

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22

Radiologist observations of chest X-rays (CXR) predict sputum smear microscopy status in TB
Portals, a real-world database of tuberculosis (TB) cases

Gabriel Rosenfeld^{1*}, Andrei Gabrielian¹, Alyssa Meyer², Alex Rosenthal³

¹Bioinformatics and Computational Biosciences Branch, Office of Cyber Infrastructure and Computational Biology, National Institute of Allergy and Infectious Diseases, National Institutes of Health, Bethesda, Maryland, United States of America

²Software Engineering Branch, Office of Cyber Infrastructure and Computational Biology, National Institute of Allergy and Infectious Diseases, National Institutes of Health, Bethesda, Maryland, United States of America

³ Office of Cyber Infrastructure and Computational Biology, National Institute of Allergy and Infectious Diseases, National Institutes of Health, Bethesda, Maryland, United States of America

* Corresponding author

E-mail: gabriel.rosenfeld@nih.gov

23 **Abstract**

24 The Tuberculosis (TB) Portals is an international program of 14 countries connecting clinical,
25 genomic, and radiologist specialists to develop an openly available repository of deidentified TB
26 cases with multi-modal data such as case clinical characteristics, pathogen genomics, and
27 radiomics. This real-world data resource contains over 4000 TB cases, principally drug resistant
28 cases, with over 4000 chest X-rays (CXR) images. The scope of curated data offers a case-
29 focused perspective into the drivers of disease incorporating the chronological context of the
30 presented CXR data. Here, we analyze a cohort consisting of new TB cases to understand the
31 relationship between baseline sputum microscopy status and nearby Chest X rays images. The
32 Timika score, a lung biomarker of disease severity, was derived for each CXR using available
33 radiologist observations. The Timika score along with the radiologist observations were
34 compared for predictive performance of baseline sputum microscopy status. Baseline sputum
35 microscopy status is a useful marker of pre-treatment disease severity and infectiousness. The
36 modeling results support that both the radiologist observations as well as Timika score are
37 predictive of smear status and that Timika score performs similarly to the top 5 radiologist
38 features by feature selection. Moreover, inferential statistical analysis identifies the factors
39 having the greatest association with sputum smear positivity such as presence of radiologist
40 observations in both lungs, presence of cavity, presence of nodule, and Timika score itself. The
41 results are consistent with prior reports showing Timika Score utility for predicting baseline
42 sputum smear and disease status. We report testing of Timika Score on the largest, openly
43 available real-world dataset of TB cases that can serve as a reference to explore extant and new
44 TB disease severity scores bridging radiological, microbiological, and clinical data. To illustrate,
45 we visualize Timika score from images in our database with other cases characteristics

46 demonstrating that this score captures lung biomarker status consistent with known clinical risk
47 factors.

48 **Introduction**

49 Tuberculosis (TB) remains a major global pandemic with approximately 10 million new
50 cases and 1.5 million deaths each year (1, 2). With the emergence of the SARS-Cov2 global
51 pandemic in 2020, it is estimated that the TB pandemic may have worsened due to additional
52 strains and challenges encountered via healthcare systems around the world (3, 4). Concurrent to
53 those unfortunate events, drug resistant TB continues to be a persistent threat with up to ~20% of
54 TB isolates around the world estimated to be resistant to a major drug. Transmission of drug
55 resistant TB is an emerging phenomenon closely monitored by health authorities worldwide (5).
56 Drug resistant TB cases (DR-TB) are associated with poorer outcomes and more expensive cost
57 of care when compared to drug sensitive TB. DR-TB has a lower treatment success of
58 approximately 55% globally and Multi- or Extensively DR-TB care can cost up to 25 times that
59 of TB cases that are drug sensitive (6, 7). Therefore, real-world databases focusing on these DR
60 cases that span multi-domain case information are essential to identify novel relationships and
61 aspects of drug resistance to enable translational medicine to timely and efficiently address drug
62 resistance.

63 To eradicate TB, clinicians need rapid diagnostics of disease along with efficient means
64 of monitoring treatment response and completeness at discharge. Sputum smear microscopy has
65 been a primary method for diagnosis of pulmonary tuberculosis in low and middle income
66 countries (LMIC) since it is a relatively simple, rapid, and less costly approach that can identify
67 the most infectious patients and be applied in a variety of socio-economic status areas.
68 Nonetheless, this approach shows deficiencies in certain demographic groups such as extra-

69 pulmonary TB, pediatric TB, and TB patients simultaneously infected with HIV (8). Moreover,
70 the requirement for repeated sputum sampling can present obstacles to the application of the
71 approach as patients may not return for results, are lost to follow up, or have difficulty producing
72 usable sputum samples. Despite these challenges, it is still widely used throughout LMIC for
73 disease monitoring and response to therapy (9) as well as having demonstrated some ability to
74 predict treatment response albeit requiring additional clinical factors (10). Since this
75 microbiological information is sometimes unavailable or inclusive, it is important to identify
76 other modalities that may assist with diagnosis or monitoring of treatment response, and one
77 such approach is imaging of the lungs via Chest X-rays (CXRs).

78 CXR imaging is often collected during TB disease management to understand treatment
79 response and disease status. Unlike computed tomography (CT) imaging that may not always be
80 available due to the cost of associated infrastructure (11), CXRs are the primary means of
81 assessing lung status in LMIC due to their relatively lower costs (12). As such, they are more
82 widely available to clinicians for assessing lung status during TB disease management and used
83 as a decision-making clinical information point compared to CTs. Radiologist assessment of
84 CXRs have been the gold standard reference upon which CXRs have been interpreted for clinical
85 decision making historically. These observations provide an important lung biomarker that can
86 inform patient risk, disease severity, and response to treatment over the course of a TB case. For
87 example, Heo et al. tracked radiological lesions from CXRs over the course of TB treatment in a
88 prospective cohort analysis showing that presence of cavity or fibrotic lesion associated with
89 poor radiological response (13). Another example is the development of CXR-derived Timika
90 Score that has been associated with baseline sputum smear microscopy status and disease
91 severity in TB cases (14, 15).

92 The National Institute of Allergy and Infectious Diseases (NIAID) Office of Cyber
93 Infrastructure and Computational Biology leads the transnational partnership of participating
94 sites covering 14 countries with heavy DR TB burden. This partnership created the TB Portals
95 to facilitate TB data sharing and science with a goal towards a better understanding of the real-
96 world aspects of especially problematic DR TB. The TB Portals resource consists of a repository
97 of TB case data including multiple domains such as case clinical characteristics, pathogen
98 genomics, and radiomics that can support the biomedical research community's research efforts
99 towards TB. As of April 2021, the TB Portals database includes over 4000 TB cases, mostly
100 drug resistant cases, with over 4000 CXRs. Many of these cases also have radiologist
101 annotations for their CXRs to assess lung biomarker status in relation to the clinical and
102 microbiological characteristics of the case. While other resources have large numbers of chest
103 X-ray images, TB Portals provides a TB case-centered repository encapsulating the
104 chronological context associated with the CXR such as drug resistance status, regimens
105 administered so far, the genome of the pathogen, and sputum microscopy status. External
106 collaborators can apply for access to publicly shared data through an online data use agreement
107 (DUA) and then download this data to facilitate reproducibility and open science.

108 In this study, we utilize the radiologist observations for CXR images in the TB Portals
109 repository to derive Timika Score (15), a useful numerical lung biomarker, to assess its utility for
110 predicting sputum smear microscopy status. We compare Timika Score performance with other
111 features we derived from the radiologist-reported observations to determine if the additional
112 features could improve upon Timika's previously reported performance. We select a cohort of
113 cases with a case definition of new containing sputum smear microscopy results from specimens
114 taken prior to start of treatment, as well as CXRs with radiologist observations within two weeks

115 of the specimen date. We perform inferential statistical analysis of risk of positive sputum
116 microscopy from presence of various features derived from radiologist observations. We also
117 examine Timika Score in relation to other aspects of the case such as demographics, case
118 definition, and outcome. For instance, we utilize a strength of this resource in having a larger
119 number of mono drug resistant (Mono DR), poly drug resistant (Poly DR), Multi-drug resistant
120 (MDR) and Extensively drug resistant (XDR) according to WHO guidelines (16).

121 We report results consistent with prior publications regarding the utility of the Timika
122 score for predicting baseline sputum status. Importantly, we show that Timika Score offers
123 similar predictive performance compared to the top 5 features we derive from radiologist
124 observations. These results suggest that Timika Score is well-optimized for determining pre-
125 treatment disease infectiousness and severity status.

126 **Materials and Methods**

127 **Computing environment**

128 All analyses were done on a MacBook Pro laptop (x86_64-apple-darwin15.6.0 (64-bit)
129 Running under: macOS Mojave 10.14.6) using R version 4.0.2 (2020-06-22) and RStudio
130 1.2.5033. Specific versions of the R packages can be found in the associated code which
131 contains renv.lock file listing all used packages and version numbers.

132 **Cohort selection**

133 *Sputum Prediction*

134
135 To remove the potential of confounding lung biomarkers due to prior history of TB, new
136 cases of TB with CXRs containing radiologist annotations as well as sputum microscopy test

137 results from specimens prior to or on the treatment start date were selected. Those cases where
138 the specimen collection date was within 14 days of a CXR were included. For cases where
139 multiple pairs of specimens and CXRs existed, the last specimen prior to treatment was used.
140 For cases where multiple microscopy test results were present for a specimen, the last
141 microscopy test result for that specimen was used. For cases where multiple images existed, the
142 last imaging date was used. Unknown or non-standard microscopy results such as “Unknown
143 data” and “Saliva” were excluded. Code used for generating the cohort is provided in the Data
144 availability and code section. Ultimately, 572 new cases were selected from the database with
145 sputum microscopy results of Negative, 1 to 9 in 100 (1-9/100), 10 to 99 in 100 (1+), 1 to 9 in 1
146 (2+), 10 to 99 in 1 (3+), and More than 99 in 1 (4+) consisting of 259, 29, 144, 60, 58, and 22
147 cases respectively within this cohort. The cohort characteristics summarizing case details for the
148 sputum prediction cohort can be found in Supplementary Table 1.

149 *Analysis of Timika Score with regards to other case characteristics in TB portals*

150 For Figure 1 and Figure 2 visualizing the Timika Score in relation to other case
151 characteristics, we used all available images with manual radiologist annotations from the
152 February 2021 release of TB Portals data available for download from Aspera. This included
153 2058 images from 1761 cases covering not just New cases but all other types in the database.
154 The characteristics summarizing corresponding case details for the set of available images for the
155 Timika Score visualizations can be found in Supplementary Table 2.

156 **Data preprocessing and extraction of feature set**

157 Data from TB portals was downloaded as a list of .csv files from the Aspera file share
158 service using the February 2021 version of each respective file. The CXR manual annotations
159 are provided as a set of features corresponding to observations by radiologists within each

160 sextant of the lungs (dividing each lung by 1/3) as well as a set of features provided at the level
161 of the entire lung. For those features corresponding to sextant level observations, features where
162 no observations were provided by radiologists were imputed as 0 (for numerical data
163 corresponding to between 0-100% involvement of sextant) or “No” (for categorical data
164 indicating presence/absence of a specific feature within the sextant). The omission of these at
165 either the level of the entire sextant or specific sextant-level feature are interpreted as the
166 radiologist did not observe the feature.

167 After imputation, features were converted for tidy-like data processing using packages
168 from the R tidyverse. This permits various types of downstream feature engineering such as
169 identification of involvement of one or both lungs by sextant-level feature type, calculation of
170 summary statistics for numerical features across sextants, and other score calculations such as
171 Timika Score. For this analysis, involvement of both lungs as well as mean percentage of
172 sextant involvement across all sextants by specific sextant-level radiologist observation was
173 calculated along with Timika Score. Both lung features were calculated in a binary manner
174 where involvement of a left and right sextant for the features was required to indicate
175 involvement of both lungs for that feature. Timika Score was calculated like the original
176 publication (15) by a simple method of taking the overall abnormal percent of volume of the
177 lungs reported by the radiologist and adding 40 if the presence of cavity was indicated in the
178 radiologist report. Characteristics of derived features by microscopy status can be found in
179 Supplementary Table 3.

180 MLR3 framework was used to define a set of prediction tasks as well as pipelines for
181 modeling (17). 70% of the data was selected as a training set and 30% was held out as a
182 validation set. We tested two distinct prediction tasks in the MLR3 framework: 1) to predict

183 sputum positive (1 to 9 in 100, 1+, or higher) compared to negative and 2) to predict higher
184 bacterial load positive (2+, 3+, or higher) compared to negative. The positive to negative
185 prediction task was relatively well balanced between classes; however, the high bacterial load
186 positive to negative prediction task showed moderate class imbalancing so a class balancing step
187 was included in some pipelines for comparison. Machine learning (ML) pipelines using all
188 derived radiological features were compared to pipelines using only the Timika Score for
189 prediction. For ML pipelines using all derived radiological features, factor data was encoded to
190 a binary indicator (https://mlr3pipelines.mlr-org.com/reference/mlr_pipeops_encode.html), low
191 variance features were removed ([https://mlr3pipelines.mlr-](https://mlr3pipelines.mlr-org.com/reference/mlr_pipeops_removeconstants.html)
192 [org.com/reference/mlr_pipeops_removeconstants.html](https://mlr3pipelines.mlr-org.com/reference/mlr_pipeops_removeconstants.html)), features were scaled by min-max
193 scaling (https://mlr3pipelines.mlr-org.com/reference/mlr_pipeops_scalerange.html), and the top
194 5 features were selected via a variety of feature selection methods ([https://mlr3filters.mlr-](https://mlr3filters.mlr-org.com/)
195 [org.com/](https://mlr3filters.mlr-org.com/)) or Principal Component Analysis ([https://mlr3pipelines.mlr-](https://mlr3pipelines.mlr-org.com/reference/mlr_pipeops_pca.html)
196 [org.com/reference/mlr_pipeops_pca.html](https://mlr3pipelines.mlr-org.com/reference/mlr_pipeops_pca.html)). The following ML models were assessed as part of
197 the pipelines: featureless, glmnet, kkn, multinom, naïve bayes, rpart, ranger, xgboost, svm, and
198 nnet as described in the subsequent link (<https://mlr3learners.mlr-org.com/reference/index.html>).
199 For pipelines using only Timika score, only the min-max scaling step was included as part of the
200 pipeline.

201 **Model performance benchmarking**

202 5-fold cross validation was used to assess various binary classification metrics towards
203 respective prediction tasks on the training set. Both the metrics and the resampling strategies can
204 be found in the mlr3 documentation (<https://mlr3.mlr-org.com/reference/index.html>) under
205 Measures and Resampling Strategies sections. We also assess these binary classification metrics

206 in the validation dataset to ascertain performance on data which has never been used during
207 model training. To compare performance of the top radiologist observations and Timika Score,
208 the best pipelines using the top radiologist derived features were compared by bootstrapping
209 without replacement (N = 200) to the best pipelines using Timika Score alone, or a featureless
210 pipeline as a control to indicate the density of observed model performance particular to this
211 dataset that may be due to random chance.

212 **Calculation of inferential statistics**

213 To estimate the univariate Odds Ratios (OR) and multivariate adjusted Odds Ratios
214 (AOR), the finalfit R package was used. Both_hugenodule1 feature was removed from the
215 analysis as it did not show any variance. As Timika Score is highly correlated with other
216 variables in the dataset (e.g. overall abnormal volume and cavity), MRMR feature selection
217 (https://mlr3filters.mlr-org.com/reference/mlr_filters_mrmr.html) was performed for the top 5
218 features to include in the multivariate modeling. As both_hugenodule1 feature was excluded, less
219 than 5 features were selected in the multivariate modeling. The both_lungs feature was included
220 in the multivariate model to adjust for indication of involvement of both lungs from sextant level
221 features when assessing estimated odds of sputum positivity.

222 **Visualization of Timika Score with other case attributes**

223 Interesting case variables were examined for association with Timika Score to assess
224 consistency with current understanding of TB clinical risk factors. This included demographic
225 information such as age and BMI as well as case resistance status, case definition, and treatment
226 outcome. The initial CXR with available radiologist observations were selected for each case,

227 Timika Score calculated and plotted using ggplot2 to visualize associations with other case
228 attributes.

229 For calculating temporal changes in Timika Score, those cases with an initial CXR
230 identified above were filtered for cases with an additional follow up CXR with radiologist
231 observation. The log₂ transformed relative change were calculated for each image's Timika
232 Score comparing the earliest score with all subsequent scores. To account for differences in the
233 length of time between images that may impact the relative score change, the difference in
234 number of days between CXRs was calculated and the log₂ transformed relative changes were
235 divided by the number of days between pairs of images for each case to generate the log₂
236 relative change by day.

237 **Data availability and code**

238 The TB portals program necessitates all users of the data sign a DUA before access to the
239 underlying, de-identified clinical data is provided and the data can be requested at the following
240 URL (<https://tbportals.niaid.nih.gov/download-data>). Therefore, this study provides the code to
241 reproduce the analysis without the underlying raw data
242 (<https://github.com/niatd/tbportals.xray.sputum.2021>) in compliance with the DUA. To rerun the
243 analysis, interested parties can request data access by completing the DUA and then place the
244 downloaded files to the subdirectory of the data folder as provided in the GitHub repo
245 instructions. To aid reproducibility, the list of patient IDs and condition IDs used from the
246 sputum prediction analysis are provided in Supplementary Table 4. The specific record
247 identifiers for the set of images used for visualization of Timika Score in comparison with other
248 case attributes are provided in Supplementary Table 5. For cases where change in Timika was
249 calculated over time, the first and last images used in the case are shown along with the dates of

250 the image. Both supplementary tables allow those interested in examining the specific records to
251 do so after completion of required DUA irrespective of the evolution in number of available
252 cases in the database.

253 **Results**

254 **Timika score associates with case clinical characteristics, disease** 255 **severity, and risk**

256 The TB portals resource contains case information bridging across domains of interest
257 such as clinical, demographic, radiologic, and pathogen genomics. We leveraged the unique
258 value of these connections to assess any scores or other features of interest from the derived
259 radiological data. After generating the Timika Scores from all available CXRs with associated
260 radiologist annotations, we explored the relationship between Timika Score with the additional
261 information contained about the case mentioned above. We analyzed these relationships to
262 determine if they are consistent with prior TB clinical findings to assess the plausibility of the
263 derived radiological data.

264 **Age, BMI, Type of Resistance, Case Definition, and Case Outcome show** 265 **associations with Timika Scores**

266 We used only the initial image with associated radiologist observations for the
267 visualizations assessing Timika Score in relation to other aspects of the case. We visualize
268 Timika Score with age of onset, BMI, resistance type, case definition, and case outcome and
269 include a trendline in the relationships for any numerical features. We identify relationships

270 between Timika score and case characteristics of interest that are consistent with our prior
271 knowledge of TB clinical risk.

272 For instance, Timika Score of the initial image gradually increases with age of onset until
273 the relationship plateaus around age 50 and decreases although some of the decrease and initial
274 increase can likely be attributed to the lower density of observations at the two extremes of age
275 (Fig 1A). Timika Score decreases with increasing BMI until it plateaus around a BMI of ~25
276 (Fig 1B). Like age, the extremes of the BMI visualization need to be interpreted cautiously
277 given the lower densities. XDR cases, resistant to the most TB drugs and with the worst
278 outcomes, are observed to have a higher median Timika Score and interquartile range compared
279 to other case resistance types (Fig 1C). New cases of TB are shown to have a lower median
280 Timika Score and interquartile range compared to other types of cases such as Chronic TB, Prior
281 treatment failure, Relapse, Prior lost to follow up, or Other prior unknown status case definitions
282 (Fig 1D). Similarly, visualizing case outcomes reveals that detrimental treatment outcomes such
283 as Died, Treatment failure, Lost to follow up, or Unknown show higher median Timika Scores
284 compared to beneficial outcomes such as Treatment completion, Cured, or Still on treatment (Fig
285 1E). Taken together, the results demonstrate associations between Timika Score from initial
286 available image and case characteristics that reflects TB clinical risk.

287 **Fig 1 Association of Timika Score from initial CXR with radiologist observations with**
288 **other case attributes.** Timika Score derived from initial available CXRs associated with cases
289 in the TB Portals repository are visualized along with a variety of salient case characteristics
290 with missing observations dropped according to variable (N = 1757). For A) and B), the age of
291 onset (N = 1757) and BMI (N = 1268) from the case are shown with blue trend line respectively.
292 For C), D), and E), boxplots with interquartile range showing Timika Score by the type of drug
293 resistance, status of case at start, and status of case at end are shown. In C), MDR non XDR (N =
294 752), Mono DR (N = 118), Poly DR (N = 78), Sensitive (N = 514), and XDR (N = 295) case
295 drug resistance statuses are shown with the associated Timika Score from initial CXR with the
296 case. XDR cases tend to show relatively higher Timika Scores. In D), Chronic TB (N = 18),
297 Failure (N = 179), Lost to follow up (N = 65), New (N = 1141), Other (N = 45), Relapse (N =
298 304), and Unknown (N = 5) case definitions are shown with the associated Timika Score from

299 initial CXR with the case. Undesirable case definitions such as Failure, Lost to follow up,
300 Relapse, Chronic TB, or Unknown from prior history show higher Timika Score compared to
301 New cases. In E), Completed (N = 170), Cured (N = 984), Died (N = 126), Failure (N = 128),
302 Lost to follow up (N = 151), Still on treatment (N = 169), and Unknown (N = 29) case outcomes
303 are shown with the associated Timika Score from initial CXR with the case. Undesirable
304 outcomes such as Died, Failure, Lost to follow up, or Unknown show higher Timika Score
305 compared to beneficial outcomes such as Completed, Cured, or Still on treatment.

306 **Age, BMI, Type of Resistance, Case Definition, and Case Outcome show** 307 **associations with the temporal changes in Timika Scores**

308 Of those cases with initial CXR images visualized above, we next examined changes in
309 Timika Score whenever follow up CXR images were available. To do so, we filter on cases with
310 this additional imaging information. We calculate log₂ transformed Timika Score from initial
311 image to last available image per case dividing by the number of days between images to account
312 for the relative amount of time between each image. We use ggplot2 to visualize log₂
313 transformed change in Timika Score by day with age of onset, BMI, resistance type, case
314 definition, and case outcome and include a trend line in the associations for any numerical
315 features. Interestingly, most cases have a negative relative change in Timika Score by day
316 indicating improvement in lung status over the course of the case. Such a decrease over time
317 would be expected given these cases would have been undergoing clinical management. We
318 observed associations between the relative change by day and case characteristics of interest that
319 are consistent with prior knowledge of TB clinical risk.

320 For instance, the log₂ transformed change in Timika Score by day steadily decreases with
321 age of onset such that younger age shows greater relative change whereas older age shows less
322 relative change (Fig 2A). Conversely, log₂ transformed change in Timika Score by day
323 increases as BMI increases (Fig 2B). Lower BMI demonstrates a smaller relative change as
324 compared to higher BMI. When examining relative change by resistance type of the case, Drug

325 sensitive cases are observed to have a larger relative change compared to drug resistant types
326 (Fig 2C). Similarly, new cases of TB show greater relative change compared to other types of
327 cases such as Prior treatment failure, Prior lost to follow up, Relapse, or Other prior status at the
328 start of the case (Fig 2D). Visualizing by case outcomes demonstrates that undesirable outcomes
329 such as Died, Treatment failure, or Lost to follow up show decreased relative change by day
330 compared to beneficial outcomes such as Treatment completion or Cured (Fig 2E). Like earlier
331 visualizations of the Timika Score from available initial image, observations of relative changes
332 are consistent our prior understanding of TB clinical risk factors.

333 **Fig 2 Association of relative change in Timika Score by day from initial CXR with**
334 **radiologist observations to last available CXR with radiologist observations with other case**
335 **attributes.** Log₂ relative change in Timika Score by day from initial available CXR to last
336 available CXR associated with cases in the TB Portals repository are visualized along with a
337 variety of salient case characteristics with missing observations dropped according to variable (N
338 = 297). For A) and B), the age of onset (N = 297) and BMI (N = 292) from the case are
339 visualized with log₂ relative change in Timika Score by day with blue trend line respectively.
340 To aid in visualizing the trendline, the y axis was limited to between -0.1 and 0.1 resulting in an
341 additional 9 outliers being removed for age of onset and BMI case numbers above respectively.
342 For C), D), and E), boxplots with interquartile range showing log₂ relative change in Timika
343 Score by day compared by the type of drug resistance, status of case at start, and status of case at
344 end are shown. In C), MDR non XDR (N = 160), Mono DR (N = 13), Poly DR (N = 1),
345 Sensitive (N = 24), and XDR (N = 99) case drug resistance statuses are shown with the log₂
346 relative change in Timika Score by day from initial CXR to last available CXR. To aid in
347 visualization, y axis was limited to -0.1, and 0.1 resulting in additional 5, 2, 0, 2, and 0 outliers
348 being removed from MDR non XDR, Mono DR, Poly DR, Sensitive, and XDR case numbers
349 above respectively. In general, drug resistant cases show lower relative change in Timika Score
350 although several groups show low N and must be interpreted cautiously. In D), Failure (N = 35),
351 Lost to follow up (N = 8), New (N = 184), Other (N = 4), and Relapse (N = 66) case definitions
352 are shown with the log₂ relative change in Timika Score by day from initial CXR to last
353 available CXR. To aid in visualization, y axis was limited to -0.1, and 0.1 resulting in additional
354 0, 0, 8, 1, and 0 outliers being removed from Failure, Lost to follow up, New, Other, and Relapse
355 case numbers above respectively. Deleterious case definitions such as Failure, Lost to follow up,
356 Relapse, or Other from prior history show less change in relative Timika Score compared to New
357 cases. In E), Completed (N = 27), Cured (N = 223), Died (N = 15), Failure (N = 10), Lost to
358 follow up (N = 20), Still on treatment (N = 1), and Unknown (N = 1) case outcomes are shown
359 with the log₂ relative change in Timika Score from initial CXR to last available CXR. To aid in
360 visualization, y axis was limited to -0.1, and 0.1 resulting in additional 4, 5, 0, 0, 0 and 0 outliers
361 being removed from Completed, Cured, Died, Failure, Lost to follow up, Still on treatment, and
362 Unknown case numbers above respectively. Deleterious outcomes such as Died, Failure, Lost to

363 follow up, or Unknown show smaller relative changes compared to beneficial outcomes such as
364 Completed and Cured. As above for other visualizations, caution is warranted given the small
365 N's associated with certain subgroups.
366

367 **Sputum smear microscopy results associate with Timika score in this cohort** 368 **of TB portals cases**

369 We next assessed the previously reported role of Timika Score for predicting sputum
370 smear status by analyzing the selected cohort of new cases having microscopy results from
371 sputum specimens taken prior to treatment with associated images within two weeks of the
372 specimen (N = 572). Mean Timika Score is lowest amongst new cases with a sputum microscopy
373 result of negative and increases for microscopy results indicating a higher burden of bacteria
374 within sputum [1 to 9 in 100, 1+, 2+, etc.] (Fig 3). Only the 4+ level shows slightly lower mean
375 Timika Score compared to the next highest level of 3+, which may be due to variance from the
376 lower number of available cases in this 4+ level (N = 22). This clear trend in the TB portals
377 dataset is consistent with previously reported role of Timika Score for predicting baseline
378 sputum status.

379 **Fig 3 Timika score derived from radiologist observations of CXR within two weeks of**
380 **specimen taken prior to treatment start.** Timika score is visualized by the smear microscopy
381 results of specimens taken prior to treatment start for which CXRs were available within two
382 weeks of specimen date (N = 572). Boxplots show median Timika Score for associated images
383 with interquartile range. Images associated with negative smear microscopy status have lower
384 Timika Scores while those with positive statuses (1 to 9 in 100, 1+, and higher) show
385 progressively higher Timika Scores that appear to plateau around 2+ or higher. The following
386 number of Timika Scores derived from radiologist observations of images are available for
387 Negative, 1 to 9 in 100 (1–9/100), 10 to 99 in 100 (1+), 1 to 9 in 1 (2+), 10 to 99 in 1 (3+), and
388 More than 99 in 1 (4+) groups respectively: 259, 29, 144, 60, 58, 22.
389

390 **Inferential statistics associated with sputum microscopy status**

391 Given the association of Timika Score with sputum smear microscopy, we continued our
392 investigation by assessing the risk of a positive sputum microscopy status (1 to 9 in 100, 1+, or
393 higher) compared to a negative status using Timika score along with the other derived features
394 from radiologist observations. To do so, we performed univariate and multivariate logistic
395 regression removing any feature with no variance that caused univariate or multivariate modeling
396 to fail. We leveraged MRMR feature selection to select the top 5 features for multivariate
397 models. Timika Score is derived from the radiological features (e.g., presence of cavity and
398 overall abnormal volume) so we wanted to select additional features that would not directly
399 correlate with Timika Score but still potentially correlate with sputum microscopy status.

400 We observe multiple features with evidence of involvement of both lungs showing higher
401 risk of sputum microscopy positivity including calcified nodules, fibrotic nodules, low density
402 nodules, involvement of both lungs by indication of any type of sextant feature, medium density
403 nodules, medium cavities, small cavities, multiple cavities, and small nodules (Table 1). For
404 numeric variables, large cavities, low density nodules, medium cavities, small cavities, and
405 overall percent of abnormal volume showed statistically significant increases in risk of positive
406 sputum result (active pathogen detected in the sputum) per each unit increase in percentage
407 whereas pleural effusion percent of hemithorax involved showed the opposite. Each unit
408 increase in Timika Score showed an increased risk in pre-treatment sputum positive microscopy
409 status consistent with its prior reported role. In the multivariate model, the Timika Score showed
410 a higher risk of positive sputum microscopy status after adjusting for indication of involvement
411 of both lungs. This suggests that risk of positive sputum microscopy status does not require
412 evidence of both lung involvement but rather greater percentage abnormal regions or cavity area

413 may be sufficient for the increased risk indicated by Timika Score. Interestingly, the other
 414 MRMR selected features of the multivariate model were all indicators of involvement of both
 415 lungs for the respective features. Only the indication of calcified nodules in both lungs
 416 demonstrated a statistically significant increase in risk of positive sputum microscopy status
 417 adjusting for other covariates in the multivariate model. This feature could suggest cases with
 418 unreported prior history of pulmonary TB.

419 **Table 1. Risk of positive sputum status (1+ or higher) by univariate or multivariate logistic**
 420 **regression analysis on derived features from radiologist observations of CXRs within two**
 421 **weeks of specimens taken prior to treatment.**
 422

Dependent: event2		Negative	Positive	OR (univariable)	OR (multivariable)
aremediastinallymphnodespresent	no	233 (44.6)	289 (55.4)	-	-
	yes	26 (52.0)	24 (48.0)	0.74 (0.41-1.33, p=0.319)	-
both_calcnod	no	241 (47.2)	270 (52.8)	-	-
	yes	18 (29.5)	43 (70.5)	2.13 (1.22-3.88, p=0.010)	2.14 (1.14-4.13, p=0.020)
both_calcsequella1	no	254 (45.9)	299 (54.1)	-	-
	yes	5 (26.3)	14 (73.7)	2.38 (0.90-7.44, p=0.101)	-
both_clustnod	no	229 (44.9)	281 (55.1)	-	-
	yes	30 (48.4)	32 (51.6)	0.87 (0.51-1.48, p=0.603)	-
both_collapse1	no	259 (45.4)	311 (54.6)	-	-
	yes	0 (0.0)	2 (100.0)	1764014.84 (0.00-NA, p=0.982)	1485353.66 (0.00-NA, p=0.981)
both_fibroticnodule1	no	201 (48.4)	214 (51.6)	-	-
	yes	58 (36.9)	99 (63.1)	1.60 (1.10-2.35, p=0.014)	-
both_highden1	no	257 (4	304 (5	-	-

		5.8)	4.2)		
	yes	2 (18.2)	9 (81.8)	3.80 (0.97-25.10, p=0.089)	-
both_isanylargecavitymult	no	259 (45.4)	311 (54.6)	-	-
	yes	0 (0.0)	2 (100.0)	1764014.84 (0.00-NA, p=0.982)	-
both_largecavity1	no	259 (45.4)	312 (54.6)	-	-
	yes	0 (0.0)	1 (100.0)	646864.01 (0.00-NA, p=0.980)	1056464.82 (0.00-NA, p=0.987)
both_largenodule1	no	257 (45.5)	308 (54.5)	-	-
	yes	2 (28.6)	5 (71.4)	2.09 (0.45-14.65, p=0.382)	-
both_lowden1	no	231 (47.4)	256 (52.6)	-	-
	yes	28 (32.9)	57 (67.1)	1.84 (1.14-3.02, p=0.014)	-
both_lowgroundglassdensity activefreshnodules1	no	193 (46.1)	226 (53.9)	-	-
	yes	66 (43.1)	87 (56.9)	1.13 (0.78-1.64, p=0.534)	-
both_lungs	no	149 (52.3)	136 (47.7)	-	-
	yes	110 (38.3)	177 (61.7)	1.76 (1.27-2.46, p=0.001)	0.87 (0.57-1.31, p=0.504)
both_medden1	no	231 (46.9)	262 (53.1)	-	-
	yes	28 (35.4)	51 (64.6)	1.61 (0.99-2.66, p=0.060)	-
both_mediumcavity1	no	258 (45.9)	304 (54.1)	-	-
	yes	1 (10.0)	9 (90.0)	7.64 (1.42-141.33, p=0.055)	-
both_mediumnodule1	no	229 (45.7)	272 (54.3)	-	-
	yes	30 (42.3)	41 (57.7)	1.15 (0.70-1.92, p=0.584)	-
both_multiplecavitiesbeseen	no	256 (46.2)	298 (53.8)	-	-
	yes	3 (16.7)	15 (83.3)	4.30 (1.40-18.69, p=0.022)	-

both_multnod	no	220 (4 6.9)	249 (5 3.1)	-	-
	yes	39 (37 .9)	64 (62 .1)	1.45 (0.94- 2.26, p=0.096)	-
both_noncalcnod	no	200 (4 7.3)	223 (5 2.7)	-	-
	yes	59 (39 .6)	90 (60 .4)	1.37 (0.94- 2.01, p=0.106)	-
both_smallcavity1	no	253 (4 6.9)	287 (5 3.1)	-	-
	yes	6 (18. 8)	26 (81 .2)	3.82 (1.65- 10.39, p=0.004)	-
both_smallnodule1	no	180 (4 8.9)	188 (5 1.1)	-	-
	yes	79 (38 .7)	125 (6 1.3)	1.51 (1.07- 2.15, p=0.019)	-
ispleuraleffusionbilateral	no	257 (4 5.2)	312 (5 4.8)	-	-
	yes	2 (66. 7)	1 (33. 3)	0.41 (0.02- 4.32, p=0.470)	-
mean_calcsequella	Mean (SD)	0.1 (0. 3)	0.2 (0. 8)	1.20 (0.90- 1.83, p=0.316)	-
mean_collapse	Mean (SD)	0.3 (3. 2)	0.3 (2. 1)	1.00 (0.93- 1.07, p=0.938)	-
mean_fibroticnodule	Mean (SD)	2.2 (3. 6)	2.4 (3. 8)	1.02 (0.97- 1.06, p=0.505)	-
mean_highden	Mean (SD)	0.2 (0. 8)	0.5 (3. 1)	1.08 (0.98- 1.27, p=0.225)	-
mean_hugenodule	Mean (SD)	0.0 (0. 2)	0.0 (0. 2)	0.87 (0.35- 2.04, p=0.729)	-
mean_largecavity	Mean (SD)	0.1 (1. 2)	0.6 (2. 6)	1.18 (1.05- 1.38, p=0.015)	-
mean_largenodule	Mean (SD)	0.2 (0. 6)	0.1 (0. 8)	0.93 (0.71- 1.18, p=0.532)	-
mean_lowden	Mean (SD)	2.7 (4. 3)	4.6 (6. 6)	1.07 (1.04- 1.11, p<0.001)	-
mean_lowgroundglassdensit yactivefreshnodules	Mean (SD)	4.7 (8. 7)	4.6 (7. 4)	1.00 (0.98- 1.02, p=0.948)	-
mean_medden	Mean (SD)	2.3 (5. 4)	3.0 (5. 9)	1.02 (0.99- 1.06, p=0.150)	-
mean_mediumcavity	Mean (SD)	0.3 (1. 2)	1.1 (2. 7)	1.28 (1.15- 1.46, p<0.001)	-
mean_mediumnodule	Mean	1.9 (4.)	1.4 (3.)	0.96 (0.92-	-

	(SD)	0)	2)	1.01, p=0.107)	
mean_smallcavity	Mean (SD)	0.4 (1.1)	1.2 (2.4)	1.40 (1.23-1.61, p<0.001)	-
mean_smallnodule	Mean (SD)	4.9 (8.3)	5.7 (8.0)	1.01 (0.99-1.03, p=0.265)	-
othernontbabnormalities	no	217 (45.5)	260 (54.5)	-	-
	yes	42 (44.2)	53 (55.8)	1.05 (0.68-1.65, p=0.819)	-
overall_timika	Mean (SD)	30.9 (26.6)	47.9 (30.2)	1.02 (1.01-1.03, p<0.001)	1.02 (1.01-1.03, p<0.001)
overallpercentofabnormalvolume	Mean (SD)	22.1 (18.2)	26.6 (18.7)	1.01 (1.00-1.02, p=0.004)	-
pleuraleffusionpercentofhemithoraxinvolved	Mean (SD)	1.9 (7.5)	0.8 (4.1)	0.96 (0.93-0.99, p=0.029)	-

423

424 Odds ratios and adjusted odds ratios from univariate and multivariate logistic regression analysis
 425 on derived features from radiologist observations of images taken within two weeks of
 426 specimens prior to treatment start. Shown are the individual derived features from radiologist
 427 observations along with the summary statistics across the prediction categories of Negative or
 428 Positive (1 to 9 in 100, 1+, or higher sputum smear microscopy status). Odds ratios with the 95%
 429 confidence intervals with unadjusted P-values are shown for each derived feature for the
 430 univariate and if applicable, multivariate models. The - sign shows reference categories in the
 431 univariate Odds Ratio column or reference as well as excluded variables in the multivariate Odds
 432 Ratio column.
 433

434 **Assessing predictive capacity of machine learning models**

435 The TB portals offers a variety of radiologist observations of CXRs that may provide
 436 additional information towards the prediction of baseline sputum smear microscopy status.
 437 Therefore, we investigated the additional features comparing predictive performance to Timika
 438 Score alone. By doing so, we sought to identify any additional features that might improve upon
 439 Timika Score performance and to evaluate how well Timika Score itself could predict sputum
 440 status when compared to various feature selection or dimensionality reduction techniques that
 441 summarize the radiological features within the data.

442 **Comparison of predicting 2+ versus 1+ sputum smear status in training and** 443 **validation sets**

444 We noted that sputum smear scores of 2+ or greater showed a higher mean Timika Score
445 so we tried two predictive approaches: task one involved predicting positive (1 to 9 in 100, 1+, or
446 greater smear status indicating any active pathogen in the sputum) versus negative smear status
447 while task two involved predicting higher bacterial load positive (2+, 3+, or greater smear status)
448 versus negative sputum status. We split the data into a 70% training set and 30% validation set
449 for the model training and validation. All pipelines were created and run using the MLR3
450 package which allows for unbiased assessment of model performance by encapsulating the pre-
451 processing steps within the cross-validation approach. All prediction tasks included featureless
452 pipelines showing a non-informative model that only predicts the most prevalent class or
453 randomly selects a class in case of a tie. Thus, the featureless model can be considered a control
454 and predictive models should perform significantly better than this featureless control model.

455 In the first prediction problem attempting to discriminate positive from negative status, 5-
456 fold cross validation results on the training data showed that most models demonstrated
457 relatively similar predictive performance across pipelines (Table 2). Pipelines using top 5
458 components by principal component analysis (which captured ~ 50% of variability in the dataset)
459 tended to show slightly decreased performance in general. Since this prediction problem did not
460 have a large class imbalance, the addition of a class balancing step did not make a significant
461 impact to prediction performance for the pipelines tested. Of note, pipelines only including the
462 Timika Score showed equivalent performance in general to workflows using the top 5 predictive
463 features from various feature selection algorithms. There is a slight decrease in performance of
464 Timika-only models to the best top 5 feature selection model, which reflect that additional

465 features may provide some minimal improvement over Timika Score. To test predictive
 466 performance on data which the models had not seen before, we trained the above models on the
 467 entire training set and tested on 30% of held out validation data (Table 3). We observed similar
 468 findings to the cross-validated results we obtained from the training data.

469 **Table 2. Comparison of machine learning pipeline performance via 5-fold cross-validation**
 470 **for predicting positive (1 to 9 in 100, 1+, or higher) versus negative sputum microscopy**
 471 **status on training data.**

pipeline	classif.au c	classif.mc c	classif.sensit ivity	classif.specif icity
ffact.cb.enc.zv.num_scale.flt.auc.classif.m ultinom	0.69 +/- 0.02	0.37 +/- 0.02	0.73 +/- 0.05	0.63 +/- 0.05
ffact.cb.enc.zv.num_scale.flt.auc.classif.lo g_reg	0.7 +/- 0.02	0.34 +/- 0.05	0.78 +/- 0.03	0.61 +/- 0.03
ffact.cb.enc.zv.num_scale.flt.njmim.classi f.ranger	0.69 +/- 0.03	0.33 +/- 0.03	0.69 +/- 0.04	0.66 +/- 0.09
timika.num_scale.classif.rpart	0.7 +/- 0.01	0.32 +/- 0.05	0.51 +/- 0.23	0.75 +/- 0.08
timika.num_scale.cb.classif.rpart	0.67 +/- 0.04	0.32 +/- 0.03	0.62 +/- 0.07	0.65 +/- 0.08
ffact.cb.enc.zv.num_scale.flt.auc.classif.c v_glmnet	0.68 +/- 0.03	0.32 +/- 0.02	0.72 +/- 0.08	0.6 +/- 0.05
ffact.cb.enc.zv.num_scale.flt.auc.classif.gl mnet	0.69 +/- 0.01	0.31 +/- 0.06	0.78 +/- 0.04	0.61 +/- 0.02
ffact.cb.enc.zv.num_scale.flt.auc.classif.n net	0.66 +/- 0.06	0.31 +/- 0.04	0.78 +/- 0.01	0.57 +/- 0.05
ffact.enc.zv.num_scale.flt.njmim.classif.lo g_reg	0.7 +/- 0.01	0.31 +/- 0.04	0.75 +/- 0	0.6 +/- 0.08
ffact.enc.zv.num_scale.flt.njmim.classif.m ultinom	0.7 +/- 0.01	0.31 +/- 0.04	0.75 +/- 0	0.6 +/- 0.08
ffact.enc.zv.num_scale.flt.njmim.classif.ra nger	0.68 +/- 0.04	0.31 +/- 0.02	0.53 +/- 0.12	0.74 +/- 0.04
ffact.cb.enc.zv.num_scale.flt.njmim.classi f.svm	0.68 +/- 0.01	0.31 +/- 0	0.68 +/- 0.05	0.63 +/- 0.05
ffact.cb.enc.zv.num_scale.flt.auc.classif.sv m	0.66 +/- 0.01	0.3 +/- 0.04	0.73 +/- 0.05	0.51 +/- 0.08
ffact.cb.enc.zv.num_scale.flt.mrmr.classif. cv_glmnet	0.67 +/- 0.02	0.3 +/- 0.04	0.69 +/- 0.08	0.6 +/- 0.05
ffact.cb.enc.zv.num_scale.flt.njmim.classi f.log_reg	0.65 +/- 0.03	0.3 +/- 0.04	0.76 +/- 0.07	0.56 +/- 0.02
ffact.cb.enc.zv.num_scale.flt.njmim.classi	0.69 +/-	0.3 +/-	0.81 +/- 0.01	0.42 +/- 0.07

f.naive_bayes	0.05	0.04		
ffact.cb.enc.zv.num_scale.flt.njmim.classi	0.69 +/-	0.3 +/-	0.75 +/- 0.04	0.55 +/- 0.05
f.multinom	0.02	0.03		
timika.num_scale.cb.classif.log_reg	0.67 +/-	0.29 +/-	0.69 +/- 0.11	0.6 +/- 0.05
	0.02	0.03		
timika.num_scale.cb.classif.multinom	0.67 +/-	0.29 +/-	0.69 +/- 0.11	0.6 +/- 0.05
	0.02	0.03		
timika.num_scale.classif.naive_bayes	0.67 +/-	0.29 +/-	0.69 +/- 0.12	0.6 +/- 0.08
	0.02	0.03		
ffact.cb.enc.zv.num_scale.pca.classif.nnet	0.62 +/-	0.28 +/-	0.58 +/- 0.08	0.77 +/- 0.11
	0.09	0.09		
ffact.cb.enc.zv.num_scale.flt.auc.classif.n	0.71 +/-	0.28 +/-	0.83 +/- 0.08	0.39 +/- 0.06
aive_bayes	0.04	0.06		
ffact.cb.enc.zv.num_scale.flt.auc.classif.r	0.64 +/-	0.28 +/-	0.68 +/- 0.19	0.6 +/- 0.21
part	0.06	0.04		
ffact.enc.zv.num_scale.flt.njmim.classif.sv	0.66 +/-	0.27 +/-	0.56 +/- 0.02	0.68 +/- 0.05
m	0.01	0.09		
ffact.cb.enc.zv.num_scale.pca.classif.svm	0.67 +/-	0.27 +/-	0.59 +/- 0.02	0.6 +/- 0.17
	0.03	0.08		
timika.num_scale.classif.kknn	0.64 +/-	0.27 +/-	0.75 +/- 0.08	0.5 +/- 0.02
	0.03	0.08		
ffact.cb.enc.zv.num_scale.pca.classif.mul	0.71 +/-	0.27 +/-	0.76 +/- 0.07	0.49 +/- 0
tinom	0.02	0.05		
timika.num_scale.classif.log_reg	0.67 +/-	0.27 +/-	0.56 +/- 0.08	0.66 +/- 0.03
	0.02	0.05		
timika.num_scale.classif.multinom	0.67 +/-	0.27 +/-	0.56 +/- 0.08	0.66 +/- 0.03
	0.02	0.05		
ffact.cb.enc.zv.num_scale.flt.mrmr.classif.	0.66 +/-	0.27 +/-	0.68 +/- 0.05	0.58 +/- 0.07
xgboost	0.03	0.03		
ffact.cb.enc.zv.num_scale.flt.mrmr.classif.	0.66 +/-	0.27 +/-	0.72 +/- 0.08	0.55 +/- 0.08
ranger	0.01	0.02		
timika.num_scale.cb.classif.naive_bayes	0.67 +/-	0.27 +/-	0.72 +/- 0.04	0.56 +/- 0.03
	0.02	0.02		
timika.num_scale.classif.xgboost	0.66 +/-	0.27 +/-	0.61 +/- 0.08	0.63 +/- 0.02
	0.04	0.01		
ffact.cb.enc.zv.num_scale.flt.auc.classif.ra	0.61 +/-	0.26 +/-	0.61 +/- 0.08	0.63 +/- 0.03
nger	0.02	0.17		
ffact.cb.enc.zv.num_scale.flt.mrmr.classif.	0.65 +/-	0.26 +/-	0.69 +/- 0.08	0.6 +/- 0.1
svm	0.02	0.07		
timika.num_scale.classif.svm	0.65 +/-	0.26 +/-	0.61 +/- 0.12	0.65 +/- 0.07
	0.03	0.07		
ffact.cb.enc.zv.num_scale.flt.mrmr.classif.	0.66 +/-	0.26 +/-	0.69 +/- 0.04	0.6 +/- 0.1
multinom	0.03	0.06		

ffact.enc.zv.num_scale.flt.njmim.classif.glmnet	0.7 +/- 0.01	0.26 +/- 0.05	0.61 +/- 0.08	0.65 +/- 0.12
ffact.cb.enc.zv.num_scale.flt.mrmr.classif.rpart	0.61 +/- 0.04	0.26 +/- 0.04	0.67 +/- 0.12	0.6 +/- 0.07
timika.num_scale.classif.nnet	0.66 +/- 0.02	0.26 +/- 0.04	0.69 +/- 0.08	0.56 +/- 0.01
ffact.cb.enc.zv.num_scale.flt.mrmr.classif.glmnet	0.64 +/- 0.01	0.26 +/- 0.03	0.69 +/- 0.08	0.6 +/- 0.05
ffact.cb.enc.zv.num_scale.flt.njmim.classif.cv_glmnet	0.68 +/- 0.03	0.26 +/- 0.03	0.69 +/- 0.11	0.6 +/- 0.05
timika.num_scale.cb.classif.nnet	0.67 +/- 0.01	0.26 +/- 0.03	0.72 +/- 0.04	0.56 +/- 0.07
timika.num_scale.cb.classif.svm	0.68 +/- 0.01	0.26 +/- 0.03	0.69 +/- 0.08	0.6 +/- 0.05
ffact.enc.zv.num_scale.flt.njmim.classif.naive_bayes	0.69 +/- 0.01	0.25 +/- 0.1	0.78 +/- 0.13	0.55 +/- 0.2
ffact.cb.enc.zv.num_scale.flt.njmim.classif.nnet	0.65 +/- 0.04	0.25 +/- 0.09	0.65 +/- 0.06	0.53 +/- 0.09
ffact.cb.enc.zv.num_scale.pca.classif.log_reg	0.69 +/- 0.08	0.24 +/- 0.06	0.75 +/- 0.08	0.48 +/- 0.02
ffact.cb.enc.zv.num_scale.flt.mrmr.classif.nnet	0.65 +/- 0.02	0.24 +/- 0.05	0.69 +/- 0.08	0.51 +/- 0.03
ffact.enc.zv.num_scale.flt.njmim.classif.cv_glmnet	0.68 +/- 0.03	0.24 +/- 0.04	0.56 +/- 0.04	0.68 +/- 0.2
ffact.cb.enc.zv.num_scale.flt.mrmr.classif.log_reg	0.65 +/- 0.01	0.24 +/- 0.02	0.69 +/- 0.08	0.6 +/- 0.05
ffact.cb.enc.zv.num_scale.pca.classif.kknn	0.65 +/- 0.03	0.24 +/- 0.02	0.61 +/- 0.04	0.63 +/- 0.03
ffact.cb.enc.zv.num_scale.flt.njmim.classif.kknn	0.64 +/- 0.05	0.23 +/- 0.08	0.64 +/- 0.16	0.53 +/- 0.12
ffact.cb.enc.zv.num_scale.flt.njmim.classif.glmnet	0.66 +/- 0.02	0.23 +/- 0.02	0.64 +/- 0.08	0.58 +/- 0.08
ffact.enc.zv.num_scale.flt.njmim.classif.rpart	0.65 +/- 0.06	0.22 +/- 0.13	0.44 +/- 0.08	0.73 +/- 0.06
ffact.enc.zv.num_scale.flt.njmim.classif.kknn	0.64 +/- 0.06	0.22 +/- 0.11	0.67 +/- 0.04	0.58 +/- 0.05
ffact.cb.enc.zv.num_scale.pca.classif.ranger	0.62 +/- 0.05	0.22 +/- 0.1	0.56 +/- 0.12	0.63 +/- 0.08
timika.num_scale.cb.classif.kknn	0.63 +/- 0.03	0.21 +/- 0.04	0.73 +/- 0.01	0.44 +/- 0.07
ffact.cb.enc.zv.num_scale.flt.njmim.classif.xgboost	0.67 +/- 0.05	0.2 +/- 0.19	0.67 +/- 0.09	0.56 +/- 0.12
ffact.cb.enc.zv.num_scale.pca.classif.glm	0.7 +/- 0.01	0.2 +/- 0.01	0.75 +/- 0.04	0.49 +/- 0.02

net	0.08	0.11		
timika.num_scale.classif.ranger	0.64 +/- 0.06	0.2 +/- 0.11	0.59 +/- 0.11	0.63 +/- 0.05
timika.num_scale.cb.classif.ranger	0.6 +/- 0.09	0.19 +/- 0.13	0.65 +/- 0.07	0.55 +/- 0.03
ffact.cb.enc.zv.num_scale.flt.auc.classif.xgboost	0.6 +/- 0.01	0.19 +/- 0.03	0.7 +/- 0.03	0.51 +/- 0.07
ffact.enc.zv.num_scale.flt.njmim.classif.xgboost	0.61 +/- 0.05	0.18 +/- 0.05	0.57 +/- 0.14	0.63 +/- 0.05
timika.num_scale.cb.classif.xgboost	0.62 +/- 0.06	0.17 +/- 0.12	0.57 +/- 0.06	0.53 +/- 0.02
ffact.cb.enc.zv.num_scale.flt.njmim.classif.rpart	0.63 +/- 0.03	0.16 +/- 0.14	0.53 +/- 0.12	0.63 +/- 0.14
ffact.cb.enc.zv.num_scale.flt.mrmr.classif.kknn	0.58 +/- 0.02	0.16 +/- 0.03	0.75 +/- 0.08	0.37 +/- 0.1
ffact.enc.zv.num_scale.flt.njmim.classif.net	0.66 +/- 0.07	0.15 +/- 0.11	0.49 +/- 0.1	0.7 +/- 0.02
ffact.cb.enc.zv.num_scale.pca.classif.xgboost	0.62 +/- 0.04	0.13 +/- 0.16	0.53 +/- 0.04	0.58 +/- 0.02
ffact.cb.enc.zv.num_scale.flt.auc.classif.kknn	0.61 +/- 0.07	0.12 +/- 0.08	0.58 +/- 0.04	0.51 +/- 0.03
ffact.cb.enc.zv.num_scale.pca.classif.rpart	0.57 +/- 0	0.12 +/- 0.08	0.53 +/- 0.04	0.58 +/- 0.03
ffact.cb.enc.zv.num_scale.pca.classif.naive_bayes	0.55 +/- 0.06	0.1 +/- 0.18	0.86 +/- 0.04	0.23 +/- 0.06
ffact.cb.enc.zv.num_scale.flt.mrmr.classif.naive_bayes	0.65 +/- 0.01	0.1 +/- 0.15	1 +/- 0	0.02 +/- 0.03
ffact.cb.enc.zv.num_scale.pca.classif.cv_glmnet	0.59 +/- 0.13	0.09 +/- 0.17	0.64 +/- 0.16	0.41 +/- 0.05
ffact.cb.enc.zv.num_scale.pca.classif.featureless	0.5 +/- 0	0.09 +/- 0.02	0.57 +/- 0.02	0.52 +/- 0.05
ffact.cb.enc.zv.num_scale.flt.mrmr.classif.featureless	0.5 +/- 0	0.05 +/- 0.09	0.56 +/- 0.1	0.53 +/- 0.05
ffact.cb.enc.zv.num_scale.flt.njmim.classif.featureless	0.5 +/- 0	0.03 +/- 0.09	0.53 +/- 0.08	0.56 +/- 0.1
ffact.enc.zv.num_scale.flt.njmim.classif.featureless	0.5 +/- 0	0 +/- 0	0 +/- 0	1 +/- 0
timika.num_scale.classif.featureless	0.5 +/- 0	0 +/- 0	0 +/- 0	1 +/- 0
timika.num_scale.cb.classif.featureless	0.5 +/- 0	-0.12 +/- 0.04	0.36 +/- 0.12	0.52 +/- 0.09
ffact.cb.enc.zv.num_scale.flt.auc.classif.featureless	0.5 +/- 0	-0.02 +/- 0.03	0.43 +/- 0.06	0.45 +/- 0.05

473 The machine learning pipeline is shown in the pipeline column. Model performance metrics
 474 include Area under the curve (classif.auc), Balanced accuracy (classif.bacc), Matthew's
 475 Correlation Coefficient (classif.mcc), Sensitivity (classif.sensitivity), and Specificity
 476 (classif.specifcity). Each cell represents the median metric +/- the MAD for the 5-fold cross
 477 validation testing on the training data representing 70% of the entire dataset. Pipelines using only
 478 the Timika Score for prediction start with Timika in the pipeline name. Each pipeline shows all
 479 steps used in the pipeline ending with the machine learning algorithm used and are ordered by
 480 classif.mcc.
 481

482 **Table 3. Comparison of machine learning pipeline performance for predicting positive (1 to**
 483 **9 in 100, 1+, or higher) versus negative sputum microscopy status on validation data.**
 484

pipeline	classif.auc	classif.mcc	classif.sensitivity	classif.specifcity
timika.num_scale.classif.rpart	0.66905 26	0.3686316 25	0.65333333	0.71578947
timika.num_scale.cb.classif.rpart	0.66905 26	0.3686316 25	0.65333333	0.71578947
ffact.cb.enc.zv.num_scale.flt.mrmr.classif.log_reg	0.66533 33	0.3426355 61	0.73333333	0.61052632
ffact.cb.enc.zv.num_scale.flt.njmim.classif.multinom	0.67747 37	0.3300229 46	0.76	0.56842105
timika.num_scale.classif.naive_bayes	0.65649 12	0.3290517 73	0.72	0.61052632
timika.num_scale.cb.classif.log_reg	0.65649 12	0.3290517 73	0.72	0.61052632
timika.num_scale.cb.classif.multinom	0.65649 12	0.3290517 73	0.72	0.61052632
timika.num_scale.cb.classif.svm	0.65494 74	0.3290517 73	0.72	0.61052632
ffact.cb.enc.zv.num_scale.flt.mrmr.classif.cv_glmnet	0.65649 12	0.3290517 73	0.72	0.61052632
ffact.cb.enc.zv.num_scale.flt.mrmr.classif.rpart	0.66526 32	0.3290517 73	0.72	0.61052632
ffact.cb.enc.zv.num_scale.flt.mrmr.classif.svm	0.65382 46	0.3290517 73	0.72	0.61052632
ffact.cb.enc.zv.num_scale.flt.mrmr.classif.ranger	0.67389 47	0.3260949 51	0.7466667	0.57894737
ffact.cb.enc.zv.num_scale.flt.mrmr.classif.multinom	0.65080 7	0.3200582 31	0.76	0.55789474
timika.num_scale.classif.nnet	0.65649 12	0.3188608 49	0.72	0.6
ffact.cb.enc.zv.num_scale.flt.njmim.classif.lo	0.68343	0.3188608	0.72	0.6

g_reg	86	49		
ffact.enc.zv.num_scale.flt.njmim.classif.m ultinom	0.68343 86	0.3188608 49	0.72	0.6
ffact.cb.enc.zv.num_scale.flt.mrmr.classif. glmnet	0.64807 02	0.3160651 51	0.7466667	0.56842105
ffact.cb.enc.zv.num_scale.flt.njmim.classif. glmnet	0.68414 04	0.3101006 39	0.76	0.54736842
ffact.cb.enc.zv.num_scale.flt.mrmr.classif. nnet	0.65396 49	0.3101006 39	0.76	0.54736842
ffact.cb.enc.zv.num_scale.flt.auc.classif.lo g_reg	0.68315 79	0.3052985 67	0.7066667	0.6
ffact.cb.enc.zv.num_scale.flt.auc.classif.sv m	0.64322 81	0.3052985 67	0.7066667	0.6
ffact.cb.enc.zv.num_scale.flt.auc.classif.gl mnet	0.68526 32	0.3020841 6	0.6933333	0.61052632
ffact.cb.enc.zv.num_scale.flt.auc.classif.rp art	0.69536 84	0.3020841 6	0.6933333	0.61052632
timika.num_scale.cb.classif.xgboost	0.66764 91	0.2985495 67	0.72	0.57894737
ffact.cb.enc.zv.num_scale.flt.auc.classif.m ultinom	0.68154 39	0.2950861 36	0.7066667	0.58947368
ffact.enc.zv.num_scale.flt.njmim.classif.nn et	0.68336 84	0.2901905	0.56	0.72631579
ffact.cb.enc.zv.num_scale.flt.auc.classif.cv glmnet	0.67214 04	0.2884210 53	0.72	0.56842105
ffact.cb.enc.zv.num_scale.flt.njmim.classif. cv_glmnet	0.66701 75	0.2884210 53	0.72	0.56842105
ffact.cb.enc.zv.num_scale.flt.njmim.classif. log_reg	0.65985 96	0.2820708 94	0.7333333	0.54736842
ffact.cb.enc.zv.num_scale.flt.njmim.classif. svm	0.66940 35	0.2820708 94	0.7333333	0.54736842
ffact.cb.enc.zv.num_scale.flt.njmim.classif. kknn	0.66722 81	0.2783051 73	0.72	0.55789474
timika.num_scale.classif.ranger	0.65319 3	0.2747216 68	0.7066667	0.56842105
ffact.cb.enc.zv.num_scale.flt.auc.classif.ra nger	0.66694 74	0.2747216 68	0.7066667	0.56842105
ffact.cb.enc.zv.num_scale.pca.classif.rang er	0.67024 56	0.2712772 65	0.5866667	0.68421053
ffact.enc.zv.num_scale.flt.njmim.classif.na ive_bayes	0.67838 6	0.2645800 31	0.8	0.45263158
ffact.cb.enc.zv.num_scale.pca.classif.log_r eg	0.66631 58	0.2619571 6	0.7333333	0.52631579

timika.num_scale.classif.xgboost	0.63726 32	0.2610956 07	0.6933333	0.56842105
ffact.enc.zv.num_scale.flt.njmim.classif.kknn	0.67607 02	0.2602752 47	0.76	0.49473684
timika.num_scale.cb.classif.naive_bayes	0.65649 12	0.2580948 68	0.72	0.53684211
timika.num_scale.classif.svm	0.63221 05	0.2517322 49	0.6	0.65263158
ffact.cb.enc.zv.num_scale.flt.njmim.classif.naive_bayes	0.65656 14	0.2499759 12	0.8133333	0.42105263
timika.num_scale.cb.classif.nnet	0.65649 12	0.2479921 39	0.72	0.52631579
ffact.cb.enc.zv.num_scale.flt.njmim.classif.nnet	0.64463 16	0.2460671 29	0.6266667	0.62105263
ffact.cb.enc.zv.num_scale.flt.njmim.classif.rpart	0.62301 75	0.2454016 78	0.56	0.68421053
ffact.enc.zv.num_scale.flt.njmim.classif.rpart	0.68603 51	0.2388899 24	0.2933333	0.89473684
timika.num_scale.classif.log_reg	0.65649 12	0.2386959 84	0.5866667	0.65263158
timika.num_scale.classif.multinom	0.65649 12	0.2386959 84	0.5866667	0.65263158
ffact.enc.zv.num_scale.flt.njmim.classif.cv_glmnet	0.67614 04	0.2386959 84	0.5866667	0.65263158
ffact.enc.zv.num_scale.flt.njmim.classif.svm	0.67319 3	0.2386959 84	0.5866667	0.65263158
timika.num_scale.cb.classif.ranger	0.66084 21	0.2372960 08	0.68	0.55789474
ffact.cb.enc.zv.num_scale.flt.auc.classif.kknn	0.62736 84	0.2340389 47	0.6666667	0.56842105
ffact.cb.enc.zv.num_scale.flt.mrmr.classif.xgboost	0.64238 6	0.2324484 72	0.5466667	0.68421053
timika.num_scale.classif.kknn	0.66477 19	0.2313835 79	0.84	0.36842105
ffact.enc.zv.num_scale.flt.njmim.classif.glmnet	0.68470 18	0.2303696 5	0.6	0.63157895
ffact.cb.enc.zv.num_scale.flt.auc.classif.xgboost	0.63298 25	0.2270658 68	0.68	0.54736842
timika.num_scale.cb.classif.kknn	0.65410 53	0.2254002 87	0.8266667	0.37894737
ffact.cb.enc.zv.num_scale.flt.auc.classif.nnet	0.65578 95	0.2237619 99	0.6666667	0.55789474
ffact.cb.enc.zv.num_scale.pca.classif.xgbo	0.61922	0.2194776	0.5333333	0.68421053

ost	81	21		
ffact.cb.enc.zv.num_scale.flt.auc.classif.naive_bayes	0.6817544	0.219443561	0.7866667	0.42105263
ffact.enc.zv.num_scale.flt.njmim.classif.ranger	0.6828772	0.202914311	0.49333333	0.70526316
ffact.cb.enc.zv.num_scale.flt.njmim.classif.ranger	0.6682807	0.197449626	0.53333333	0.66315789
ffact.cb.enc.zv.num_scale.pca.classif.svm	0.678386	0.197449626	0.53333333	0.66315789
ffact.cb.enc.zv.num_scale.flt.njmim.classif.xgboost	0.6248421	0.193460333	0.5066667	0.68421053
ffact.cb.enc.zv.num_scale.pca.classif.multinom	0.6685614	0.18616679	0.68	0.50526316
ffact.cb.enc.zv.num_scale.pca.classif.glmnet	0.6400702	0.183968835	0.77333333	0.4
ffact.cb.enc.zv.num_scale.pca.classif.rpart	0.6038596	0.182617134	0.41333333	0.75789474
ffact.cb.enc.zv.num_scale.flt.mrmr.classif.kknn	0.7187368	0.179349807	0.84	0.31578947
ffact.cb.enc.zv.num_scale.pca.classif.kknn	0.6215439	0.169277792	0.49333333	0.67368421
ffact.cb.enc.zv.num_scale.pca.classif.naive_bayes	0.6355789	0.163928656	0.85333333	0.28421053
ffact.cb.enc.zv.num_scale.pca.classif.nnet	0.6329123	0.155776562	0.3866667	0.75789474
ffact.cb.enc.zv.num_scale.pca.classif.cv_glmnet	0.6386667	0.139639461	0.73333333	0.4
ffact.cb.enc.zv.num_scale.flt.mrmr.classif.featureless	0.5	0.130316753	0.57333333	0.55789474
ffact.enc.zv.num_scale.flt.njmim.classif.xgboost	0.6146667	0.10097508	0.41333333	0.68421053
ffact.cb.enc.zv.num_scale.flt.njmim.classif.featureless	0.5	0.091519708	0.5866667	0.50526316
ffact.cb.enc.zv.num_scale.pca.classif.featureless	0.5	0.059234888	0.53333333	0.52631579
timika.num_scale.classif.featureless	0.5	0	0	1
ffact.enc.zv.num_scale.flt.njmim.classif.featureless	0.5	0	0	1
ffact.cb.enc.zv.num_scale.flt.auc.classif.featureless	0.5	-0.000699988	0.5466667	0.45263158
timika.num_scale.cb.classif.featureless	0.5	-0.0097895	0.45333333	0.53684211

		88		
ffact.cb.enc.zv.num_scale.flt.mrmr.classif.	0.65789	-	0.9866667	0.01052632
naive_bayes	47	0.0129261		
		12		

485

486 The machine learning pipeline is shown in the pipeline column. Model performance metrics
487 include Area under the curve (classif.auc), Balanced accuracy (classif.bacc), Matthew's
488 Correlation Coefficient (classif.mcc), Sensitivity (classif.sensitivity), and Specificity
489 (classif.specificity). Each cell represents the metric for the performance on the 30% held-out
490 validation test set after training on the 70% training data. Pipelines using only the Timika Score
491 for prediction start with Timika in the pipeline name. Each pipeline shows all steps used in the
492 pipeline ending with the machine learning algorithm used and are ordered by classif.mcc.
493

494 We hypothesized that the second prediction task might demonstrate a performance boost
495 in predictive power since the 2+ sputum status or higher (very high pathogen load in the sputum)
496 showed a larger difference in Timika Score compared to negative. For this prediction task, we
497 removed the borderline 1 to 9 in 100 and 1+ sputum test results from the analysis. 5-fold cross
498 validation results on training set confirmed our hypothesis as we observed increases in
499 performance for reported metrics (Table 4) for both top 5 features pipelines as well as Timika
500 Score only workflows. In general, we observed equivalent performance from Timika Score only
501 pipelines to workflows using the top 5 predictive features from various feature selection
502 algorithms. The benchmarking suggests that while possible to achieve additional gains from the
503 set of derived radiologist observations, these would likely be minimal. Interestingly, though this
504 prediction problem shows a moderate class imbalance, the incorporation of class balancing did
505 not significantly increase or impair performance for the pipelines. When we tested models
506 trained on the entire training set on a validation set of 30% held out data, we saw similar
507 predictive performance to our observation of the 5-fold cross-validated results on the training set
508 (Table 5).

509 **Table 4. Comparison of machine learning pipeline performance via 5-fold cross-validation**
 510 **for predicting high bacterial load positive (2+ or higher) versus negative sputum**
 511 **microscopy status on training data.**
 512

pipeline	classif.au c	classif.mc c	classif.sensit ivity	classif.specif icity
timika.num_scale.cb.classif.naive_bayes	0.77 +/- 0.06	0.5 +/- 0.09	0.72 +/- 0.04	0.74 +/- 0.13
ffact.enc.zv.num_scale.flt.njmim.classif.ranger	0.78 +/- 0.07	0.49 +/- 0.09	0.84 +/- 0.05	0.65 +/- 0.15
ffact.cb.enc.zv.num_scale.flt.njmim.classif.glmnet	0.76 +/- 0.06	0.48 +/- 0.1	0.75 +/- 0.04	0.7 +/- 0.07
ffact.cb.enc.zv.num_scale.flt.njmim.classif.multinom	0.78 +/- 0.04	0.48 +/- 0.1	0.81 +/- 0	0.7 +/- 0.07
ffact.cb.enc.zv.num_scale.flt.mrmr.classif.cv_glmnet	0.75 +/- 0.05	0.47 +/- 0.14	0.69 +/- 0.08	0.75 +/- 0.07
ffact.cb.enc.zv.num_scale.flt.njmim.classif.cv_glmnet	0.77 +/- 0.04	0.47 +/- 0.14	0.69 +/- 0.08	0.75 +/- 0.1
timika.num_scale.cb.classif.log_reg	0.77 +/- 0.06	0.47 +/- 0.14	0.69 +/- 0.08	0.75 +/- 0.07
timika.num_scale.cb.classif.multinom	0.77 +/- 0.06	0.47 +/- 0.14	0.69 +/- 0.08	0.75 +/- 0.07
timika.num_scale.cb.classif.svm	0.74 +/- 0.05	0.47 +/- 0.14	0.69 +/- 0.08	0.75 +/- 0.07
timika.num_scale.classif.nnet	0.77 +/- 0.06	0.47 +/- 0.12	0.75 +/- 0.03	0.68 +/- 0.08
ffact.cb.enc.zv.num_scale.flt.auc.classif.glmnet	0.75 +/- 0.08	0.47 +/- 0.06	0.75 +/- 0.04	0.7 +/- 0.15
ffact.cb.enc.zv.num_scale.flt.auc.classif.log_reg	0.75 +/- 0.03	0.44 +/- 0.16	0.81 +/- 0.03	0.65 +/- 0.15
ffact.cb.enc.zv.num_scale.flt.njmim.classif.nnet	0.74 +/- 0.06	0.44 +/- 0.11	0.76 +/- 0.07	0.63 +/- 0.2
ffact.cb.enc.zv.num_scale.flt.auc.classif.nnet	0.77 +/- 0.04	0.44 +/- 0.07	0.72 +/- 0.12	0.74 +/- 0.16
ffact.cb.enc.zv.num_scale.flt.auc.classif.rpart	0.75 +/- 0.09	0.43 +/- 0.23	0.69 +/- 0.04	0.8 +/- 0.15
ffact.cb.enc.zv.num_scale.flt.mrmr.classif.glmnet	0.7 +/- 0.08	0.43 +/- 0.15	0.69 +/- 0.08	0.75 +/- 0.02
ffact.cb.enc.zv.num_scale.flt.auc.classif.multinom	0.75 +/- 0.06	0.42 +/- 0.19	0.75 +/- 0.04	0.7 +/- 0.1
ffact.cb.enc.zv.num_scale.flt.auc.classif.svm	0.73 +/- 0.06	0.42 +/- 0.19	0.75 +/- 0.04	0.7 +/- 0.1
ffact.cb.enc.zv.num_scale.flt.njmim.classif	0.78 +/-	0.42 +/-	0.78 +/- 0.04	0.7 +/- 0.18

f.log_reg	0.04	0.19		
ffact.cb.enc.zv.num_scale.flt.njmim.classif.svm	0.75 +/- 0.04	0.42 +/- 0.16	0.72 +/- 0.08	0.75 +/- 0.07
timika.num_scale.cb.classif.nnet	0.72 +/- 0.07	0.41 +/- 0.1	0.69 +/- 0.04	0.74 +/- 0.09
ffact.cb.enc.zv.num_scale.flt.njmim.classif.ranger	0.78 +/- 0.06	0.41 +/- 0.02	0.72 +/- 0.07	0.68 +/- 0.1
ffact.cb.enc.zv.num_scale.flt.auc.classif.ranger	0.73 +/- 0.11	0.4 +/- 0.17	0.67 +/- 0.03	0.75 +/- 0.07
ffact.cb.enc.zv.num_scale.pca.classif.ranger	0.74 +/- 0.04	0.37 +/- 0.06	0.65 +/- 0.06	0.75 +/- 0.02
ffact.cb.enc.zv.num_scale.flt.auc.classif.naive_bayes	0.74 +/- 0.03	0.36 +/- 0.14	0.81 +/- 0.05	0.5 +/- 0.07
ffact.enc.zv.num_scale.flt.njmim.classif.log_reg	0.75 +/- 0.06	0.36 +/- 0.14	0.89 +/- 0	0.5 +/- 0.07
ffact.enc.zv.num_scale.flt.njmim.classif.multinom	0.75 +/- 0.06	0.36 +/- 0.14	0.89 +/- 0	0.5 +/- 0.07
timika.num_scale.cb.classif.ranger	0.7 +/- 0.08	0.36 +/- 0.08	0.69 +/- 0.04	0.58 +/- 0.16
timika.num_scale.classif.svm	0.75 +/- 0	0.36 +/- 0.06	0.81 +/- 0.04	0.5 +/- 0.12
ffact.cb.enc.zv.num_scale.flt.mrmr.classif.multinom	0.67 +/- 0.13	0.35 +/- 0.26	0.68 +/- 0.05	0.7 +/- 0.13
ffact.cb.enc.zv.num_scale.flt.mrmr.classif.log_reg	0.66 +/- 0.1	0.35 +/- 0.2	0.69 +/- 0.04	0.68 +/- 0.12
ffact.cb.enc.zv.num_scale.flt.auc.classif.cv_glmnet	0.72 +/- 0.07	0.35 +/- 0.19	0.69 +/- 0.08	0.7 +/- 0.07
ffact.enc.zv.num_scale.flt.njmim.classif.svm	0.78 +/- 0.04	0.35 +/- 0.19	0.86 +/- 0.04	0.5 +/- 0.15
ffact.cb.enc.zv.num_scale.flt.mrmr.classif.nnet	0.66 +/- 0.07	0.35 +/- 0.17	0.75 +/- 0.08	0.74 +/- 0.02
ffact.cb.enc.zv.num_scale.flt.njmim.classif.rpart	0.75 +/- 0.02	0.34 +/- 0.17	0.72 +/- 0.04	0.8 +/- 0.06
timika.num_scale.cb.classif.kknn	0.71 +/- 0.08	0.34 +/- 0.1	0.75 +/- 0.07	0.6 +/- 0.11
ffact.enc.zv.num_scale.flt.njmim.classif.glmnet	0.75 +/- 0.06	0.34 +/- 0.09	0.86 +/- 0.03	0.45 +/- 0.04
ffact.enc.zv.num_scale.flt.njmim.classif.rpart	0.7 +/- 0.04	0.34 +/- 0.08	0.78 +/- 0.09	0.47 +/- 0.11
timika.num_scale.classif.ranger	0.68 +/- 0.05	0.34 +/- 0.01	0.78 +/- 0.03	0.53 +/- 0.04
timika.num_scale.cb.classif.xgboost	0.73 +/- 0.03	0.33 +/- 0.11	0.61 +/- 0.04	0.63 +/- 0.08

ffact.cb.enc.zv.num_scale.flt.auc.classif.kknn	0.73 +/- 0.05	0.33 +/- 0.01	0.7 +/- 0.03	0.63 +/- 0.03
ffact.cb.enc.zv.num_scale.flt.njmim.classif.net	0.65 +/- 0.09	0.32 +/- 0.13	0.81 +/- 0.07	0.5 +/- 0.04
timika.num_scale.classif.rpart	0.69 +/- 0.07	0.32 +/- 0.1	0.81 +/- 0.12	0.53 +/- 0.08
ffact.cb.enc.zv.num_scale.flt.mrmr.classif.ranger	0.68 +/- 0.14	0.31 +/- 0.16	0.75 +/- 0.11	0.65 +/- 0.11
ffact.cb.enc.zv.num_scale.flt.njmim.classif.kknn	0.72 +/- 0.07	0.31 +/- 0.06	0.84 +/- 0.03	0.45 +/- 0.2
timika.num_scale.classif.xgboost	0.72 +/- 0.03	0.31 +/- 0.03	0.83 +/- 0.07	0.47 +/- 0.04
ffact.cb.enc.zv.num_scale.flt.mrmr.classif.svm	0.71 +/- 0.07	0.3 +/- 0.3	0.72 +/- 0.08	0.7 +/- 0.05
ffact.cb.enc.zv.num_scale.pca.classif.nnet	0.65 +/- 0.04	0.3 +/- 0.17	0.57 +/- 0.1	0.68 +/- 0.1
timika.num_scale.classif.kknn	0.72 +/- 0.1	0.3 +/- 0.1	0.92 +/- 0	0.37 +/- 0.08
ffact.cb.enc.zv.num_scale.flt.auc.classif.xgboost	0.7 +/- 0.12	0.3 +/- 0.08	0.72 +/- 0.08	0.7 +/- 0.07
timika.num_scale.cb.classif.rpart	0.71 +/- 0.04	0.29 +/- 0.19	0.67 +/- 0.04	0.65 +/- 0.11
timika.num_scale.classif.log_reg	0.77 +/- 0.06	0.29 +/- 0.06	0.86 +/- 0.04	0.35 +/- 0.05
timika.num_scale.classif.multinom	0.77 +/- 0.06	0.29 +/- 0.06	0.86 +/- 0.04	0.35 +/- 0.05
ffact.cb.enc.zv.num_scale.flt.mrmr.classif.rpart	0.71 +/- 0.05	0.28 +/- 0.06	0.58 +/- 0.06	0.75 +/- 0.06
ffact.cb.enc.zv.num_scale.flt.njmim.classif.naive_bayes	0.72 +/- 0.06	0.27 +/- 0.07	0.83 +/- 0	0.45 +/- 0.15
ffact.cb.enc.zv.num_scale.pca.classif.log_reg	0.74 +/- 0.06	0.26 +/- 0.1	0.69 +/- 0.08	0.5 +/- 0.04
ffact.cb.enc.zv.num_scale.flt.mrmr.classif.xgboost	0.7 +/- 0.07	0.25 +/- 0.11	0.69 +/- 0.04	0.58 +/- 0.23
timika.num_scale.classif.naive_bayes	0.77 +/- 0.06	0.25 +/- 0.04	0.86 +/- 0.04	0.35 +/- 0.05
ffact.cb.enc.zv.num_scale.pca.classif.cv_glmnet	0.68 +/- 0.13	0.24 +/- 0.2	0.69 +/- 0.12	0.5 +/- 0.12
ffact.cb.enc.zv.num_scale.flt.njmim.classif.naive_bayes	0.73 +/- 0.06	0.24 +/- 0.04	0.83 +/- 0.04	0.45 +/- 0.12
ffact.cb.enc.zv.num_scale.pca.classif.kknn	0.65 +/- 0.02	0.24 +/- 0.02	0.65 +/- 0.03	0.65 +/- 0.07
ffact.cb.enc.zv.num_scale.pca.classif.xgb	0.64 +/-	0.22 +/-	0.64 +/- 0.07	0.53 +/- 0.19

oost	0.07	0.17		
ffact.cb.enc.zv.num_scale.pca.classif.rpart	0.6 +/- 0.07	0.22 +/- 0.13	0.61 +/- 0.08	0.53 +/- 0.04
ffact.cb.enc.zv.num_scale.pca.classif.svm	0.69 +/- 0.13	0.22 +/- 0.12	0.67 +/- 0.08	0.74 +/- 0.09
ffact.cb.enc.zv.num_scale.pca.classif.multinom	0.73 +/- 0.1	0.21 +/- 0.12	0.69 +/- 0.04	0.47 +/- 0.04
ffact.cb.enc.zv.num_scale.flt.njmim.classif.xgboost	0.68 +/- 0.06	0.21 +/- 0.08	0.69 +/- 0.12	0.6 +/- 0.11
ffact.cb.enc.zv.num_scale.flt.njmim.classif.xgboost	0.74 +/- 0.04	0.21 +/- 0.06	0.83 +/- 0.04	0.5 +/- 0.2
ffact.cb.enc.zv.num_scale.pca.classif.glmnet	0.74 +/- 0.07	0.2 +/- 0.02	0.72 +/- 0.04	0.47 +/- 0.04
ffact.cb.enc.zv.num_scale.flt.njmim.classif.kknn	0.65 +/- 0.1	0.19 +/- 0.07	0.67 +/- 0.08	0.55 +/- 0.04
ffact.cb.enc.zv.num_scale.pca.classif.naive_bayes	0.7 +/- 0.01	0.17 +/- 0.14	0.83 +/- 0.04	0.3 +/- 0.1
ffact.cb.enc.zv.num_scale.flt.njmim.classif.cv_glmnet	0.77 +/- 0.06	0.17 +/- 0.14	0.97 +/- 0.04	0.11 +/- 0.08
ffact.cb.enc.zv.num_scale.flt.mrmr.classif.kknn	0.63 +/- 0.12	0.15 +/- 0.14	0.72 +/- 0.12	0.42 +/- 0.04
ffact.cb.enc.zv.num_scale.flt.mrmr.classif.naive_bayes	0.72 +/- 0.01	0.1 +/- 0.11	0.94 +/- 0.08	0.16 +/- 0.23
ffact.cb.enc.zv.num_scale.pca.classif.featureless	0.5 +/- 0	0.04 +/- 0.14	0.5 +/- 0.04	0.47 +/- 0.08
ffact.cb.enc.zv.num_scale.flt.mrmr.classif.featureless	0.5 +/- 0	0.01 +/- 0.2	0.5 +/- 0.04	0.55 +/- 0.2
ffact.cb.enc.zv.num_scale.flt.njmim.classif.featureless	0.5 +/- 0	0 +/- 0	1 +/- 0	0 +/- 0
timika.num_scale.classif.featureless	0.5 +/- 0	0 +/- 0	1 +/- 0	0 +/- 0
ffact.cb.enc.zv.num_scale.flt.njmim.classif.featureless	0.5 +/- 0	-0.1 +/- 0.1	0.53 +/- 0.08	0.47 +/- 0.11
timika.num_scale.cb.classif.featureless	0.5 +/- 0	-0.03 +/- 0.05	0.5 +/- 0.04	0.47 +/- 0.04
ffact.cb.enc.zv.num_scale.flt.auc.classif.featureless	0.5 +/- 0	-0.02 +/- 0.23	0.47 +/- 0.12	0.5 +/- 0.15

513

514 The machine learning pipeline is shown in the pipeline column. Model performance metrics
515 include Area under the curve (classif.auc), Balanced accuracy (classif.bacc), Matthew's
516 Correlation Coefficient (classif.mcc), Sensitivity (classif.sensitivity), and Specificity
517 (classif.specificity). Each cell represents the median metric +/- the MAD for the 5-fold cross
518 validation testing on the training data representing 70% of the entire dataset. Pipelines using only

519 the Timika Score for prediction start with Timika in the pipeline name. Each pipeline shows all
 520 steps used in the pipeline ending with the machine learning algorithm used and are ordered by
 521 `classif.mcc`.
 522

523 **Table 5. Comparison of machine learning pipeline performance for predicting high**
 524 **bacterial load positive (2+ or higher) versus negative sputum microscopy status on**
 525 **validation data.**
 526

pipeline	classif.a uc	classif.m cc	classif.sensiti vity	classif.specifi city
ffact.cb.enc.zv.num_scale.flt.njmim.classif. log_reg	0.72511 63	0.455404 19	0.7466667	0.72093023
ffact.cb.enc.zv.num_scale.flt.mrmr.classif.l og_reg	0.72387 6	0.441316 61	0.7333333	0.72093023
ffact.cb.enc.zv.num_scale.flt.mrmr.classif. multinom	0.72527 13	0.441316 61	0.7333333	0.72093023
ffact.cb.enc.zv.num_scale.flt.njmim.classif. multinom	0.72434 11	0.433969 04	0.7466667	0.69767442
timika.num_scale.cb.classif.rpart	0.74108 53	0.427453 16	0.6533333	0.79069767
timika.num_scale.cb.classif.log_reg	0.71054 26	0.427426 16	0.72	0.72093023
timika.num_scale.cb.classif.multinom	0.71054 26	0.427426 16	0.72	0.72093023
timika.num_scale.cb.classif.nnet	0.71100 78	0.427426 16	0.72	0.72093023
timika.num_scale.cb.classif.svm	0.69813 95	0.427426 16	0.72	0.72093023
ffact.cb.enc.zv.num_scale.flt.mrmr.classif.c v_glmnet	0.71054 26	0.427426 16	0.72	0.72093023
ffact.cb.enc.zv.num_scale.flt.njmim.classif. cv_glmnet	0.71054 26	0.427426 16	0.72	0.72093023
ffact.cb.enc.zv.num_scale.flt.mrmr.classif.s vm	0.71271 32	0.427426 16	0.72	0.72093023
ffact.cb.enc.zv.num_scale.flt.auc.classif.rpa rt	0.74620 16	0.427066 17	0.76	0.6744186
ffact.cb.enc.zv.num_scale.flt.njmim.classif. svm	0.71875 97	0.419769 12	0.7333333	0.69767442
ffact.cb.enc.zv.num_scale.flt.auc.classif.log _reg	0.75922 48	0.412539 25	0.7466667	0.6744186
ffact.cb.enc.zv.num_scale.flt.auc.classif.mu ltinom	0.75426 36	0.412539 25	0.7466667	0.6744186

ffact.cb.enc.zv.num_scale.flt.auc.classif.nn et	0.73689 92	0.412539 25	0.7466667	0.6744186
timika.num_scale.cb.classif.naive_bayes	0.71054 26	0.405770 29	0.72	0.69767442
ffact.cb.enc.zv.num_scale.flt.auc.classif.cv_ glmnet	0.72914 73	0.405770 29	0.72	0.69767442
ffact.cb.enc.zv.num_scale.flt.mrmr.classif.r anger	0.72589 15	0.405770 29	0.72	0.69767442
ffact.cb.enc.zv.num_scale.flt.auc.classif.gl mnet	0.73968 99	0.398233 97	0.7333333	0.6744186
ffact.cb.enc.zv.num_scale.flt.njmim.classif. glmnet	0.71519 38	0.398233 97	0.7333333	0.6744186
ffact.cb.enc.zv.num_scale.flt.auc.classif.sv m	0.69472 87	0.398233 97	0.7333333	0.6744186
ffact.cb.enc.zv.num_scale.flt.mrmr.classif.n net	0.67395 35	0.391957 6	0.7066667	0.69767442
ffact.cb.enc.zv.num_scale.flt.mrmr.classif.k knn	0.73612 4	0.384348 12	0.76	0.62790698
ffact.cb.enc.zv.num_scale.flt.mrmr.classif. glmnet	0.72294 57	0.378316 79	0.6933333	0.69767442
ffact.cb.enc.zv.num_scale.flt.auc.classif.ran ger	0.71922 48	0.378316 79	0.6933333	0.69767442
ffact.cb.enc.zv.num_scale.flt.njmim.classif. ranger	0.70480 62	0.378316 79	0.6933333	0.69767442
ffact.cb.enc.zv.num_scale.flt.auc.classif.nai ve_bayes	0.74496 12	0.372004 82	0.7866667	0.58139535
ffact.cb.enc.zv.num_scale.flt.mrmr.classif.r part	0.73317 83	0.366981 06	0.6133333	0.76744186
ffact.cb.enc.zv.num_scale.flt.njmim.classif. naive_bayes	0.73891 47	0.361219 78	0.8133333	0.53488372
timika.num_scale.cb.classif.kknn	0.74899 22	0.356440 32	0.8266667	0.51162791
ffact.cb.enc.zv.num_scale.flt.njmim.classif. kknn	0.72	0.348094 66	0.7466667	0.60465116
ffact.cb.enc.zv.num_scale.flt.auc.classif.xg boost	0.74341 09	0.340846 12	0.72	0.62790698
ffact.cb.enc.zv.num_scale.flt.mrmr.classif.x gboost	0.67798 45	0.338291 99	0.6533333	0.69767442
timika.num_scale.cb.classif.ranger	0.69968 99	0.334571 34	0.6266667	0.72093023
ffact.cb.enc.zv.num_scale.flt.njmim.classif.nai ve_bayes	0.73364 34	0.317888 57	0.8133333	0.48837209
ffact.cb.enc.zv.num_scale.flt.njmim.classif.	0.72263	0.316467	0.5333333	0.79069767

rpart	57	3		
ffact.cb.enc.zv.num_scale.flt.njmim.classif.nnet	0.644186	0.312863	0.6933333	0.62790698
timika.num_scale.classif.nnet	6	89		
ffact.cb.enc.zv.num_scale.flt.njmim.classif.xgboost	0.7105426	0.3118412	0.7333333	0.58139535
ffact.cb.enc.zv.num_scale.pca.classif.glmnet	0.7010853	0.3090816	0.6	0.72093023
ffact.cb.enc.zv.num_scale.pca.classif.range	0.7029457	0.30499956	0.7066667	0.60465116
ffact.cb.enc.zv.num_scale.pca.classif.range	0.7091473	0.30302833	0.64	0.6744186
ffact.enc.zv.num_scale.flt.njmim.classif.log_reg	0.7203101	0.30141157	0.8	0.48837209
ffact.enc.zv.num_scale.flt.njmim.classif.multinom	0.7203101	0.30141157	0.8	0.48837209
ffact.enc.zv.num_scale.flt.njmim.classif.range	0.7215504	0.30141157	0.8	0.48837209
timika.num_scale.classif.svm	0.6826357	0.29154447	0.7733333	0.51162791
ffact.cb.enc.zv.num_scale.pca.classif.cv_glmnet	0.6948837	0.29008201	0.7333333	0.55813953
ffact.cb.enc.zv.num_scale.pca.classif.kknn	0.7207752	0.28582435	0.84	0.41860465
ffact.cb.enc.zv.num_scale.pca.classif.rpart	0.7221705	0.27722881	0.68	0.60465116
timika.num_scale.classif.log_reg	0.7105426	0.27363304	0.8133333	0.44186047
timika.num_scale.classif.multinom	0.7105426	0.27363304	0.8133333	0.44186047
timika.num_scale.classif.naive_bayes	0.7105426	0.27363304	0.8133333	0.44186047
ffact.cb.enc.zv.num_scale.pca.classif.svm	0.7128682	0.27363304	0.8133333	0.44186047
ffact.cb.enc.zv.num_scale.pca.classif.svm	0.7190698	0.25879866	0.64	0.62790698
ffact.cb.enc.zv.num_scale.flt.njmim.classif.glmnet	0.7210853	0.25701037	0.8	0.44186047
ffact.cb.enc.zv.num_scale.flt.auc.classif.kknn	0.6857364	0.25366746	0.72	0.53488372
ffact.cb.enc.zv.num_scale.pca.classif.multinom	0.7032558	0.24622423	0.7333333	0.51162791
ffact.cb.enc.zv.num_scale.pca.classif.kknn	0.6351938	0.24178566	0.6	0.65116279

timika.num_scale.classif.ranger	0.69426 36	0.240785 07	0.7866667	0.44186047
ffact.enc.zv.num_scale.flt.njmim.classif.xgboost	0.70325 58	0.240785 07	0.7866667	0.44186047
ffact.enc.zv.num_scale.flt.njmim.classif.rpart	0.67147 29	0.234748 13	0.8533333	0.34883721
timika.num_scale.cb.classif.xgboost	0.68093 02	0.232501 03	0.6133333	0.62790698
ffact.cb.enc.zv.num_scale.pca.classif.nnet	0.69023 26	0.231627 91	0.72	0.51162791
ffact.enc.zv.num_scale.flt.njmim.classif.cv_glmnet	0.70976 74	0.217485 43	0.9066667	0.25581395
ffact.enc.zv.num_scale.flt.njmim.classif.nnet	0.66728 68	0.198541 5	0.8266667	0.34883721
ffact.cb.enc.zv.num_scale.pca.classif.logreg	0.70387 6	0.197459 88	0.6666667	0.53488372
ffact.cb.enc.zv.num_scale.pca.classif.naive_bayes	0.67100 78	0.186109 83	0.8533333	0.30232558
timika.num_scale.classif.xgboost	0.70341 09	0.181277 05	0.8133333	0.34883721
timika.num_scale.classif.rpart	0.68682 17	0.171862 41	0.7866667	0.37209302
timika.num_scale.classif.kknn	0.70744 19	0.127279 1	0.8666667	0.23255814
ffact.cb.enc.zv.num_scale.flt.auc.classif.featureless	0.5	0.097651 87	0.52	0.58139535
ffact.cb.enc.zv.num_scale.pca.classif.featureless	0.5	0.050307 08	0.5866667	0.46511628
ffact.cb.enc.zv.num_scale.pca.classif.xgboost	0.58046 51	0.049608 19	0.4933333	0.55813953
ffact.cb.enc.zv.num_scale.flt.mrmr.classif.naive_bayes	0.71674 42	0.036994 84	0.9866667	0.02325581
ffact.cb.enc.zv.num_scale.flt.mrmr.classif.featureless	0.5	0.027175 29	0.4933333	0.53488372
ffact.cb.enc.zv.num_scale.flt.njmim.classif.featureless	0.5	0.011382 35	0.5466667	0.46511628
timika.num_scale.classif.featureless	0.5	0	1	0
ffact.enc.zv.num_scale.flt.njmim.classif.featureless	0.5	0	1	0
timika.num_scale.cb.classif.featureless	0.5	- 0.097651 87	0.48	0.41860465

528 The machine learning pipeline is shown in the pipeline column. Model performance metrics
529 include Area under the curve (classif.auc), Balanced accuracy (classif.bacc), Matthew's
530 Correlation Coefficient (classif.mcc), Sensitivity (classif.sensitivity), and Specificity
531 (classif.specificity). Each cell represents the metric for the performance on the 30% held-out
532 validation test set after training on the 70% training data. Pipelines using only the Timika Score
533 for prediction start with Timika in the pipeline name. Each pipeline shows all steps used in the
534 pipeline ending with the machine learning algorithm used and are ordered by classif.mcc.
535

536 **Comparison of best performing feature selection and Timika-only models by** 537 **bootstrapping**

538 Though our initial training and validation testing suggested that Timika Score only
539 pipelines showed minimal differences with the top 5 feature models, we wanted to further test
540 this for statistical significance. We chose the best performing class-balanced pipelines from the
541 5-fold cross validated results obtained from the training data. Thus, the best top 5 feature
542 pipeline, Timika Score only pipeline, and a featureless pipeline were selected for further testing.
543 The featureless workflow is a control that would reveal the density of results expected due to
544 random chance regardless of any upstream preprocessing given that classes are balanced. We
545 perform a bootstrapping without replacement ($N = 200$) on the entire dataset using 70% of the
546 data for training and 30% for testing on each split.

547 The bootstrapping result on the prediction problem attempting to distinguish positive
548 from negative sputum results showed that the performance of the best top 5 feature pipeline was
549 slightly better than the Timika-only pipeline. Both showed significantly improved performance
550 from the workflow using a featureless model (Fig 4A). The bootstrapping on the prediction
551 problem attempting to distinguish higher bacterial load sputum results (2+ or higher compared to
552 negative) did not show any difference in performance for the best Timika-only pipeline
553 compared to the best top 5 workflow (Fig 4B). As before, both pipelines performed significantly

554 better than the workflow using the featureless model. These results are consistent with the idea
555 that Timika offers generally equivalent predictive performance to using the top 5 features from
556 the dataset. Inclusion of other features may offer minimal improvements in predictive
557 performance depending upon the model or set of features selected.

558 **Fig 4. Comparison of best class-balanced pipelines via bootstrapping without replacement**
559 **(N = 200).** The best class-balanced pipelines via 5-fold cross-validation performance on training
560 data were selected for comparison to assess using top 5 features via feature selection or
561 dimensionality reduction as compared to using only Timika Score. A featureless pipeline is used
562 as a control to show expected performance via random selection of outcome. Box plots with
563 interquartile range are overlaid on density plots showing the density of results of Matthew's
564 Correlation Coefficient across all bootstrapping results per pipeline. Performance is compared
565 across all groups by Kruskal-Wallis and by individual pairs using Wilcoxon test which are
566 shown by the brackets. In A), the performance of the best pipelines for predicting positive (1 to 9
567 in 100, 1+, or higher) versus negative is shown whereas B) the performance of the best pipelines
568 for predicting high bacterial load positive (2+ or higher) versus negative is shown. Interestingly
569 there is a statistically significant difference between top 5 feature pipeline versus Timika Score
570 pipeline in A) and no significant difference in performance for Timika Score pipeline in B). This
571 shows that additional derived features from radiologist observations may result in small
572 performance gains although Timika Score alone provides generally equivalent prediction
573 performance for the identified best models.
574

575 Discussion

576 X ray imaging is a useful approach to diagnose and monitor disease progression and
577 status during routine TB clinical management. X ray imaging cannot discern the type of
578 resistance of tuberculosis as well as characterize the amount of pathogen in sputum, which only
579 microbiological methods can provide. These approaches assist clinicians with a more complete
580 understanding of the case and understanding their relationship is important. CXR is relatively
581 less expensive than other imaging modalities such as CT permitting its wider use especially in
582 LMIC that may face challenges with infrastructure cost to support routine CT use. Using CXRs,
583 radiologists can report on a variety of observations that determine lung biomarker status such as
584 overall abnormal volume of the lungs, presence of cavity, and presence of nodule, which the TB

585 Portals resource collects, standardizes, and provides as part of the patient record. Here we
586 investigated the previously reported Timika Score that can be derived from CXR radiologist
587 observations to characterize pre-treatment severity of disease. TB portals provides a unique real-
588 world repository of TB cases, especially drug resistant cases to bridge across distinct domains
589 including radiological, pathogen genomic, microbiological, and clinical features; it is especially
590 suited to serve as a large reference resource for assessing derived scores like Timika Score for
591 testing in a real-world database of especially challenging TB cases. Our goal was to assess the
592 plausibility and utility of the derived Timika Score within this real-world resource by studying its
593 relationships to the other available case characteristics.

594 We demonstrate that Timika Score associates with other case characteristics consistent
595 with prior reporting of TB clinical risk factors. For instance, we show that images from patients
596 with a lower BMI tended to have a higher Timika Score and less of a change from the initial
597 CXR to the last available CXR. TB and BMI have been reported to show a strongly logarithmic
598 association and there was reported fivefold increase in age-adjusted incidence of new pulmonary
599 TB in lowest BMI group compared to highest in a study of 1.7 million Norwegians (18, 19). In
600 the same report, Tverdal mention an interesting U-shaped association with BMI and all-cause
601 mortality which is strikingly like the U-shaped association we observed with Timika Score and
602 BMI. In our analysis, increasing age tended to associate with higher Timika Score and less of a
603 change from the initial CXR to the last available CXR. This is consistent with higher mortality,
604 morbidity, and risk of TB with increasing age especially since the symptoms of TB may be
605 confused with other age-related illnesses (20) resulting in delayed diagnosis or treatment.
606 Moreover, when comparing Timika Score with other clinical factors associated with the case, we
607 observe both higher Timika Scores and lower relative changes in Timika Score in cases with

608 higher-risk clinical factors (XDR, Relapse, etc.) or poor reported outcomes (e.g. Treatment
609 failure, Died, etc.). This is consistent with prior reports examining predictors that affect change
610 in radiological lesions over the course of treatment monitoring (13).

611 Finally, we demonstrated that Timika Score show a clear and statistically significant
612 predictive capacity for baseline, pre-treatment sputum microscopy status in the cohort of new
613 cases we identified from the TB Portals repository. This is important because the original reports
614 on Timika Score suggested the same association on a smaller dataset (15) that did not span
615 across the wide-range of participating sites from 14 countries. Taken together, these
616 observations support that the Timika Scores we are calculating reflect lung biomarker status
617 consistent with accumulated knowledge of the radiological and clinical associations in TB
618 disease.

619 This analysis has several limitations and caveats when interpreting results. The TB
620 Portals is a real-world data repository to better understand DR TB so it is challenging to separate
621 identified associations with other observed or unobserved variables from the case. Moreover,
622 respective images and test results for each case are not collected uniformly in time but rather as
623 clinical management of the case allows. We select cases for inclusion into the analysis cohort
624 based upon criteria we believe will accurately represent associations between images and
625 microbiological test results but we cannot rule out timing or other aspects of the case impacting
626 the associations we observe. For example, we noted both lungs involvement of calcified nodule
627 as showing higher risk of pre-treatment sputum positivity. Such a marker suggests a long-term
628 prior history of pulmonary TB that might not be reflected in the “New” case definition from
629 WHO. Collecting a prior history of chronic lung symptoms around the baseline sample
630 collection might allow us to see if the relationships we identified remain after stratifying by

631 these symptoms that could suggest a period of prior disease burden. Given these caveats, the
632 modeling and visualizations need to be interpreted as hypothesis-generating.

633 A key goal of this study was to identify a risk score (either new or previously reported)
634 that could encapsulate a temporal snapshot of case dynamics relating to disease risk such as poor
635 outcome or infectiousness. The CXR derived Timika Score may provide a useful score in this
636 regard from the initial testing we performed. One caution with regards to the utility of Timika
637 Score from this analysis is that the risks and associations with sputum microscopy positivity
638 were calculated for samples taken prior to treatment. The dynamics of microbiological status in
639 sputum might not reflect the same dynamics of lung biomarkers in response to treatment. Given
640 these dynamics, applying the same risks from Timika Score for sputum positivity (i.e. presence
641 of pathogen in sputum) after treatment is not advisable and may require new approaches that
642 stratify by type of resistance and different treatment regimens. It also may require additional data
643 collection to support the requisite number of cases given the real-world nature of the resource
644 where only subsets of cases may meet the inclusion and exclusion criteria for analysis.

645 The temporal relationships between lung status (as observed in CXR images) and
646 bacterial pathogen load in the sputum (as observed by microscopy) are complex. For instance, at
647 the beginning of disease, radiological features may not be detectable in the lungs despite the
648 presence of bacteria in sputum. Meanwhile, towards the end of treatment, sputum may no longer
649 contain TB pathogen indicating a non-infectious case; however, the pathogen may remain in
650 certain areas of the lung such as cavities. Given these intricacies, additional research is
651 warranted to determine if improvements can be made in Timika Score to account for these
652 situations. Moreover, there may be limitations to the granularity of a clinical score like Timika
653 that can be generated from CXRs given the limitations of the modality with regards to imaging

654 detail. It may be necessary for other modalities such as CT to account for more complex features
655 such as the size of cavities, nodules or other aspects identified in the lung, which may not be
656 obvious on a CXR. The assessment of clinical scores such as Timika Score coming from these
657 various modalities is especially important for hard-to-treat cases such as MDR and XDR TB
658 where scores can be compared in the context of other features of the case.

659 We observed the best predictive performance in models predicting higher bacterial load
660 sputum status. This improvement in predictive performance for high bacterial load sputum
661 statuses (e.g. 2+ or higher) illustrates the nuances associated with predicting sputum status from
662 Timika Score. The borderline cases such as 1 to 9 in 100 or 1+ are more challenging to predict
663 as pathogen load is only slightly higher than the negative status specimens and furthermore some
664 negative samples may be false negatives. These false negatives may suffer from issues such as
665 sensitivity due to the challenges of acquiring a usable sputum sample. Despite these nuances,
666 our results confirm prior reported Timika Score utility for predicting baseline sputum positivity
667 albeit with better performance for high-bacterial load sputum samples. The high pathogen load
668 cases would also be among the most infectious and challenging to treat; therefore, we were
669 satisfied to observe the higher predictive performance in this clinically important group. We
670 believe that adding CXR-derived Timika Score to the TB Portals resource will open
671 opportunities to other researchers to utilize this score to understand TB in a real-world setting.
672 The score can serve as a reference from which to test against additional clinical scores that could
673 be derived from the available set of features captured in TB portals.

674 **Acknowledgements**

675 We would like to thank Qinlu Wang, Jingwen Gu, and Ziv Yaniv for helpful suggestions; the
676 MLR3 team for development of the MLR3 suite of packages. This research was supported in

677 part by the Office of Science Management and Operations (OSMO) of the NIAID. For their
678 contributions to the vision and requirements of TB Portals, we would like to thank: Mike
679 Tartakovsky, Darrell Hurt, Alina Grinev, and members of the TB Portals team.

680 References

681

- 682 1. Dheda K, Barry CE, & Maartens G (2016) Tuberculosis. *Lancet* 387(10024):1211-1226.
- 683 2. Chakaya J, *et al.* (2021) Global Tuberculosis Report 2020 - Reflections on the Global TB
684 burden, treatment and prevention efforts. *Int J Infect Dis*.
- 685 3. Magro P, *et al.* (2020) Impact of the SARS-CoV-2 epidemic on tuberculosis treatment
686 outcome in Northern Italy. *Eur Respir J* 56(4).
- 687 4. Migliori GB, *et al.* (2020) Worldwide Effects of Coronavirus Disease Pandemic on
688 Tuberculosis Services, January-April 2020. *Emerg Infect Dis* 26(11):2709-2712.
- 689 5. Dheda K, *et al.* (2017) The epidemiology, pathogenesis, transmission, diagnosis, and
690 management of multidrug-resistant, extensively drug-resistant, and incurable
691 tuberculosis. *Lancet Respiratory Medicine* 5(4):291-360.
- 692 6. (WHO) WHO (2018) Global Tuberculosis report 2018.
- 693 7. Manjelienskaia J, Erck D, Piracha S, & Schragger L (2016) Drug-resistant TB: deadly,
694 costly and in need of a vaccine. *Trans R Soc Trop Med Hyg* 110(3):186-191.
- 695 8. Desikan P (2013) Sputum smear microscopy in tuberculosis: is it still relevant? *Indian J*
696 *Med Res* 137(3):442-444.
- 697 9. Olaru ID, Heyckendorf J, Grossmann S, & Lange C (2014) Time to culture positivity and
698 sputum smear microscopy during tuberculosis therapy. *PLoS One* 9(8):e106075.
- 699 10. Su WJ, Feng JY, Chiu YC, Huang SF, & Lee YC (2011) Role of 2-month sputum smears
700 in predicting culture conversion in pulmonary tuberculosis. *Eur Respir J* 37(2):376-383.
- 701 11. Saul EE, *et al.* (2020) The challenges of implementing low-dose computed tomography
702 for lung cancer screening in low- and middle-income countries. *Nature Cancer*
703 1(12):1140-1152.
- 704 12. Miller C, Lonroth K, Sotgiu G, & Migliori GB (2017) The long and winding road of
705 chest radiography for tuberculosis detection. *Eur Respir J* 49(5).
- 706 13. Heo EY, *et al.* (2009) Radiographic improvement and its predictors in patients with
707 pulmonary tuberculosis. *Int J Infect Dis* 13(6):e371-376.
- 708 14. Chakraborty A, Shivananjaiah AJ, Ramaswamy S, & Chikkavenkatappa N (2018) Chest
709 X ray score (Timika score): an useful adjunct to predict treatment outcome in
710 tuberculosis. *Adv Respir Med* 86(5):205-210.
- 711 15. Ralph AP, *et al.* (2010) A simple, valid, numerical score for grading chest x-ray severity
712 in adult smear-positive pulmonary tuberculosis. *Thorax* 65(10):863-869.
- 713 16. Anonymous (2020) *WHO consolidated guidelines on tuberculosis: Module 4: Treatment*
714 *- Drug-resistant tuberculosis treatment*, WHO Guidelines Approved by the Guidelines
715 Review Committee, Geneva).
- 716 17. Lang M, *et al.* (2019) mlr3: A modern object-oriented machine learning framework in R.
717 *Journal of Open Source Software* 4(44).
- 718 18. Tverdal A (1986) Body mass index and incidence of tuberculosis. *Eur J Respir Dis*
719 69(5):355-362.
- 720 19. Lonroth K, Williams BG, Cegielski P, & Dye C (2010) A consistent log-linear
721 relationship between tuberculosis incidence and body mass index. *Int J Epidemiol*
722 39(1):149-155.

- 723 20. Rajagopalan S (2001) Tuberculosis and aging: a global health problem. *Clin Infect Dis*
724 33(7):1034-1039.
725

726 **Supporting Information**

727 **S1 Table. Case characteristics of the cohort of cases selected for evaluation of baseline**
728 **sputum microscopy status (N = 572).** Case characteristics were compared by baseline sputum
729 microscopy status. P-values were calculated for continuous variables (age_of_onset, bmi,
730 overall_timika) using analysis of variance test. P-values for categorical variables
731 (registration_date, gender, country, type_of_resistance, outcome, current_smoker) were
732 calculated using Chi-squared test.
733

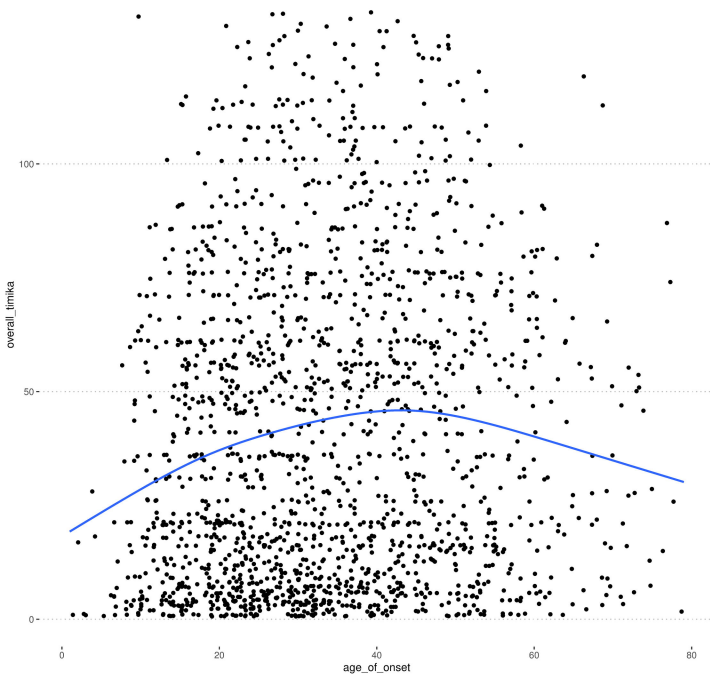
734 **S2 Table. Case characteristics of the cases with CXR radiologist observations used for**
735 **assessing Timika Score in relation to other case characteristics (N = 1761).** Cases in the TB
736 portals publicly shared dataset having a CXR with available radiologist were selected. The case
737 characteristics are shown.
738

739 **S3 Table. CXR derived features from radiologist observations in the cohort of cases**
740 **selected for evaluation of baseline sputum microscopy status (N = 572).** The derived features
741 from the available radiologist observations from the cohort of selected cases used for evaluation
742 of baseline sputum microscopy status were compared by baseline sputum status. P-values were
743 calculated for continuous variables (mean_collapse, mean_smallcavity, mean_mediumcavity,
744 mean_largecavity, mean_lowden, mean_medden, mean_highden, mean_smallnodule,
745 mean_mediumnodule, mean_largenodule, mean_hugenodule,
746 mean_lowgroundglassdensityactivefreshnodules, mean_fibroticnodule, mean_calcsequella,
747 overall_timika) using analysis of variance test. P-values for categorical variables (both_lungs,
748 both_collapse1, both_smallcavity1, both_mediumcavity1, both_largecavity1,
749 both_isanylargetcavitymult, both_multiplecavitiesbeseen, both_lowden1, both_medden1,
750 both_highden1, both_smallnodule1, both_mediumnodule1, both_largenodule1,
751 both_hugenodule1, both_calcnod, both_noncalcnod, both_clustnod, both_multnod,
752 both_lowgroundglassdensityactivefreshnodules1, both_fibroticnodule1, both_calcsequella1)
753 were calculated using Chi-squared test.
754
755

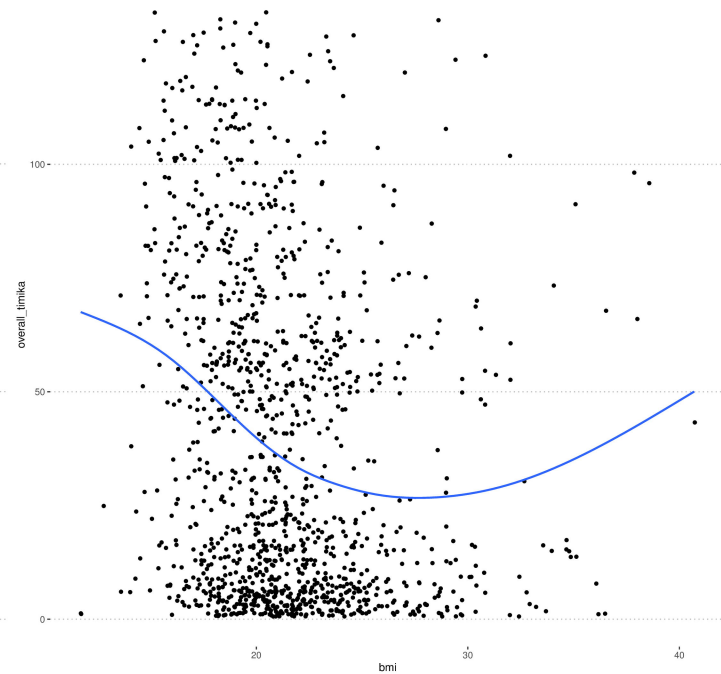
756 **S4 Table. Patient and condition ids for the cohort of cases selected for evaluation of**
757 **baseline sputum microscopy status (N = 572).** A table of patient and condition ids is provided
758 for the de-identified records that were used for evaluation of baseline sputum microscopy status.
759

760 **S5 Table. Patient, condition, and imaging ids for the cases having CXRs with radiologist**
761 **observations used for assessing Timika Score in relation to other case characteristics (N =**
762 **1761).** A table of patient, condition, and imaging ids with associated relative date of imaging is
763 provided for the de-identified records that were used for Timika visualizations. For images used
764 for temporal analysis of changes in Timika Score over time, the temporal_analysis column
765 provides a filter equal to “yes” to select only those sets of images.

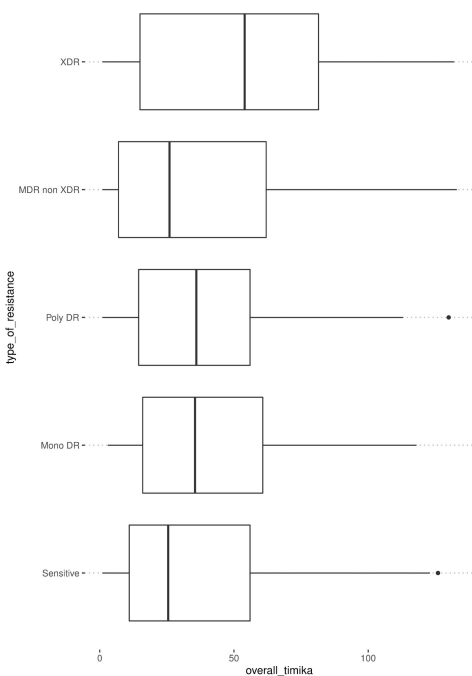
Overall timika stratified by age_of_onset



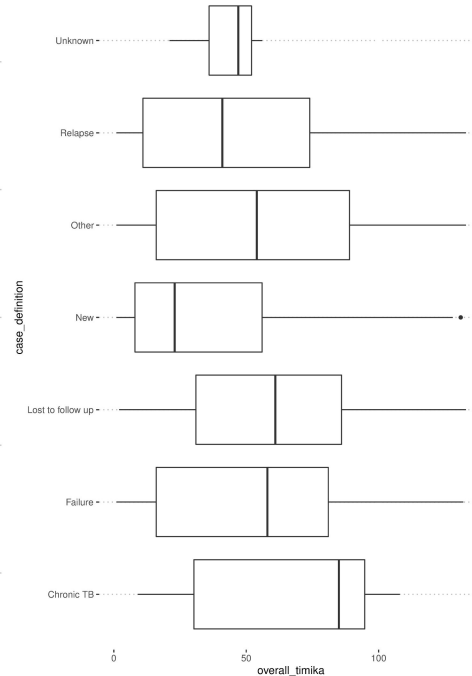
Overall timika stratified by bmi



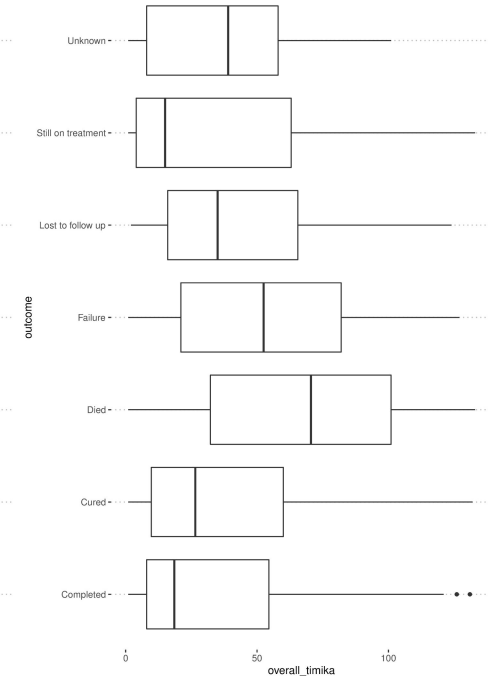
C Overall timika stratified by type_of_resistance



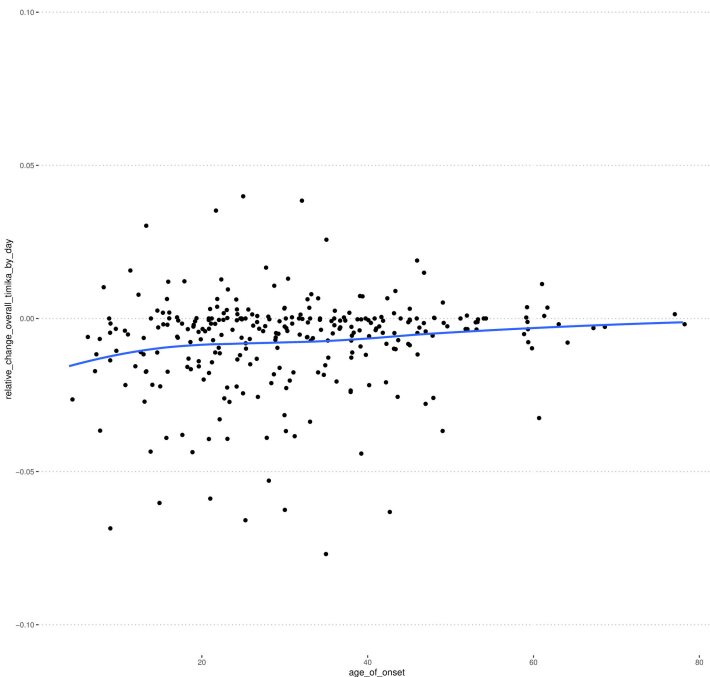
D Overall timika stratified by case_definition



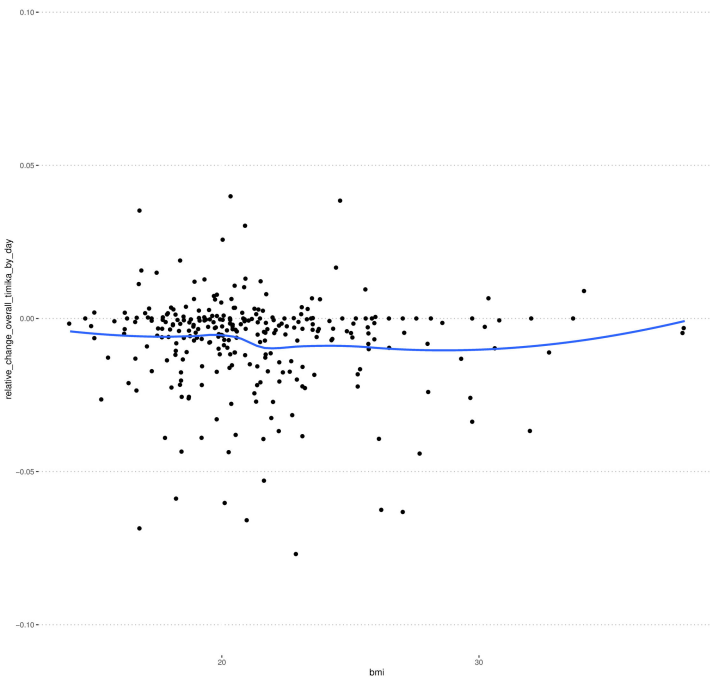
E Overall timika stratified by outcome



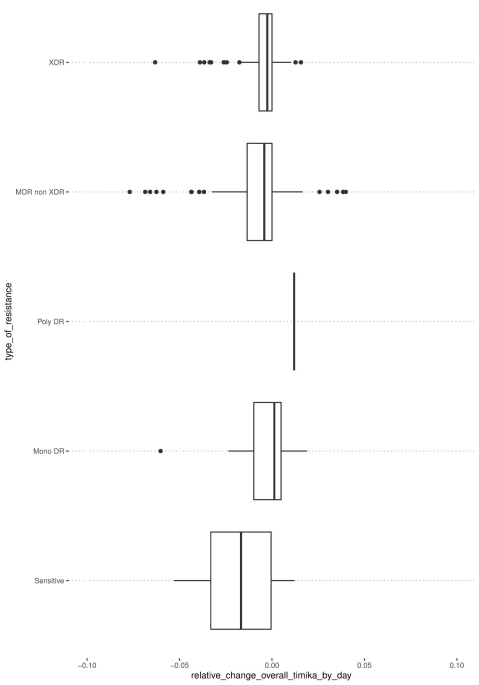
A log2 change in timika stratified by age_of_onset



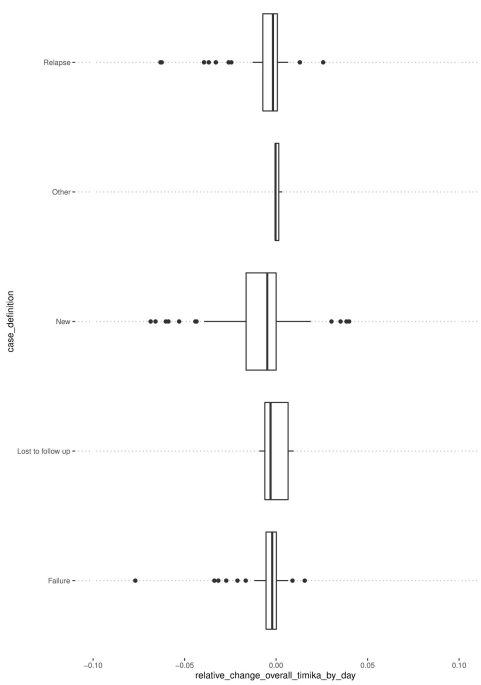
B log2 change in timika stratified by bmi



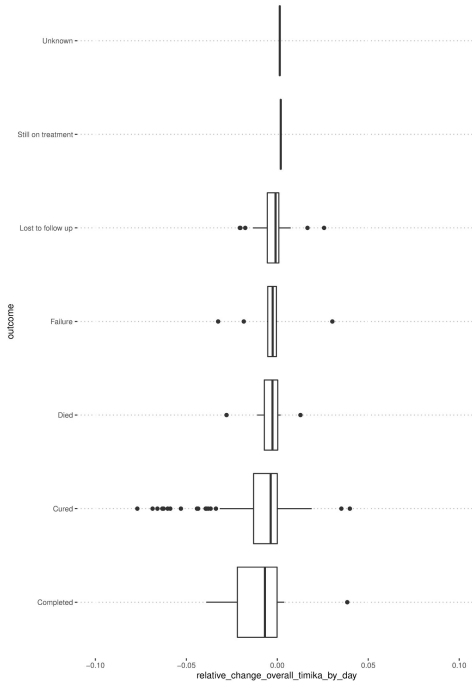
C log2 change in timika stratified by type_of_resistance



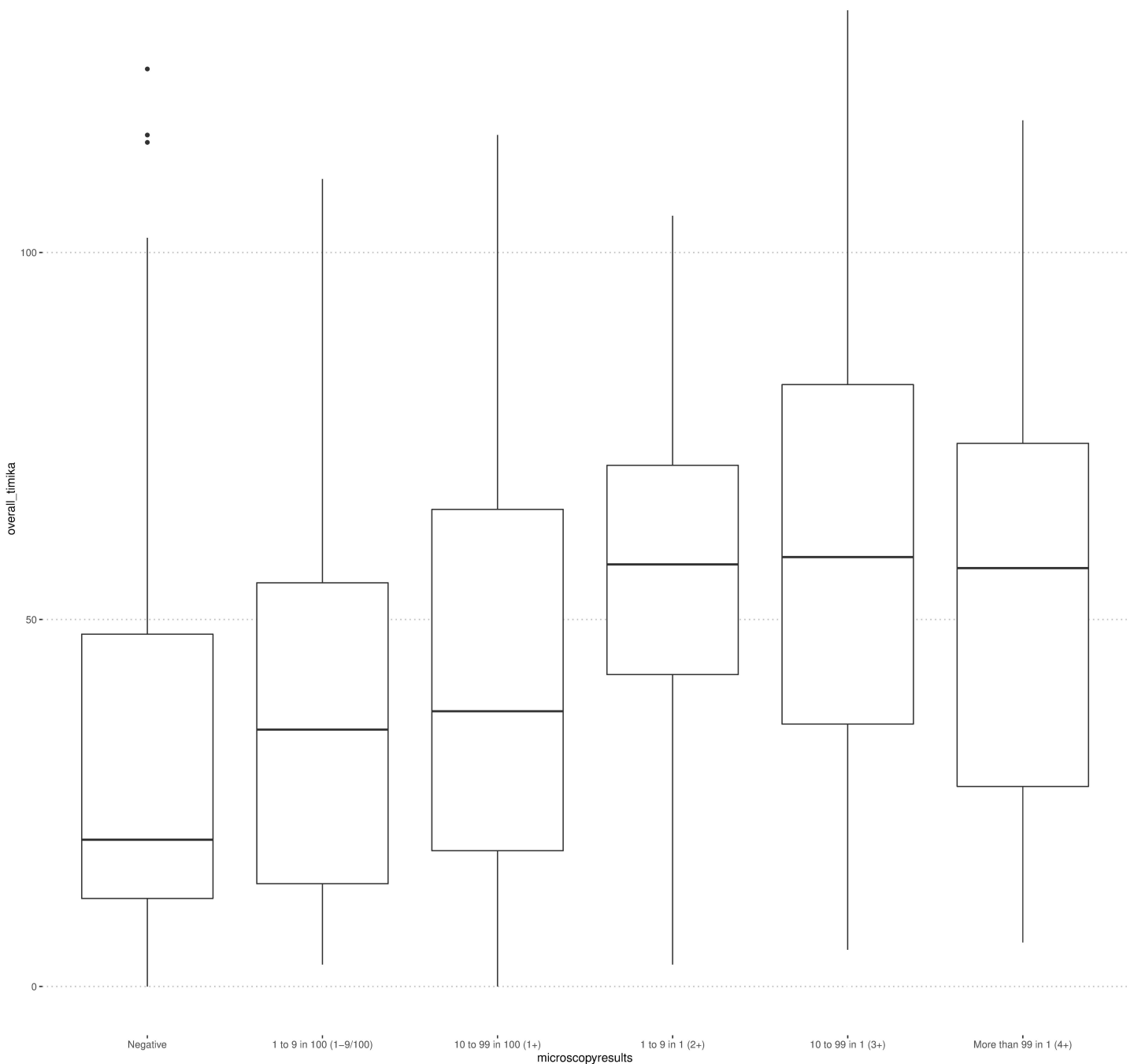
D log2 change in timika stratified by case_definition



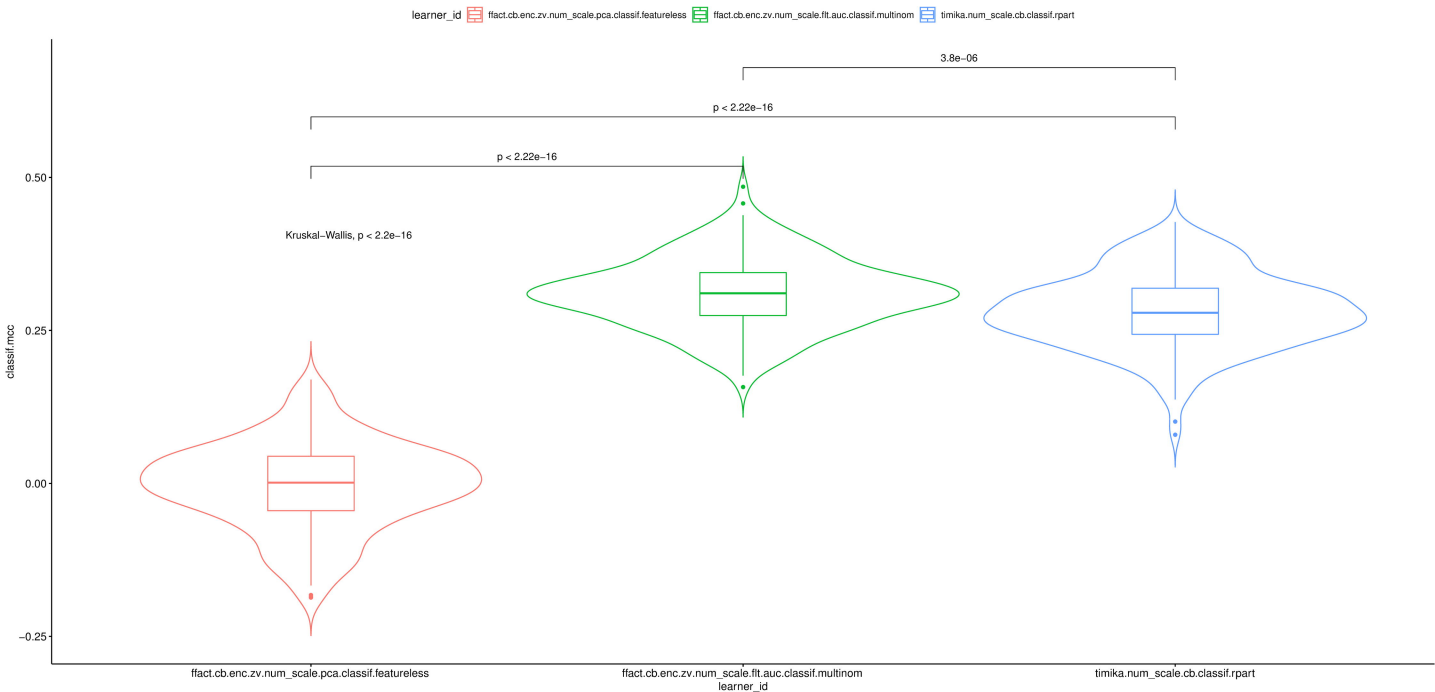
E log2 change in timika stratified by outcome



Timika score stratified by sputum smear status of specimen prior to treatment start from CXR within two weeks of specimen collection date



Comparison bootstrapping without replacement (N = 200) from best models



B Comparison bootstrapping without replacement (N = 200) from best models

