symptoms  
 78 based on whether they had a strong and consistent correlation with PCR as  
 79 estimated according to Equations (2) and (3). The models were fit with both  
 80 covariates throughout the coarse round and symptoms were compared in nested  
 81 models. In the fine round of model selection, these candidate symptoms and the  
 82 covariate combinations (age and gender, age, gender and no covariates) were  
 83 permuted to more exhaustively explore the model space. Reducing the number  
 84 of possible models using the two stages of model selection was necessary to  
 85 reduce computational demand and reduce the risk of overfitting models to the  
 86 test scenarios. The large number of symptoms corresponds to a high number of  
 87 potential model configurations (>131 000 for 14 symptoms and two covariates)

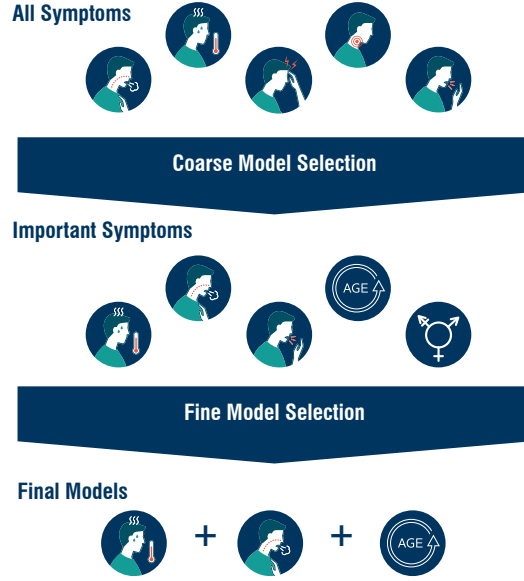


Figure 2: Schematic for rounds of model selection in the multivariate probit component of Model Classes 2 and 3. With 14 symptoms (only 5 shown here for demonstration purposes) and two covariates there are over 131000 possible model combinations. To make exploring these possible models computationally feasible and to reduce the risk of overfitting, we carried out two rounds of model selection. First, the data are divided into temporal cross-validation sets. The multivariate probit connects symptoms to the RT-PCR result through a correlation matrix. In the coarse model selection, the most complex feasible model (all symptoms and covariates) is fit to the training data. The estimated correlations between each symptom and the RT-PCR result are compared for each cross-validation set. The symptoms that have non-zero correlations in a systematic direction (i.e. all positively or all negatively correlated with RT-PCR result) are retained. The process is then repeated on each retained set of symptoms until the four symptoms in each model class with the strongest correlation to RT-PCR result. We then conduct a more exhaustive model selection on all the possible permutations of the four symptoms and two covariates. In this round, each model is fit to training data and used to predict for the test set, and the quality of those predictions is measured using cross-entropy scoring. The cross-entropy score is then used to select the best predictive model for each level of model complexity. Only these final models are then used for classification. This reduces the set of models tested as classifiers from  $>131\,000$  to just four per model class.

which might perform well on the test sets (even under the challenging conditions of temporal cross-validation) but lack transferability.

By using general predictive power to narrow down the number of candidate models and then testing those models, we are more likely to choose models which generalise well to new data. The number of candidate models used was not pre-determined but it was clear when fitting the models that there were “jumps” in performance (as defined below) between models containing five and four symptoms, so the models with one to four symptoms were used as the candidate models. Zero symptom models were not included in the analysis as they do not correspond to a feasible policy (with covariates they would require governments to ask individuals of a given gender and age as COVID-19 positive, and without covariates they would involve randomly assigning individuals as COVID-19 positive).

Model selection carried out in this way is quite robust. However, if more time were available for model development we would advocate more sophisticated methods of variable selection, such as penalised complexity priors [4] or auxiliary variable selection [5]. These methods, once developed for the use-case, are superior but require significant testing that would prevent getting this work out in a policy-relevant timeframe. We also believe that these approaches may be harder for non-statisticians to adapt.

#### 0.1.3. Predictive Performance

We scored the models’ predictive power using binary cross-entropy (hereafter, cross-entropy). Cross-entropy measures the accuracy of models that generate probabilities of binary outcomes, rather than make binary classifications, similar in concept to a mean square error for normally-distributed data, but adapted for binary data [6]. A cross-entropy value close to zero corresponds to high levels of accuracy, with larger values indicating lower accuracy. More specifically, the metric allows us to compare a binary vector,  $\mathbf{y} \in [0, 1]$ , with a vector of probabilistic predictions ( $p(\mathbf{y}) \in (0, 1)$ ) as follows:

$$\mathbf{H}_p(q) = \frac{1}{N} \sum_{i=1}^N y_i \cdot \log(p(y_i)) + (1 - y_i) \cdot \log(1 - p(y_i)) \quad (4)$$

The resulting score is comparable across all methods for assigning predictions where the same test data are used, allowing us to compare predictions from Model Classes 1-3.  $\mathbf{H}_p(q) \in 0, \mathbf{R}_+$  with zero indicating perfect prediction (assigning probabilities of ones and zeroes to outcomes of ones and zeros exactly) and larger values indicating worse predictions.

#### 0.1.4. Classification Performance

In applied settings, models must often be evaluated on their performance as classifiers rather than just as prediction engines (i.e. their ability to say a patient is COVID-19 positive or negative, not simply the probability the patient might be COVID-19 positive or negative). To generate a classification,  $\hat{Y}$ , a

127 probability threshold,  $\hat{p}$ , must be chosen over which patients are classified as  
 128 COVID-19 positive:

$$\hat{Y} = \begin{cases} 1, & \text{if } p(y) \geq \hat{p} \\ 0 & \text{otherwise} \end{cases} \quad (5)$$

129 Receiver operating characteristics (ROCs) are a way to measure the per-  
 130 formance of a set of classifications in terms of true and false positives and  
 131 negatives (TP, FP, TN and FN) and the rates of each of these classification types  
 132 (e.g.  $TPR = \frac{TP}{TP+FN}$ , and  $FPR = \frac{FP}{FP+TN}$ ). The error rates are calculated with  
 133 respect to a particular threshold,  $\hat{p}$ , or across the range of possible  $\hat{p}$ s to generate  
 134 a ROC curve [7]. In our epidemiological scenarios (outlined below) we use our  
 135 ROC curve calculations to identify single thresholds which yield a required error  
 136 rate.

137 We strongly emphasise that generic performance here is only used to show  
 138 the flexibility of the model classes; the best model for a local situation can only  
 139 be determined if the relative cost of false positives and false negatives is known.  
 140 Here, we choose three representative scenarios. Each scenario has a requirement  
 141 and error rate (defined in Table 1).

142 In Scenario 1, we do not consider epidemiological context but simply minimise  
 143 false negative and false positive rates equally. We do this by maximising the two  
 144 correct classification rates both individually and in total, as measured by the  
 145 harmonic mean (as opposed to the arithmetic mean which would only maximise  
 146 the rates in total).

$$\text{Find } \hat{p} \text{ which } \max \left( \frac{2 \cdot TPR \cdot TNR}{TPR + TNR} \right), \epsilon = FPR + FNR$$

147  
 148 Scenario 2 corresponds to the current situation in Bangladesh at time of  
 149 writing (July 2021), with COVID-19 cases beginning to rapidly increase again.  
 150 Under these circumstances, false negatives are extremely costly relative to false  
 151 positives due to the exponential growth of the disease.

$$\text{Find } \hat{p} \text{ which } \max (FNR - 0.2 \leq 0), \epsilon = FPR$$

152 In Scenario 3, the pandemic is not declining but maintaining a steady rate  
 153 of cases. In this situation, policy-makers may be keen to keep false positive  
 154 diagnoses low to prevent lockdown fatigue and to keep the workforce active.

$$\text{Find } \hat{p} \text{ which } \max (FPR - 0.2 \leq 0), \epsilon = FNR$$

155 The requirements in Scenario 2 and 3 were developed in discussion with the  
 156 Institute of Epidemiology, Disease Control and Research (IEDR), Bangladesh,  
 157 for illustrative purposes.

158 [1] Albert JH, Chib S. Bayesian analysis of binary and polychotomous re-  
 159 sponse data. Journal of the American Statistical Association 1993;88:669–  
 79.

I'm not sure  
if any of  
these equa-  
tions are  
correct or  
sufficient or  
clear

Table 1: For each epidemiological scenario there is a requirement and a performance criterion. The requirement refers to a base level of performance the model must achieve; in general this will be a maximum acceptable error rate of some kind. These requirements were determined in discussion with members of the Institute of Epidemiology, Disease Control and Research, Ministry of Health, Bangladesh (IEDCR). The requirement determines a probability threshold for each model which most closely meets that requirement. The performance criterion is then used to determine which model performs the 'best' given that the requirement has been met.

| Scenario Name     | Requirement                           | Performance Criterion (Error) |
|-------------------|---------------------------------------|-------------------------------|
| 1 Agnostic        | Maximise correct classification rates | Sum of error rates            |
| 2 Rising Cases    | Max. 20% false negative rate          | False negative rate           |
| 3 Low-Level Cases | Max. 20% false positive rate          | False positive rate           |

- [2] Lewandowski D, Kurowicka D, Joe H. Generating random correlation matrices based on vines and extended onion method. *Journal of Multivariate Analysis* 2009;100:1989–2001.
- [3] Carpenter B, Gelman A, Hoffman MD, Lee D, Goodrich B, Betancourt M, et al. Stan: A probabilistic programming language. *Journal of Statistical Software* 2017;76:1–32.
- [4] Simpson D, Rue H, Riebler A, Martins TG, Sørbye SH. Penalising model component complexity: A principled, practical approach to constructing priors. *Statistical Science* 2017;32:1–28.
- [5] Held L, Holmes CC. Bayesian auxiliary variable models for binary and multinomial regression. *Bayesian Analysis* 2006;1:145–68.
- [6] Gneiting T, Raftery AE. Strictly proper scoring rules, prediction, and estimation. *Journal of the American Statistical Association* 2007;102:359–78.
- [7] Hoo ZH, Candlish J, Teare D. What is an ROC curve? 2017.