

1 **Symptom-based prediction model of SARS-CoV-2 infection developed from**
2 **self-reported symptoms of SARS-CoV-2-infected individuals in an online**
3 **survey**

4

5

6

7 Hansjörg Schulze¹, Daniel Hoffmann², Wibke Bayer^{1, *}

8

9

10 ¹ Institute for Virology, University Hospital Essen, University Duisburg-Essen, Hufelandstr.
11 55, 45147 Essen, Germany

12

13 ² Bioinformatics and Computational Biophysics, Faculty of Biology, University Duisburg-
14 Essen, Universitätsstr. 2, 45141 Essen, Germany

15

16 * corresponding author:

17 Email: wibke.bayer@uni-due.de

18 phone: +49-201-723-83034

19

20 **Abstract**

21 **Background:** Infections with the newly emerged severe acute respiratory syndrome virus 2
22 (SARS-CoV-2) have quickly reached pandemic proportions and are causing a global health
23 crisis. First recognized for the induction of severe disease, the virus also causes asymptomatic
24 infections or infections with mild symptoms that can resemble common colds. Since
25 infections with mild course are probably a major contributor to the spread of SARS-CoV-2,
26 better detection of such cases is important. To provide better understanding of these mild
27 SARS-CoV-2 infections and to improve information for potentially infected individuals, we
28 performed a detailed analysis of self-reported symptoms of SARS-CoV-2 positive and SARS-
29 CoV-2 negative individuals.

30 **Methods:** In an online-based survey, 963 individuals provided information on symptoms
31 associated with an acute respiratory infection, 336 of the participants had tested positive for
32 SARS-CoV-2 infection, 107 had tested negative, and 520 had not been tested for SARS-CoV-
33 2 infection.

34 **Results:** The symptoms reported most frequently by SARS-CoV-2 infected individuals were
35 tiredness, loss of appetite, impairment of smell or taste and dry cough. The symptoms with the
36 highest odds ratios between SARS-CoV-2 positive and negative individuals were loss of
37 appetite and impairment of smell or taste. Based on the most distinguishing symptoms, we
38 developed a Bayesian prediction model, which had a positive predictive value of 0.80 and a
39 negative predictive value of 0.72 on the SARS-CoV-2 tested individuals. The model predicted
40 56 of 520 non-tested individuals to be SARS-CoV-2 positive with more than 75% probability,
41 and another 84 to be SARS-CoV-2 positive with probability between 50% and 75%.

42 **Conclusions:** A combination of symptoms can provide a good estimate of the probability of
43 SARS-CoV-2 infection.

44

45 **Keywords**

46 SARS-CoV-2, COVID-19, symptoms, mild disease course, Bayesian model, Bayes model,

47 prediction, symptom model

48

49 **Background**

50 The outbreak of severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2), the virus
51 associated with the coronavirus infectious disease 19 (COVID-19), at the end of 2019 [1] has
52 quickly developed into pandemic proportions and was declared a public health emergency of
53 international concern by the World Health Organization on January 30th 2020 [2] . The first
54 infections in Europe occurred in Italy and Germany in late January 2020 [3]. By March 2020,
55 SARS-CoV-2 infection rose rapidly in many countries, and have started to do so again with
56 the fall season in the northern hemisphere; as of November 6th 2020, more than 48 million
57 people have been infected by SARS-CoV-2, and more than 1.2 million SARS-CoV-2 related
58 deaths have been reported [4].

59 Due to the rapid rise in demand for SARS-CoV-2 diagnostic tests and a resulting shortness of
60 supplies, not all patients could be given access to SARS-CoV-2 testing in spring 2020, and a
61 pre-screening was performed on the basis of symptoms, and prior residence in SARS-CoV-2
62 risk areas or prior contact to individuals who had tested SARS-CoV-2 positive. The lack of
63 access to testing led to many uncertainties for patients with cold symptoms. From the first
64 cases in China, the typical symptoms of SARS-CoV-2 infection were described as dry cough,
65 fever, and pneumonia [5, 6]. These first reports of symptoms associated with SARS-CoV-2
66 infection / COVID-19 were based on symptoms found in hospitalized patients with severe
67 disease. Notably, in the first reports on COVID-19, disease course was considered mild if
68 patients did not require ventilation. It soon became obvious though that many SARS-CoV-2
69 infected individuals experienced asymptomatic or now commonly considered non-severe
70 disease course that resembled common colds and did not require hospitalization [7]. Since the
71 hospitalized and critically ill patients were of the greatest concern in the early days of the
72 pandemic, the range of symptoms of SARS-CoV-2 infected individuals with mild symptoms
73 was not the focus of attention. In following studies, it has been found that many SARS-CoV-2

74 infected patients experienced altered or completely lost sense of smell and / or taste [8, 9],
75 which has since been regarded as a key indicator of SARS-CoV-2 infection.

76 To get a more detailed picture of the symptoms of SARS-CoV-2 infection, we created an
77 online questionnaire and invited both tested and untested individuals to report their symptoms.
78 The symptoms included in the survey comprised general symptoms such as fever, fatigue,
79 headache and joint or muscle pains, and more specific eye, nose, throat, respiratory and
80 gastrointestinal symptoms. The self-reported symptoms of SARS-CoV-2 positive- and
81 negative-tested individuals were then used to fit a Bayesian model, which was used to predict
82 probabilities of a SARS-CoV-2 infection of the non-tested survey participants.

83

84 **Methods**

85 **Data collection**

86 Data were collected in an online questionnaire based on LimeSurvey software hosted on the
87 servers of the University of Duisburg-Essen. Participants of the online survey were recruited
88 via public health offices of the city of Hamm (Northrhine-Westphalia, Germany), of the
89 administrative district Soest (Northrhine-Westphalia, Germany), of the administrative district
90 Hochsauerland (Northrhine-Westphalia, Germany) and the SARS-CoV-2 testing center in
91 Lünen, administrative district Unna (Northrhine-Westphalia, Germany), as well as via social
92 media. Participants were invited to complete the survey in case of a positive as well as
93 negative SARS-CoV-2 test result, and also if they had not been tested but had cold symptoms.
94 The data were collected between April 6th and September 1st 2020.

95 **Data analysis**

96 Data were analyzed using R (version 3.6.3) and RStudio software and fsmb, plyr, splyr,
97 tidyverse, randomForest, patchwork, rstanarm, ROCR, ggplot2, viridis and bayesplot
98 packages. Only data from individuals who had completed the survey were included in the
99 analysis.

100 For graphic representation of the survey results of tested individuals, data were sorted for the
101 SARS-CoV-2 test result, the severity of the SARS-CoV-2 symptoms and the severity of fever
102 using the tidyverse package, and a heatmap of the survey data was generated using ggplot2
103 package.

104 Frequencies of individual symptoms were calculated and visualized as heatmap using ggplot2
105 and viridis package. Odds ratios for individual symptoms were calculated as the conditional
106 maximum likelihood estimate with 90% confidence interval using the pairwise.fisher.test
107 command of R package fsmb, results were visualized using ggplot2.

108 Random Forest analysis of the survey data was performed using R package randomForest 4.6-
109 14 [10] with default parameters. Importance of individual symptoms was characterized as
110 mean decrease of Gini index on random permutations of case labels.

111 A Bayesian regression model was fitted with R package rstanarm version 2.21.1 [11] using
112 the function stan_glm with default priors, on all or 10 selected symptoms of the tested
113 individuals' data as indicated in the results section, using a posterior sample size of 4000. To
114 adjust for differences in the number of SARS-CoV-2 positive and negative tested participants,
115 the dataset from negative participants ($n = 107$) was increased by bootstrapping to match the
116 number of positive participants ($n = 336$). Bootstrapping was performed 4 times to exclude a
117 bias of the results, with similar outcome. The R package bayesplot [12] was used to confirm
118 the model validity by posterior predictive checks. Performance of the Bayesian and Random
119 Forest based models as ROC curves (true positive rate and false positive rate) and area under
120 the ROC curve were calculated with R package ROCR [13]. The numbers of true positives
121 (TP), false positives (FP), true negatives (TN) and false negatives (FN) was determined using
122 posterior predictions > 0.5 as positives and posterior predictions ≤ 0.5 as negative; positive
123 predictive value (PPV) was calculated as $PPV = TP / (TP + FP)$, negative predictive value
124 (NPV) was calculated as $NPV = TN / (TN + FN)$. The credibility intervals for PPV and NPV
125 were calculated using the Voila online tool [14], which is based on a Bayesian classifier
126 uncertainty estimate method [15].

127 For prediction of SARS-CoV-2 infection in untested individuals, the Bayesian model was
128 applied to the dataset from the untested individuals and the most likely SARS-CoV-2
129 infection state was calculated for each of these individuals. For visualization of the prediction
130 outcome, a heatmap was generated using ggplot2 after sorting the non-tested individuals by
131 the predicted probability of SARS-CoV-2 infection as calculated for the Bayesian model.

132

133 **Results**

134 To obtain an overview of the subjective symptoms reported by SARS-CoV-2 infected
135 individuals, an online survey was performed. Invitations to the survey were published on
136 social media, and SARS-CoV-2 tested individuals were also directly invited to participate
137 when they received notification of their test results, or by subsequent invitation, by local
138 public health offices. Participants included both tested individuals, as well as non-tested
139 individuals, who were considered separately in the analysis of the survey. The characteristics
140 of the tested study population are represented in Table 1. In total, responses of 443
141 participants who had been tested for SARS-CoV-2 infection were included in the survey, 336
142 (75.8 %) had tested positive for SARS-CoV-2 infection, while 107 (24.2 %) of those
143 participants had tested negative. 293 (66.4 %) of the tested participants were recruited via the
144 public health office, whereas 150 (33.6 %) tested participants were recruited from social
145 media.

146 Of the 336 SARS-CoV-2 positive individuals, only 18 (5.3 %) reported no symptoms,
147 whereas 141 (42.2 %) and 102 (30.1 %) reported their symptoms as mild or moderate,
148 respectively, and 49 (14.7 %) reported their symptoms as severe. 20 (5.9 %) individuals were
149 treated in the hospital due to the SARS-CoV-2 infection, and another 6 (1.8 %) required
150 intensive care treatment.

151 The symptoms reported by the tested study population are shown as a heatmap in Figure 1,
152 frequencies of the symptoms reported by SARS-CoV-2 positive and negative individuals are
153 presented in Figure 2A. Clearly, a majority in both SARS-CoV-2 positive and negative
154 individuals reported general symptoms such as tiredness and lethargy. The presence of eye
155 symptoms in both groups was low, and also most nose symptoms were reported only by a
156 minority of participants. Interestingly, about two thirds of SARS-CoV-2 positive individuals
157 reported an impaired or lost sense of smell or taste. Dry cough was reported by almost 60% of

158 both SARS-CoV-2 positive and negative individuals, whereas most SARS-CoV-2 positive
159 individuals did not report other throat symptoms including throat pain. Approximately one
160 third of individuals reported breathing problems, and also diarrhea was reported by
161 approximately a third of SARS-CoV-2 positive and negative individuals.

162 The frequency of symptoms in SARS-CoV-2 positive and negative individuals were used to
163 calculate the odds ratios for the individual symptoms (Figure 2B). Since the distribution of
164 many symptoms was similar in SARS-CoV-2 positive and negative individuals, many odds
165 ratios are close to 1. Clearly, the highest odds ratios were observed for impaired or lost sense
166 of smell (estimate: 4.7; 90% confidence interval (CI): 2.8 – 8.1) and impaired or lost sense of
167 taste (4.3; 90% CI: 2.6 – 7.2), whereas the other nose symptoms as well as eye and throat
168 symptoms were found to have odds ratios below 1. The lowest odds ratios were observed for
169 throat pain and burning sensation in the throat. The odds ratio of dry coughing is also slightly
170 below 1 (0.9; 90% CI: 0.6 – 1.5), which may likely be an underestimation due to a bias in
171 eligibility criteria at the time of testing.

172 For the development of a prediction model, we performed a Bayesian regression analysis of
173 all symptoms and plotted them with their odds ratio in combination with the overall frequency
174 (Figure 3A). As we have an imbalanced data set with less negative than positive tests, the
175 resulting model would have been biased towards positive predictions. To avoid this, we fitted
176 the model with two different approaches. In one, we increased the number of negatives by
177 bootstrapping to match the number of positives, in the other we decreased the number of
178 positives to match the number of negatives. Both approaches gave comparable results, but we
179 favored the first one as it uses all the available data and is less prone to the introduction of a
180 bias due to sample selection. For the final model, we selected symptoms that had marginal
181 posterior probabilities of 90% or more to be either positively or negatively associated with
182 positive or negative test results with mean posterior probabilities of > 0.5 or < -0.5 in the

183 Bayes regression analysis, and were reported by at least 20% of the positive or negative
184 individuals. We also performed a random forest regression analysis (Figure 3B), which
185 confirmed the importance of the selected symptoms.

186 Based on these analyses, we selected the ten symptoms impaired or lost sense of smell, throat
187 pain, loss of appetite, joint pain, runny nose, muscle pain, headache, fever higher than
188 38.5 °C, tiredness and chills. As stated above, the high representation of dry cough in SARS-
189 CoV-2 uninfected symptomatic individuals is likely due to the fact that it was communicated
190 very early on as a typical symptom of SARS-CoV-2 infection, and was therefore a criterion
191 for testing. Furthermore, we observed a strong correlation between impaired or lost sense of
192 taste and sense of smell (Pearson's $r = 0.72$). Therefore, neither dry cough nor loss of taste
193 were included in the final Bayes model.

194 We used the selected ten symptoms to fit the final Bayesian regression model (Figure 4A). On
195 the tested cohort data set, the model gave an area under the receiver operating characteristic
196 curve of 0.80 (AUC; Figure 4B), a positive predictive value (PPV) of 0.80 (95% credible
197 interval: 0.74 – 0.84) and a negative predictive value (NPV) of 0.72 (95% credible interval:
198 0.68 – 0.77). We applied the Bayesian model to the non-tested individuals who participated in
199 the survey. The characteristics of this second study population are represented in Table 1. In
200 total, responses of 520 participants who had not been tested for SARS-CoV-2 infection were
201 included in the survey. Using the Bayesian model described above, we calculated that 56
202 (10.8%) were SARS-CoV-2 positive with a probability of more than 75% (Figure 4C).

203 Another 84 individuals (16.2%) were predicted as SARS-CoV-2 positive with a probability
204 between 50% and 75%. For the other 380 individuals (73.1%), the probability of SARS-CoV-
205 2 infection based on the Bayesian model was less than 50%.

206

207 **Discussion**

208 The data presented here provide a very detailed picture of the symptoms experienced by
209 SARS-CoV-2 infected individuals. To our knowledge, we are providing the most detailed
210 symptoms overview reported so far, having queried 45 symptoms ranging from general
211 symptoms to specific eye, nose, throat, respiratory and gastrointestinal symptoms.

212 The data collected in this survey show that a range of symptoms is observed in a majority of
213 patients, such as an impaired or lost sense of smell or taste, tiredness, lethargy, loss of
214 appetite, joint or muscle pains and dry coughing. For some symptoms, we found surprisingly
215 low frequencies in infected individuals and therefore low odds ratios, including all throat
216 symptoms other than dry cough. Furthermore, only roughly one fifth of SARS-CoV-2
217 infected individuals experienced throat pain, which could be expected to be much higher since
218 virus is routinely identified in throat swabs, showing that it is indeed infecting throat tissue.

219 For some symptoms, a bias may have been introduced into the dataset at two levels: some
220 symptoms had been described early in the pandemic as typical for SARS-CoV-2 infection,
221 and have therefore been criteria for testing, and on the other hand, the presence of these
222 symptoms and the testing procedure itself may have increased the awareness of the affected
223 individuals and increased their motivation to participate in this study. We expect that such a
224 bias might exist for the symptoms “dry cough” and “difficulty breathing”. The odds ratios
225 reported here may therefore be underestimates.

226 While many of the symptoms reported by SARS-CoV-2 positive individuals taken on their
227 own are not useful for a prediction of SARS-CoV-2 infection, a prediction model based on 10
228 symptoms as we present here can provide a good estimate of the probability of a SARS-CoV-
229 2 infection. The model therefore may provide useful guidance for personal behavior when
230 symptoms are experienced as well as in the allocation of tests. It has to be noted though that
231 while the positive predictive value of the model we describe here is rather high (0.80), the

232 negative predictive value is lower (0.72). When used in the future for personal risk
233 assessment, a predicted low probability of SARS-CoV-2 infection should therefore still be
234 interpreted with caution and personal protection measures should be maintained.

235 Other studies have been performed investigating symptoms of SARS-CoV-2 and their
236 predictive value. Studies in the early days of the COVID-19 pandemic were reporting
237 symptoms of hospitalized patients and thus missed the symptoms that are associated with
238 mild disease courses [5, 6, 16]. Since then, multiple other studies have been performed that
239 focused on the symptoms experienced by SARS-CoV-2 infected individuals at a community
240 level [8, 9, 17-21]. A large study has been performed using a mobile phone app that had more
241 than seven thousand SARS-CoV-2 tested participants from the United Kingdom and the
242 United States [19]. In this study, the authors found frequencies for the symptoms “loss of
243 smell and taste”, “fever”, “skipped meals” and “diarrhea” that are comparable to the
244 frequencies reported by us and reported positive odds ratios for these symptoms, with the
245 highest odds ratio for “loss of smell and taste”. The authors report a similar frequency of
246 “persistent cough” in the SARS-CoV-2 positive individuals as we found for “dry cough”; in
247 contrast to our data, the frequency in SARS-CoV-2 negative individuals was lower than in our
248 cohort, leading to a low positive odds ratio for the persistent cough, in contrast to our
249 findings. The authors provided a prediction model based on age, sex, loss of smell and taste,
250 persistent cough, severe fatigue and skipped meals, which obtained a positive predictive value
251 of 0.69 and a negative predictive value of 0.75.

252 In another large study, only a small number of symptoms were surveyed: sore throat, cough,
253 shortness of breath, loss of smell or taste, and fever [20]. Interestingly, the authors report
254 distinctly lower frequencies for all these symptoms in their SARS-CoV-2 infected individuals,
255 which were mostly only half of the frequencies we observed in our study and as low as only
256 2% of respondents reporting a fever of 38°C or higher, compared to more than 33% observed

257 by us. This study also included a large number of individuals who were COVID-19
258 undiagnosed, who were ten times the number of diagnosed individuals, however more than
259 95% of the undiagnosed individuals reported to feel well and not have any of the queried
260 symptoms, creating an imbalance between the SARS-CoV-2 positive and negative study
261 population. The model proposed in that study that relies on this small number of symptoms
262 would likely prove less effective when applied to a cohort comprising more symptomatic
263 individuals.

264 In contrast to these two large-scale studies, we did not include age and sex in our models but
265 included more symptoms, which resulted in a higher positive predictive value. Some of the
266 symptoms included in our model are more subjective and may vary in strength, for example
267 tiredness and loss of appetite. However, we think that modeling a prediction on this higher
268 number of symptoms provides a more precise model that can be used for a first prediction of
269 the probability of SARS-CoV-2 infection, guiding decisions for self-isolation and testing.

270 **Conclusions**

271 Our data shows that many symptoms reported by SARS-CoV-2 infected individuals are rather
272 unspecific and on their own are not amenable to a clear distinction of SARS-CoV-2 infection
273 from infections with other viruses causing common cold symptoms. Nevertheless, our
274 Bayesian model based on 10 symptoms provides a good prediction of probability of SARS-
275 CoV-2 infection. The model can provide a probability estimate and provide guidance in test
276 allocation where testing capacities are limited.

277

278 **Declarations**

279 **Ethics approval and consent to participate**

280 The study was approved by the local Ethics Committee of the Medical Faculty of the
281 University Duisburg-Essen (approval number 20-9233-BO)

282 **Consent for publication**

283 Not applicable

284 **Availability of data and materials**

285 The datasets used and/or analyzed during the current study are available from the
286 corresponding author on reasonable request.

287 **Competing interests**

288 The authors declare that they have no competing interests.

289 **Funding**

290 The study was supported by a grant from the Stiftung Universitätsmedizin Essen to WB.

291 **Authors' contributions**

292 HS recruited participants, analyzed data and contributed to writing the manuscript. DH
293 analyzed data and contributed to writing the manuscript. WB conceived of the study, created
294 the online survey, analyzed data and wrote the manuscript.

295 **Acknowledgements**

296 The authors thank all study participants for taking the time to complete the online
297 questionnaire and contributing to our advancement in the recognition of SARS-CoV-2
298 symptoms. The authors also thank the employees of the public health offices in Hamm

299 (Germany), Soest (Germany), Hochsauerlandkreis (Germany) and Dr. Hüning and nurses of
300 the SARS-CoV-2 treatment center in Lünen (Germany) for distributing invitations to the
301 online survey.

302 References

- 303 1. **Undiagnosed Pneumonia - China (Hubei): Request for Information**
304 [<https://promedmail.org/promed-post/?id=6864153%20#COVID19>]
- 305 2. **WHO Director-General's statement on IHR Emergency Committee on Novel**
306 **Coronavirus (2019-nCoV)** [[https://www.who.int/dg/speeches/detail/who-director-](https://www.who.int/dg/speeches/detail/who-director-general-s-statement-on-ihr-emergency-committee-on-novel-coronavirus-(2019-ncov))
307 [general-s-statement-on-ihr-emergency-committee-on-novel-coronavirus-\(2019-ncov\)](https://www.who.int/dg/speeches/detail/who-director-general-s-statement-on-ihr-emergency-committee-on-novel-coronavirus-(2019-ncov))]
- 308 3. Spiteri G, Fielding J, Diercke M, Campese C, Enouf V, Gaymard A, Bella A,
309 Sognamiglio P, Sierra Moros MJ, Riutort AN *et al*: **First cases of coronavirus**
310 **disease 2019 (COVID-19) in the WHO European Region, 24 January to 21**
311 **February 2020.** *Euro Surveill* 2020, **25**(9).
- 312 4. **Weekly Operational Update on COVID-19, 6 November 2020**
313 [[https://www.who.int/publications/m/item/weekly-operational-update-on-covid-19---](https://www.who.int/publications/m/item/weekly-operational-update-on-covid-19---6-november-2020)
314 [6-november-2020](https://www.who.int/publications/m/item/weekly-operational-update-on-covid-19---6-november-2020)]
- 315 5. Guan WJ, Ni ZY, Hu Y, Liang WH, Ou CQ, He JX, Liu L, Shan H, Lei CL, Hui DSC
316 *et al*: **Clinical Characteristics of Coronavirus Disease 2019 in China.** *N Engl J Med*
317 2020, **382**(18):1708-1720.
- 318 6. Huang C, Wang Y, Li X, Ren L, Zhao J, Hu Y, Zhang L, Fan G, Xu J, Gu X *et al*:
319 **Clinical features of patients infected with 2019 novel coronavirus in Wuhan,**
320 **China.** *Lancet* 2020, **395**(10223):497-506.
- 321 7. Miyamae Y, Hayashi T, Yonezawa H, Fujihara J, Matsumoto Y, Ito T, Tsubota T,
322 Ishii K: **Duration of viral shedding in asymptomatic or mild cases of novel**
323 **coronavirus disease 2019 (COVID-19) from a cruise ship: A single-hospital**
324 **experience in Tokyo, Japan.** *Int J Infect Dis* 2020, **97**:293-295.
- 325 8. Giacomelli A, Pezzati L, Conti F, Bernacchia D, Siano M, Oreni L, Rusconi S,
326 Gervasoni C, Ridolfo AL, Rizzardini G *et al*: **Self-reported Olfactory and Taste**
327 **Disorders in Patients With Severe Acute Respiratory Coronavirus 2 Infection: A**
328 **Cross-sectional Study.** *Clin Infect Dis* 2020, **71**(15):889-890.
- 329 9. Maechler F, Gertler M, Hermes J, van Loon W, Schwab F, Piening B, Rojansky S,
330 Hommes F, Kausch F, Lindner AK *et al*: **Epidemiological and clinical**
331 **characteristics of SARS-CoV-2 infections at a testing site in Berlin, Germany,**
332 **March and April 2020-a cross-sectional study.** *Clin Microbiol Infect* 2020.
- 333 10. Liaw A, Wiener M: **Classification and Regression by randomForest.** *R News* 1002,
334 **2**(3):18-22.
- 335 11. **rstanarm: Bayesian applied regression modeling via Stan. R package version**
336 **2.21.1** [<https://mc-stan.org/rstanarm>]
- 337 12. Gabry J, Simpson D, Vehtari A, Betancour M, Gelman A: **Visualization in Bayesian**
338 **workflow.** *Journal of the Royal Statistical Society Series A* 2019, **182**(2):389-402.
- 339 13. Sing T, Sander O, Beerenwinkel N, Lengauer T: **ROCR: visualizing classifier**
340 **performance in R.** *Bioinformatics* 2005, **21**(20):3940-3941.
- 341 14. **Voila**
342 [[https://mybinder.org/v2/gh/niklastoe/classifier_metric_uncertainty/master?urlpath=%](https://mybinder.org/v2/gh/niklastoe/classifier_metric_uncertainty/master?urlpath=%2Fvoila%2Frender%2Finteractive_notebook.ipynb)
343 [2Fvoila%2Frender%2Finteractive_notebook.ipynb](https://mybinder.org/v2/gh/niklastoe/classifier_metric_uncertainty/master?urlpath=%2Fvoila%2Frender%2Finteractive_notebook.ipynb)]
- 344 15. Toetsch N, Hoffmann D: **Classifier uncertainty: evidence, potential impact, and**
345 **probabilistic treatment.** In. arXiv; 2020.
- 346 16. Chen N, Zhou M, Dong X, Qu J, Gong F, Han Y, Qiu Y, Wang J, Liu Y, Wei Y *et al*:
347 **Epidemiological and clinical characteristics of 99 cases of 2019 novel coronavirus**
348 **pneumonia in Wuhan, China: a descriptive study.** *Lancet* 2020, **395**(10223):507-
349 513.

- 350 17. Ahamad MM, Aktar S, Rashed-Al-Mahfuz M, Uddin S, Lio P, Xu H, Summers MA,
351 Quinn JMW, Moni MA: **A machine learning model to identify early stage**
352 **symptoms of SARS-Cov-2 infected patients.** *Expert Syst Appl* 2020, **160**:113661.
353 18. Menni C, Sudre CH, Steves CJ, Ourselin S, Spector TD: **Quantifying additional**
354 **COVID-19 symptoms will save lives.** *Lancet* 2020, **395**(10241):e107-e108.
355 19. Menni C, Valdes AM, Freidin MB, Sudre CH, Nguyen LH, Drew DA, Ganesh S,
356 Varsavsky T, Cardoso MJ, El-Sayed Moustafa JS *et al*: **Real-time tracking of self-**
357 **reported symptoms to predict potential COVID-19.** *Nat Med* 2020, **26**(7):1037-
358 1040.
359 20. Shoer S, Karady T, Keshet A, Shilo S, Rossman H, Gavrieli A, Meir T, Lavon A,
360 Kolobkov D, Kalka I *et al*: **A prediction model to prioritize individuals for SARS-**
361 **CoV-2 test built from national symptom surveys.** *Med (N Y)* 2020.
362 21. Sudre CH, Lee K, Ni Lochlainn M, Varsavsky T, Murray B, Graham MS, Menni C,
363 Modat M, Bowyer RCE, Nguyen LH *et al*: **Symptom clusters in Covid19: A**
364 **potential clinical prediction tool from the COVID Symptom study app.** *medRxiv*
365 2020:2020.2006.2012.20129056.

366

367 **Tables**

368 **Table 1. Characteristics of study cohort.**

	SARS-CoV-2 positive n (% of total)	SARS-CoV-2 negative n (% of total)	untested n (% of total)
number	336	107	520
age (years)			
18 – 29	55 (16.4%)	18 (16.8%)	92 (17.7%)
30 – 39	45 (13.4%)	36 (33.6 %)	178 (34.2%)
40 – 49	61 (18.2%)	26 (24.3%)	126 (24.2%)
50 – 59	118 (35.1%)	19 (17.8%)	88 (16.9%)
60 – 69	41 (12.2%)	7 (6.5%)	29 (5.6%)
70 – 79	9 (2.7%)	1 (0.94%)	5 (1.0%)
> 80	7 (2.1%)	0 (0%)	2 (0.4%)
gender			
female	189 (56.3%)	81 (75.7%)	409 (78.7%)
male	147 (43.8%)	25 (23.4%)	109 (21.0%)
diverse	0 (0%)	1 (0.9%)	2 (0.4%)
smoking	13.1%	23.4%	30.2%
hayfever*	19.1%	23.4%	25.6%

369 * hayfever relevant at the time of infection

370

371 **Figure Legends**

372 **Figure 1 Heatmap of symptoms reported by SARS-CoV-2 positive and negative** 373 **individuals.**

374 Symptoms reported by SARS-CoV-2 positive (top) and SARS-CoV-2 negative individuals
375 (bottom). Subjects were sorted according to responses concerning severity and presence of
376 fever. feverish/not feverish: subjective judgement, participant did not measure their body
377 temperature.

378 **Figure 2 Frequencies and odds ratio of symptoms**

379 (A) Frequencies for each symptom were calculated for SARS-CoV-2 positive and for SARS-
380 CoV-2 negative individuals, respectively. (B) Odds ratios of individual symptoms and 90%
381 confidence intervals.

382 **Figure 3 Bayesian regression and Random Forest-based selection of modelling factors**

383 (A) For visual selection of symptoms for predictive modeling, the frequency of symptoms in
384 SARS-CoV-2 negative individuals was plotted against the frequency in SARS-CoV-2 positive
385 individuals for individual symptoms, odds ratio (OR) was visualized as dot size, Bayes model
386 mean posterior probability estimate was visualized as color as indicated. Selected symptoms
387 in bold letters. (B) Random Forest logistic regression analysis was performed to identify
388 individual symptoms with importance for the development of a predictive model; the figure
389 shows the mean decrease of the Gini index for the individual symptoms. Selected symptoms
390 in bold letters.

391 **Figure 4 Heatmap of predicted probability of SARS-CoV-2 infection and of symptoms** 392 **reported by non-tested individuals**

393 (A) Bayesian regression analysis was performed using the indicated symptoms as priors,
394 mean estimates and 90% highest posterior density intervals are shown. (B) Performance of the

395 Bayesian model and the Random Forest analysis are shown as receiver operating
396 characteristic curves; AUC (area under the curve), PPV (positive predictive value) and NPV
397 (negative predictive value) with 95% credibility intervals of the Bayesian model on the model
398 dataset of tested individuals are shown. (C) The dataset of non-tested individuals (n = 520)
399 was submitted to the Bayesian model, and symptoms are shown for all individuals sorted from
400 highest predicted probability (p_B = Bayesian model prediction; top) to lowest predicted
401 probability (bottom).

Figure 1

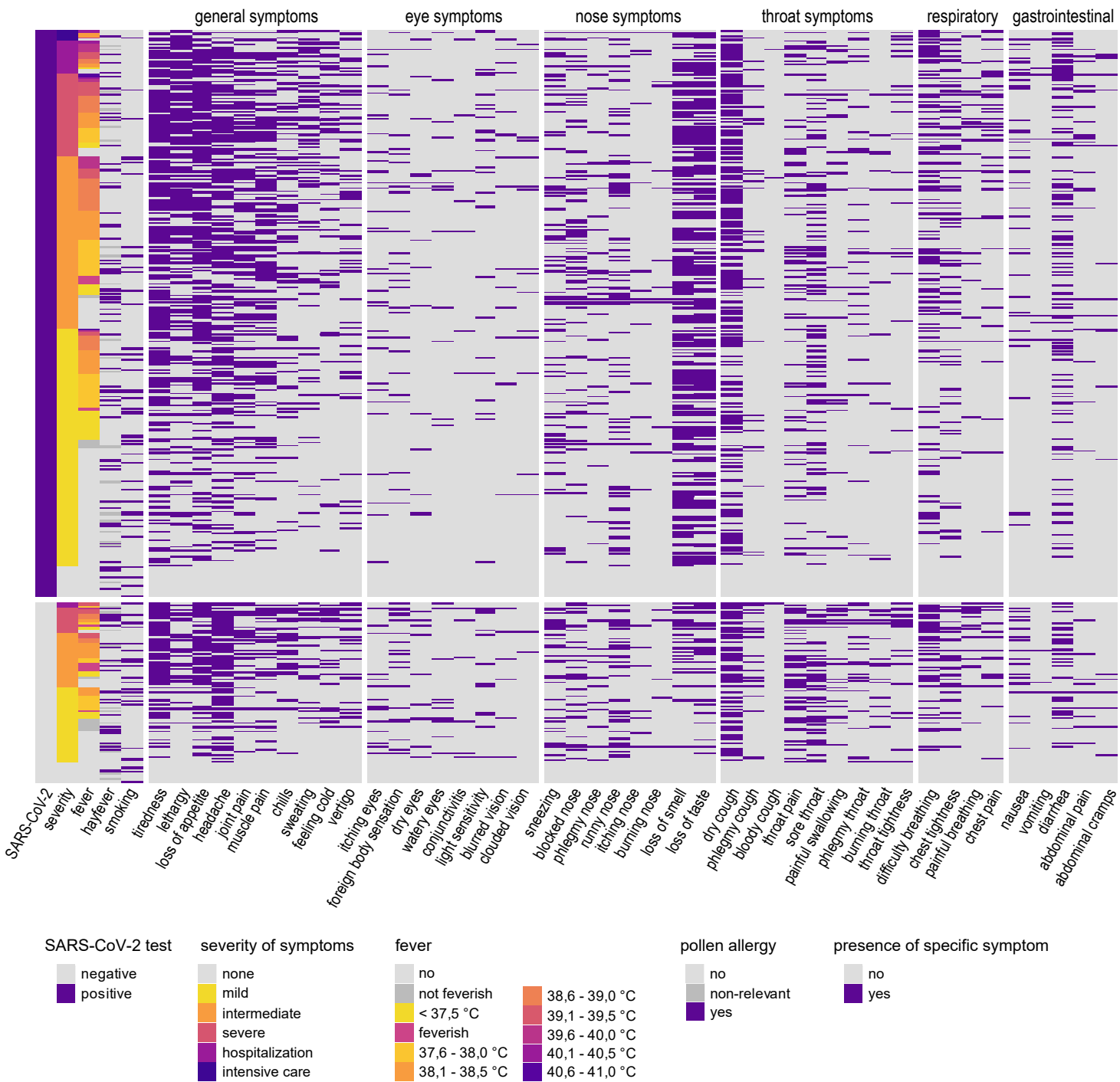


Figure 2

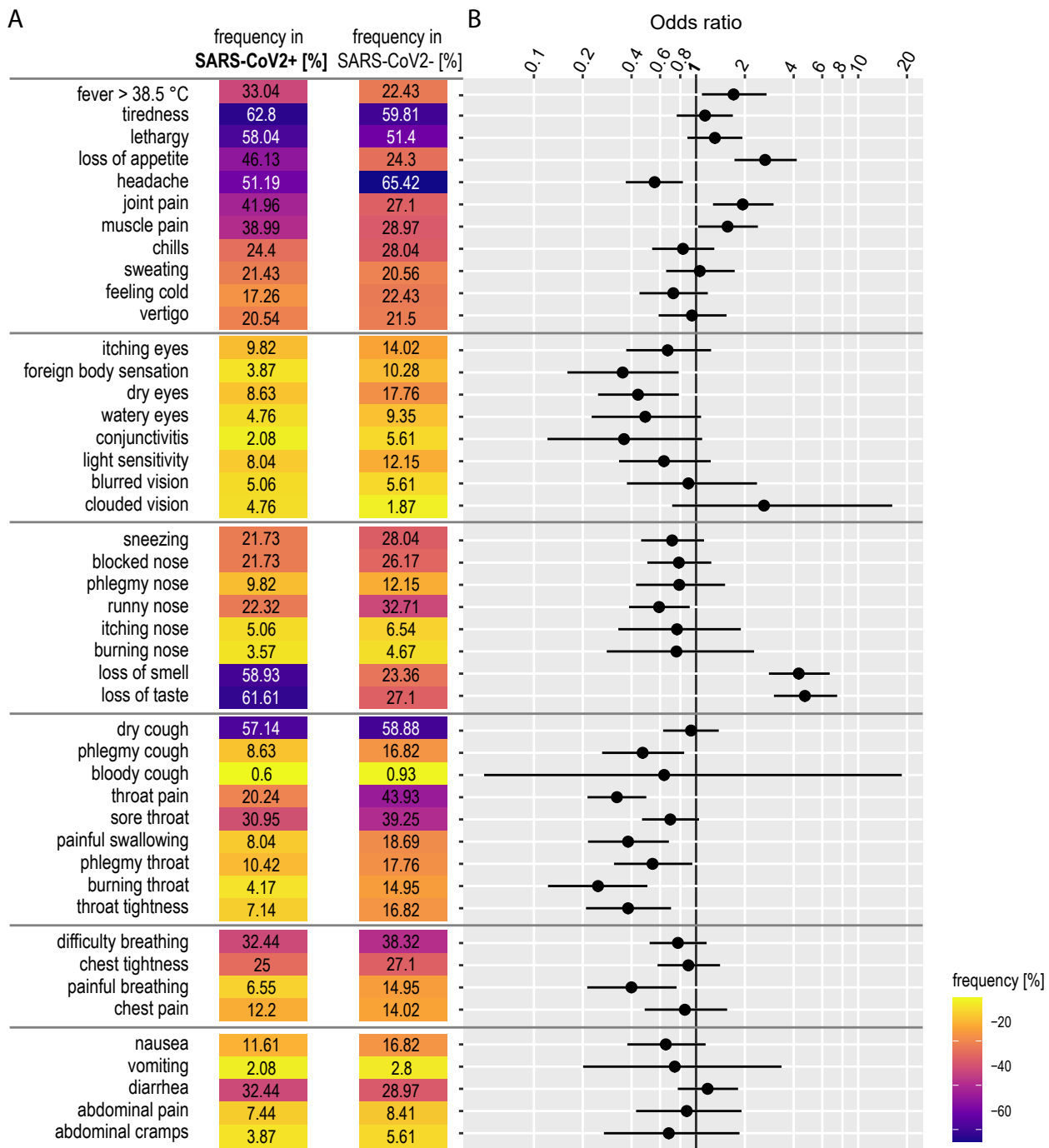
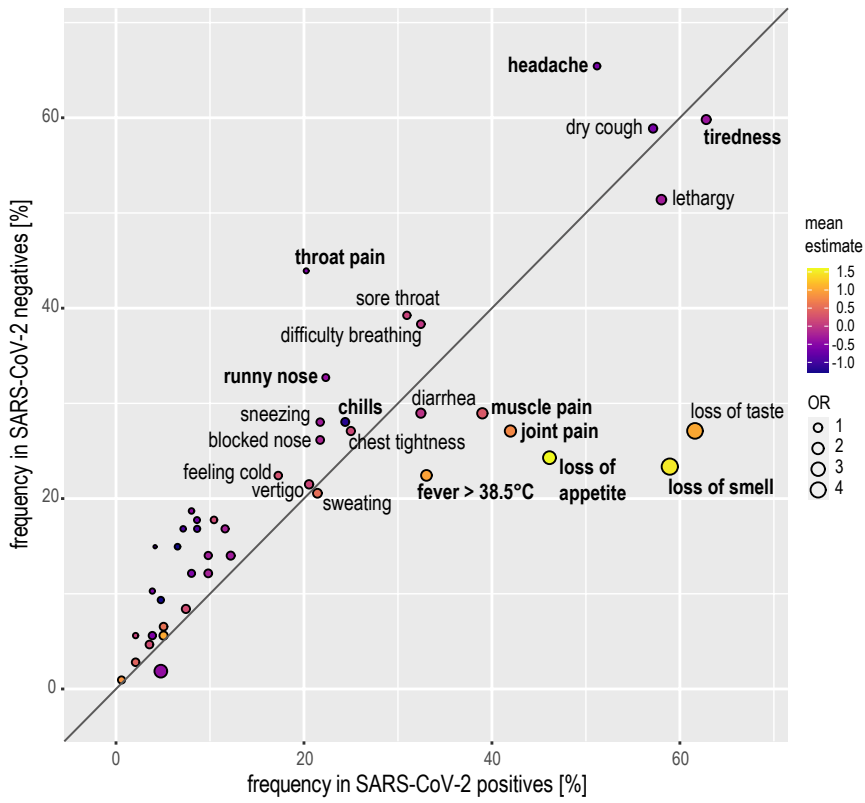


Figure 3

A



B

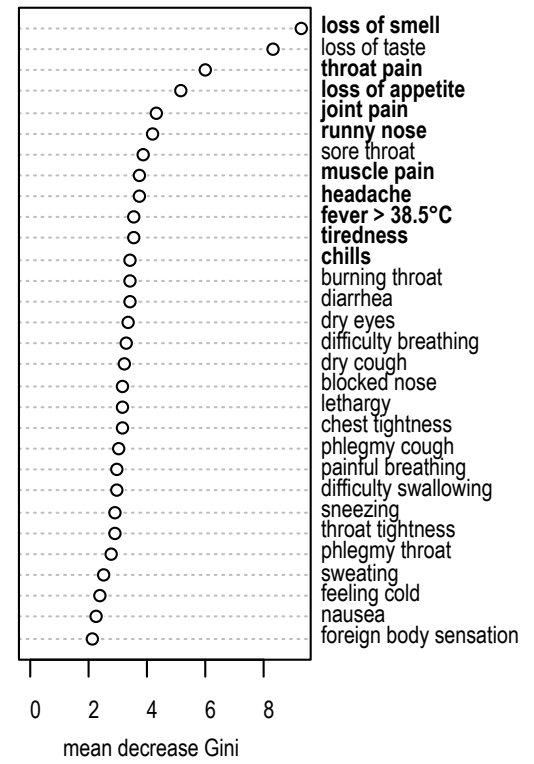


Figure 4

