

1 **Uncertainty and Inconsistency of COVID-19 Non-Pharmaceutical** 2 **Intervention Effects with Multiple Competitive Statistical Models**

3 Bernhard Müller

4 *School of Physics and Astronomy, Monash University, Clayton, VIC 3800, Australia*

5 Inken Padberg

6 *Epidemiology Unit, German Rheumatism Research Centre (DRFZ),*
7 *Charitéplatz 1, 10117, Berlin, Germany*

8 Michael Lorke

9 *Faculty of Physics, University of Duisburg-Essen, 47057 Duisburg, Germany*

10 Ralph Brinks

11 *Chair for Medical Biometry and Epidemiology Witten/Herdecke University,*
12 *Faculty of Health/School of Medicine D-58448 Witten, Germany*

13 Sally Cripps

14 *Human Technology Institute (HTI), University of Technology Sydney, Sydney, NSW, Australia*

15 M. Gabriela M. Gomes

16 *Department of Mathematics and Statistics,*
17 *University of Strathclyde, Glasgow, United Kingdom and*
18 *NOVA School of Science and Technology,*
19 *Centre for Mathematics and Applications (NOVA MATH), Caparica, Portugal.*

20 Daniel Haake

21 *Independent Researcher, D-14469 Potsdam, Germany*

22 John P. A. Ioannidis

23 *Departments of Medicine, of Epidemiology and Population Health,*
24 *and of Biomedical Data Science, and Meta-Research Innovation Center at Stanford (METRICS),*
25 *Stanford University, 3180 Porter Dr, Room A129,*
26 *Stanford Research Park, Palo Alto, CA 94304, USA*

27 (Dated: Accepted XXX. Received YYY; in original form ZZZ)

Abstract

Quantifying the effect of non-pharmaceutical interventions (NPIs) is essential for formulating lessons from the COVID-19 pandemic. To enable a more reliable and rigorous evaluation of NPIs based on time series data, we reanalyse the data for the original official evaluation of NPIs in Germany using an ensemble of 9 competitive statistical methods for estimating the effects of NPIs and other determinants of disease spread on the effective reproduction number $\mathcal{R}(t)$ and the associated error bars. A proper error analysis for time series data leads to significantly wider confidence intervals than the official evaluation. In addition to vaccination and seasonality, only few NPIs – such as restrictions in public spaces – can be confidently associated with variations in $\mathcal{R}(t)$, but even then effect sizes have large uncertainties. Furthermore, due to multicollinearity in NPI activation patterns, it is difficult to distinguish potential effects of NPIs in public spaces from other interventions that came into force early, such as physical distancing. In future, NPIs should be more carefully designed and accompanied by plans for data collections to allow for a timely evaluation of benefits and harms as a basis for an effective and proportionate response.

28 INTRODUCTION

29 The COVID-19 pandemic has arguably been the most globally disruptive event of the 21st
30 century so far. In the aftermath of the pandemic, there is now considerable interest to revisit the
31 handling of the crisis and derive lessons for better responses to similar events in the future.

32 The effectiveness of interventions to influence the spread of COVID-19 is a key piece of the
33 puzzle in deriving such lessons. Non-pharmaceutical interventions (NPIs) of an unprecedented
34 scale were implemented during the pandemic with substantial collateral effects and at the sig-
35 nificant expense of civil liberties. A proper evaluation of both the benefits and harms of such
36 interventions is required since the proportionality of the response is central to the formulation of
37 pandemic strategy [1].

38 The scientific literature on NPI effects is vast, with enormous heterogeneity in methodology,
39 quality, and reported conclusions. Official reports on the efficacy of NPIs have therefore been
40 commissioned in countries such as the UK [30], Switzerland [3], and Germany [4] to quantify
41 NPI effects by means of a literature review or a sufficiently comprehensive statistical analysis or
42 meta-analysis of available data and inference models. Similar efforts have also been undertaken
43 outside of official government reviews [5]. Such independent evaluations of NPIs are essential to
44 establish the complete picture and some form of scientific consensus view on the magnitude and

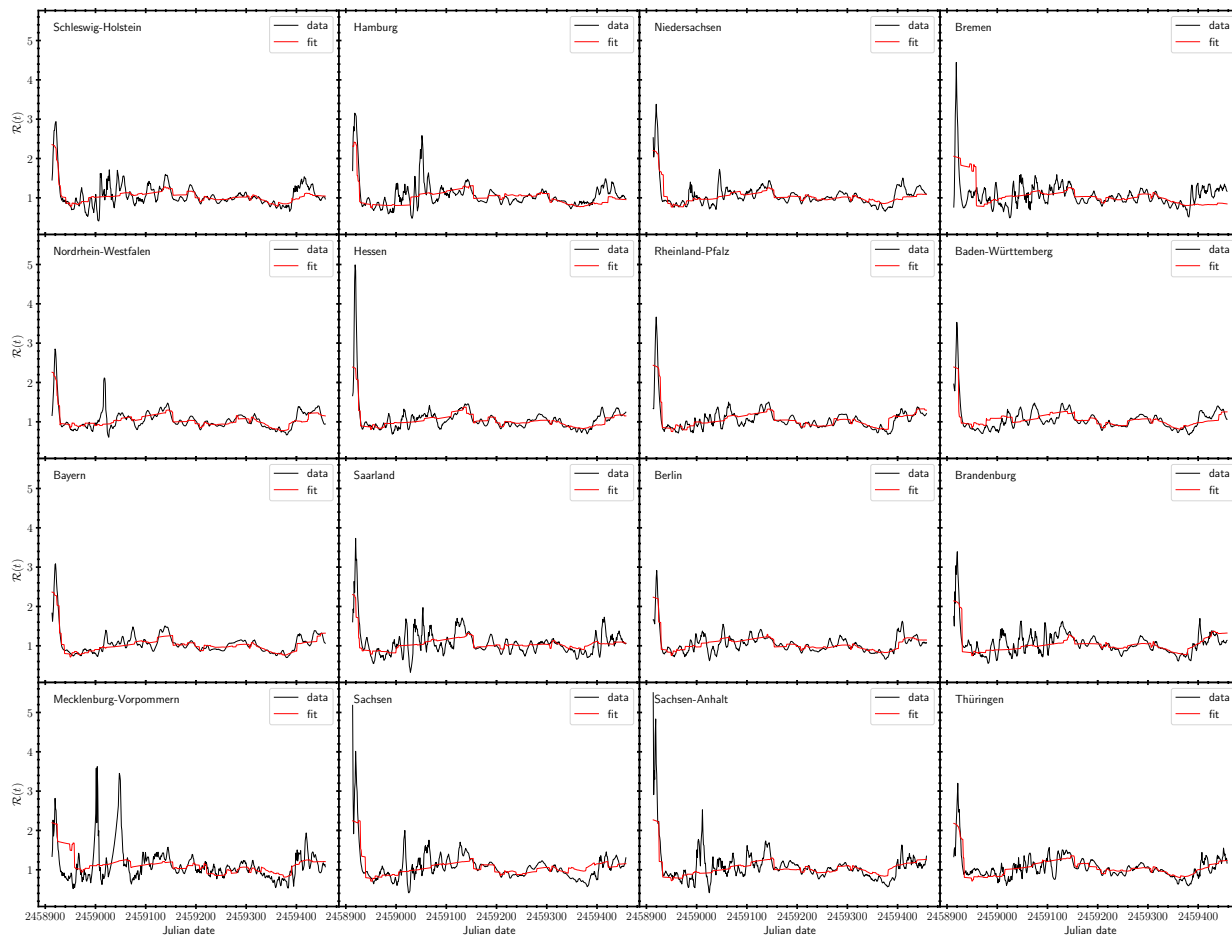


FIG. 1. $\mathcal{R}(t)$ for individual states (black) compared to the fit of the baseline model (red).

45 uncertainties of NPI effects.

46 Beyond conducting evaluations in additional countries and with large-scale evidence synthesis
47 across countries, it is also important to perform independent verification and validation (IV&V)
48 of every step of the analysis pipeline for influential government assessments conducted to-date.
49 Moreover, ensemble modelling – employing multiple models rather than just one – may help
50 achieve greater objectivity and, under appropriate circumstances, also better predictive perfor-
51 mance than any single model by model averaging [6]. Such an approach can help to better assess
52 the robustness of NPI effect estimations and systematically identify key sources of uncertainty that
53 need to be addressed by follow-up studies. Some IV&V exercises have been conducted during the
54 pandemic on epidemiological models for prediction and inference, and have provided important
55 insights on their uncertainties and sensitivities [7, 8].

56 We here conduct independent verification and validation for the evaluation of government in-
57 terventions on disease spread in Germany by the *StopptCOVID* project, which was commissioned

Model	Model type	Errors	Shrinkage
Baseline	Linear regression	Standard errors	—
DK	Linear regression	Driscoll-Kraay	—
Ebisuzaki	Linear regression	Ebisuzaki	—
BT	Linear regression	Stationary bootstrap	—
2WFE	Linear regression, two-way fixed effects	Stationary bootstrap	—
ARMA(p,q)	Linear regression	ARMA(p,q) errors	—
DYN	Renewal equation	Stationary bootstrap	—
RF	Random forest regression	Stationary bootstrap (cases)	—
Elastic net	Linear regression	Stationary bootstrap	Elastic net
PCR	Linear regression	Stationary bootstrap	Principal component regression

TABLE 1. Overview of implemented models.

58 by the German Federal Ministry of Health in 2020 [4] to be carried out by the Robert-Koch Insti-
59 tute (RKI, German Centre for Disease Control) and external collaborators. Results were published
60 as a non-refereed report and released to the press in mid-2023 [1]. Despite critiques of the study
61 methodology [10, 11] the underlying data and the analysis code were not initially made pub-
62 lic. They were finally made available to the community in early April 2024 [12] after significant
63 political pressure on the German Health Ministry for transparency in a matter of major societal
64 relevance [13].

65 After the release of the data, we set up an IV&V project to re-examine key findings of the study,
66 both in recognition of the added benefit of a completely independent reanalysis, and of RKI’s
67 limited resources to perform more extensive verification and validation. The IV&V exercise seeks
68 to provide updated estimates and error bars for effects of NPIs and other selected determinants of
69 disease spread on $\mathcal{R}(t)$, but the aims are not limited to IV&V. The project also seeks to elucidate and
70 compare capabilities and limits of commonly used inference techniques for NPI evaluation, and
71 to identify gaps in data for assessing NPI effectiveness and potential remedies (e.g., the need for
72 certain experimental studies). In keeping with the *Guidelines for Accurate and Transparent Health*
73 *Estimates Reporting* (GATHER) [14], we therefore provide exhaustive supplementary materials
74 containing a conceptual overview of the analysis methods and a description of methods to calculate
75 uncertainties (Supplementary Methods S3–S5).

76 RESULTS

77 *StopptCOVID* estimated the effect of various NPIs, holidays, vaccination, and seasonality
78 (modelled as a cosine and sine function modulation) on the logarithm of the effective reproduction
79 number $\mathcal{R}(t)$ using weighted least squares (WLS) linear regression (Equation 2). The results of
80 this baseline model were replicated to very high accuracy. However, based on the model fits and
81 the NPI activation patterns, we diagnosed two important concerns with the statistical analysis.

82 Statistical Concerns with the Baseline Model

83 The fitted time series for individual states are shown in Figure 1. We note that despite a co-
84 efficient of determination of $R^2 = 0.831$, the data show substantial dynamics that are not re-
85 produced by the fit. Visual inspection already shows that the errors display *autocorrelation* (i.e.,
86 non-independent residuals at adjacent data points), which violates the assumption made by *Stoppt-*
87 *COVID* (see also Supplementary Figure S1). Formally, strong autocorrelation is indicated by very
88 low values of the Durbin-Watson statistic [15] of the residuals around 0.2 or less for all federal
89 states. The presence of autocorrelation in the errors terms implies that standard regression errors
90 for effect sizes do not apply and may substantially underestimate the actual errors.

91 Autocorrelation in the residuals also implies that the evolution of $\mathcal{R}(t)$ is either affected
92 markedly by *unmodelled processes* (e.g., cluster effects in networks) or by *measurement arte-*
93 *facts* (e.g., ramp-up of testing) that can produce the observed autocorrelation structure in $\mathcal{R}(t)$, or
94 by both. The fact that the residuals form highly stochastic time series may point to unmodelled
95 processes rather than observational artefacts. Unmodelled processes or measurement artefacts may
96 have a bigger impact on epidemic dynamics than suggested by the residuals if the baseline model
97 is misspecified; the residuals merely define a minimum level for the magnitude of unmodelled
98 processes or observational noise.

99 A second issue is *multicollinearity*, i.e., the presence of (strong) correlations among the ex-
100 planatory variables (i.e., NPI activation variables). Strong multicollinearity leads to a highly ill-
101 conditioned regression problem, and can result in a spurious increase, decrease, or even reversal
102 of effect sizes and inflate the estimated confidence intervals [2, 17]. The degree to which estimates
103 for the regression coefficients for explanatory variables are affected by multicollinearity can be
104 quantified by the variance inflation factor (VIF; see Supplementary Discussion S1.1 for details).

105 Empirical rules-of-thumbs are typically used to identify “serious” multicollinearity, e.g., a thresh-
106 old values of $VIF > 10$. In the NPI data set, many of the included NPI variables are subject to
107 severe multicollinearity, and some of the VIFs exceed 100.

108 Given these two problems, there are concerns that the confidence intervals and point estimates
109 from *StopptCOVID* are not valid simply from a statistical perspective. Further limitations due
110 to the epidemiological assumptions are reviewed in Supplementary Discussion S1. This study
111 addresses the statistical concerns as Work Package 1 of the IV&V project.

112 **Model Ensemble**

113 To obtain more reliable point estimates and confidence intervals despite the presence of au-
114 tocorrelation and multicollinearity, we use an ensemble of 9 different competitive methods for
115 estimating effects and error bars. These were selected based on a survey of statistical analysis
116 methods for panel and time series data in the NPI studies reviewed by Murphy et al. [30], sup-
117 plemented by a wide consultation of the technical literature in relevant disciplines. In addition
118 to the baseline model – WLS with default, non-robust standard errors – the ensemble includes
119 WLS with Driscoll-Kraay errors [73], WLS with errors based on Ebisuizaki’s method [93], regres-
120 sion with autoregressive moving average (ARMA) errors, WLS with stationary bootstrap errors
121 [88], and, also with stationary bootstrap errors, two-way fixed effects WLS, a renewal equation
122 model, elastic net regression, principal component regression, and random forest regression as a
123 machine learning technique. Table 1 shows an overview of the model types, error analysis, and (if
124 applicable) shrinkage methods for handling multicollinearity.

125 **Injection-Recovery Test**

126 To gauge the sensitivity of the models, we run all the models on synthetic panel data assuming
127 that the baseline model for $\mathcal{R}(t)$ is correct (for want of experimental data with known effect sizes).
128 This approach (injection-recovery test) is a standard method for determining the sensitivity of
129 analysis methods for time series and other complex data [e.g., 21, 22]. We sort the models into
130 two groups based on whether they can recover the hypothetical effect sizes from *StopptCOVID*
131 without bias (Group A) or not (Group B).

132 Results of the injection-recovery test are shown in Figures 2 and 3. Models DK (Driscoll-Kraay

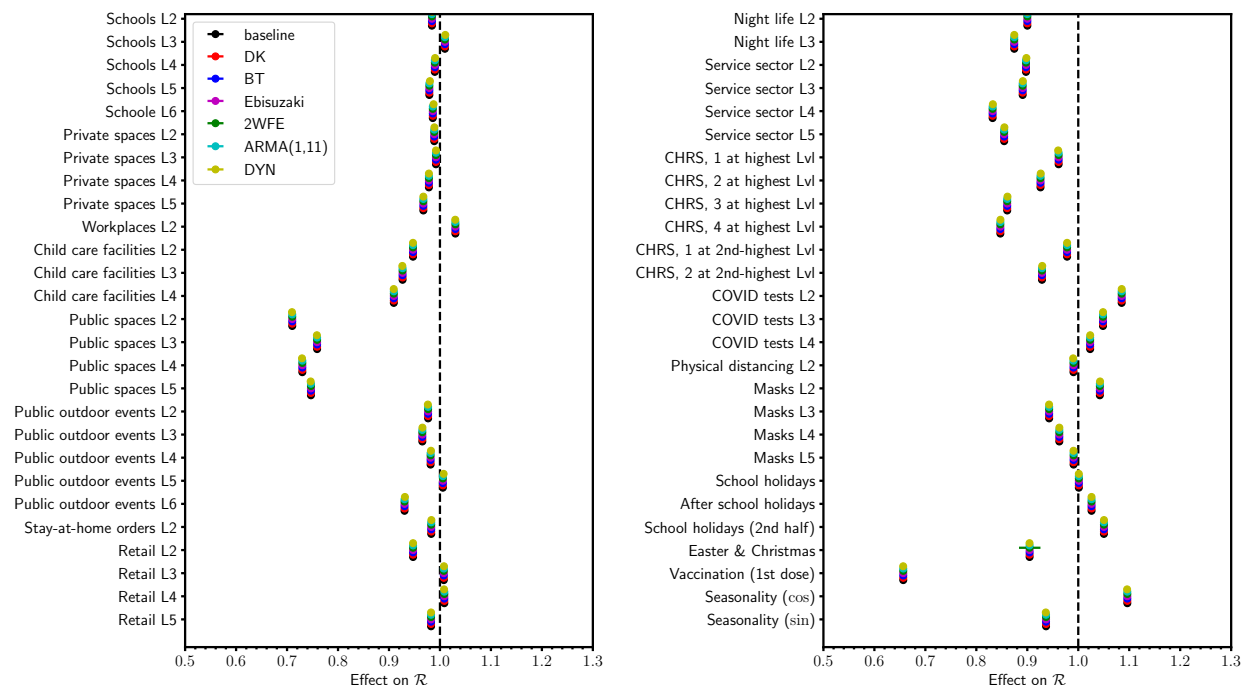


FIG. 2. Injection-recovery test for model Group A. For the models in this group, this is merely a sanity check for correct implementation. Numbers Lx indicate stringency levels of NPIs. Note that CHRS is a combined category for the cultural sector, the hotel and restaurant industry and sports.

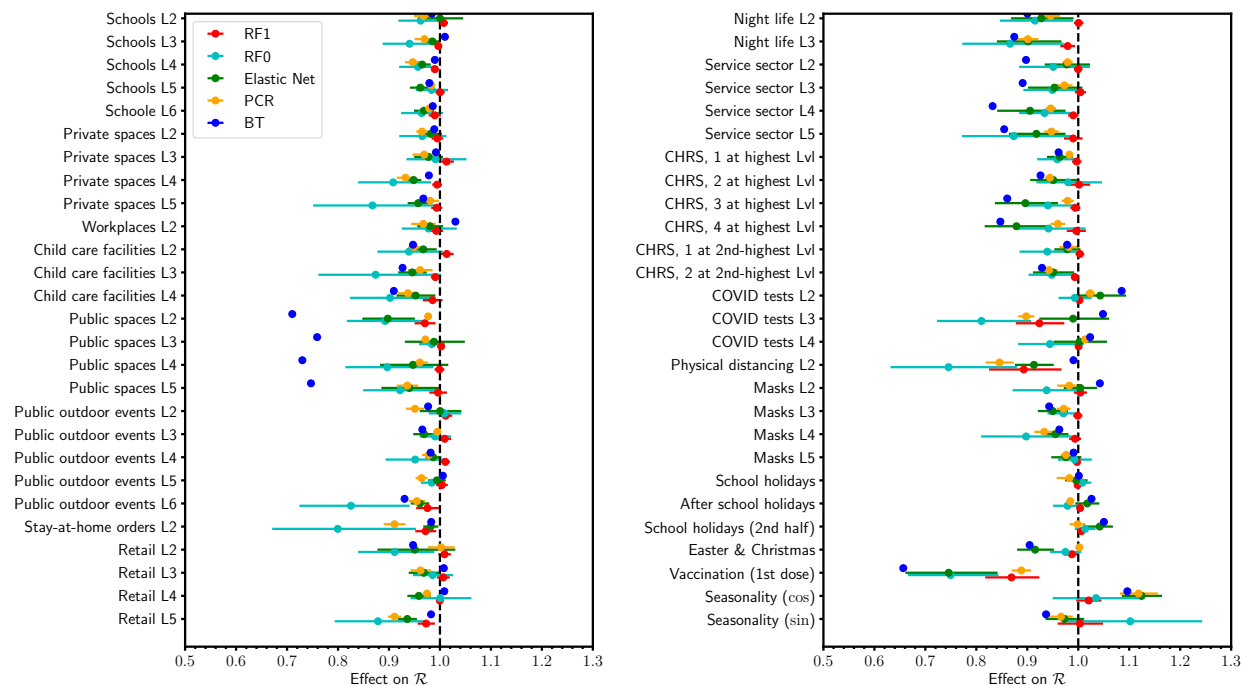


FIG. 3. Injection-recovery test for model Group B. The injection-recovery test shows that these models cannot exactly recover the effect sizes for the baseline model, but are subject to bias.

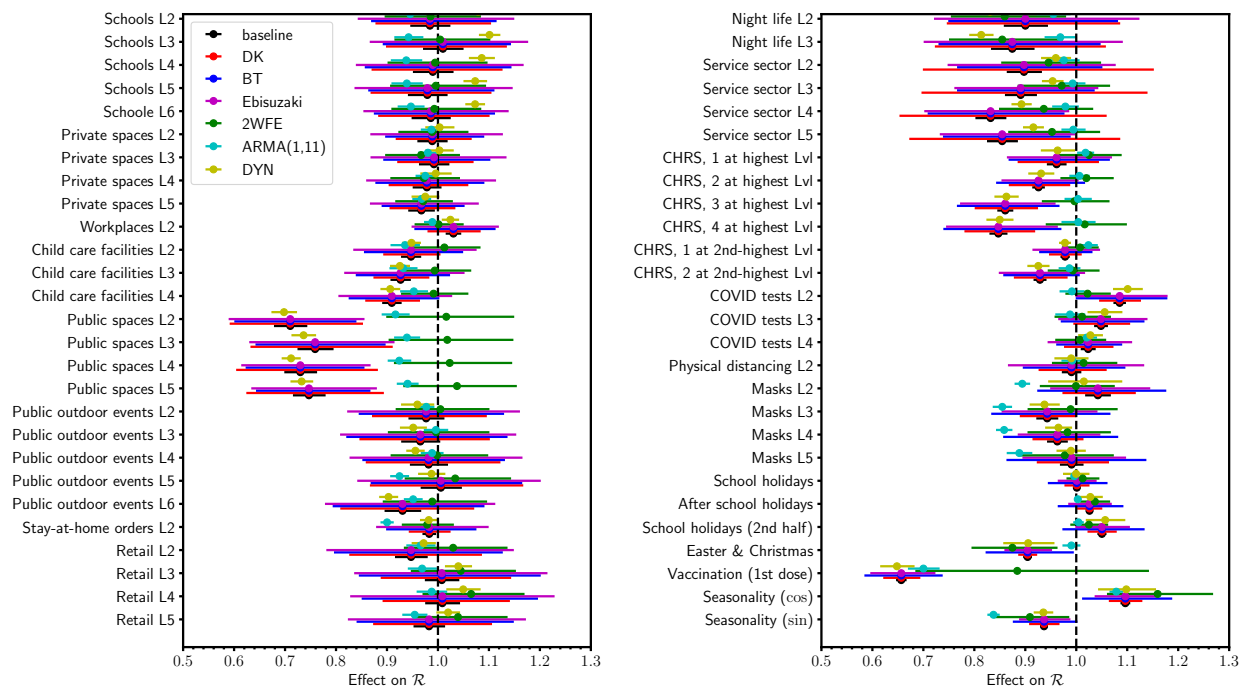


FIG. 4. Point estimates and confidence intervals for model Group A.

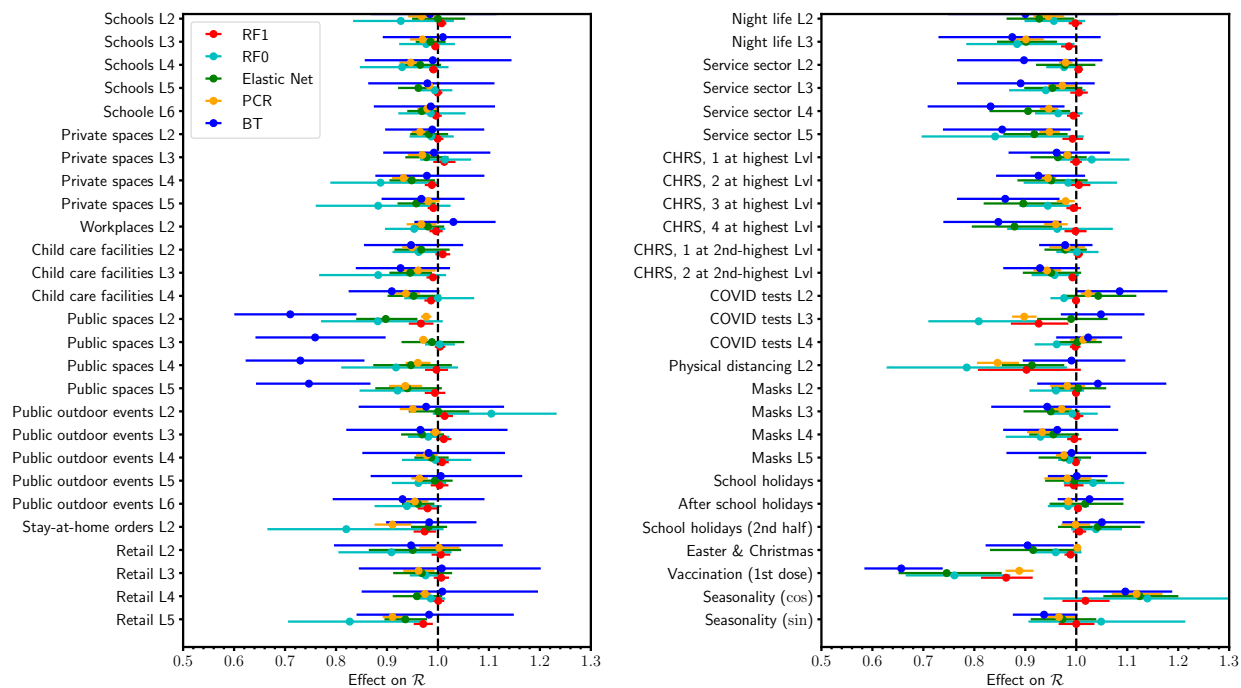


FIG. 5. Point estimates and confidence intervals for the (biased) models in group B, compared to linear regression with bootstrap errors.

133 errors), BT (bootstrap errors), Ebisuzaki, 2WFE (two-way fixed effects) and ARMA(1,11) by con-
134 struction recover the effects of the baseline model exactly, and hence this test is only a sanity check
135 for their implementation. As it is formulated for consistency with the baseline model, the dynam-
136 ical model DYN also recovers these effects without bias. With the best-fit hyperparameters from
137 cross-validation, models RF, Elastic Net and PCR all exhibit bias in the injection-recovery test and
138 are therefore assigned to Group B. For the largest effects in the baseline model (such as public
139 spaces L2-L5 and vaccination), the bias is consistently towards the null. In some cases (e.g.,
140 COVID tests L3, the sign of the effect in the baseline model is inverted). One notable exception
141 concerns physical distancing L2, to which Group B models consistently ascribe a relatively large
142 effect. Random forest regression ascribes a reduction of about 25% in $\mathcal{R}(t)$ to physical distanc-
143 ing relative to the state without interventions (case RF0). The reduction of the effects for public
144 spaces and the increased effect of physical distancing L2 is essentially a reassignment of effects
145 within groups of NPIs that were in place (at some level) quite consistently from the early phase of
146 the pandemic.

147 Random forest regression yields different linear effect estimates depending on the reference
148 state for the linearisation of the model. It tends to underestimate effects in the presence of other
149 NPIs with their actual activation profile (case RF1). By contrast, it overestimates the effect that
150 many NPIs would have as single intervention (case RF0). For example, case RF0 yields large
151 effects – though with big error bars – for public outdoor events L6, stay-at-home-orders and
152 physical distancing. Clearly, neither the RF0 nor the RF1 estimates from model RF are satis-
153 factory in the injection-recovery test. Further analysis of the behaviour of model RF is provided
154 in Supplementary Discussion S9. This analysis also reveals that the effect sizes for the sine and
155 cosine component in random forest regression *cannot* be interpreted as amplitudes of the seasonal
156 variation, and will therefore be discarded in the subsequent discussion.

157 The bias in the Group B models does not render these models incorrect for effect estimation.
158 The Elastic Net and PCR models will, by construction, achieve a more *parsimonious* fit to the
159 actual data for $\mathcal{R}(t)$ and *may* filter out noise in the effect estimates. However, because of the risk
160 of bias in Group B, Group A forms the principal basis for our conclusions, while Group B serves
161 to aid and temper the interpretation of the inferred statistical associations.

162 In particular, we assess random forest regression to be of limited usefulness for extracting linear
163 effect sizes of individual NPIs for the time being. The model *may* correctly perceive saturation
164 effects, i.e., little additional effect by single NPIs when many others are switched on already, and

165 its effect estimates in some sense correctly reflect the difficulty of distinguishing the effects of
166 NPIs with similar activation patterns. Random forest regression may be useful for generating
167 hypotheses for non-trivial interactions between NPIs, which would need to be investigated using
168 additional, independent data. More work is called for before using it routinely for determining
169 intervention effects.

170 **Effect estimates and confidence intervals**

171 Effect estimates and confidence intervals for model groups A and B are shown in Figure 4 and
172 5, respectively, and are also listed in Supplementary Table S3. To facilitate comparison between
173 the two groups, the results for model BT (linear regression with bootstrap errors) are also included
174 in Figure 5.

175 *Model Group A*

176 Within Group A, the effect sizes in the linear regression models DK, BT, and Ebisuzaki trivially
177 agree with the baseline model. The confidence intervals for these models are much wider than for
178 the baseline model from the original *StopptCOVID* study. Confidence intervals calculated using
179 Driscoll-Kraay errors, Ebisuzaki's method, or the stationary bootstrap are generally very similar.

180 Model 2WFE generally yields similar effect estimates and confidence intervals, with a few no-
181 table exceptions. First, the estimated effects of NPIs for public spaces and several NPI levels for
182 the cultural sector, the hotel and restaurant industry and sports (CHRS) are close to the null, and in
183 some cases, the confidence intervals do not overlap with the DK, BT, and Ebisuzaki models. The
184 only NPIs for which the model yields a significant beneficial effect are those for night life (L2–3).
185 Second, the estimated vaccine effect is considerably smaller, but with a very large confidence inter-
186 val that overlaps with those of models DK, BT, and Ebisuzaki. The reason for the larger error bar
187 lies in the close synchronisation of vaccination across the federal states. The two-way fixed effects
188 model is, in a sense, optimised to detect effect based on difference between the response variable
189 and explanatory variables across entities and therefore struggles to deliver a precise estimate. We
190 take this into account by constructing assessment criteria that are robust to such an outlier result.

191 Regression with ARMA(1,11) errors gives similarly narrow confidence intervals, and in a few
192 cases even narrower confidence intervals, than the baseline model, but the interval for vaccination

193 remains substantially wider. The effect estimates, however, often differ markedly from the baseline
194 model. The effect estimates for public spaces are considerably smaller. In turn, the model ascribe
195 more than a 10% reduction in \mathcal{R} to masks and stay-at-home orders and small, but statistically
196 significant effects to a few others NPIs, e.g., in schools.

197 The DYN model (renewal equation) tends to yield effect estimates within the error bars of
198 the DK, BT, and Ebisuzaki models. Thus, fitting the case data instead of $\mathcal{R}(t)$ yields relatively
199 consistent results within these “safe” error bars. However, the confidence intervals for model
200 DYN are often as narrow as for the baseline model, sometimes even narrower (e.g., for public
201 outdoor events), and sometimes wider (e.g., for masks). The confidence intervals often do *not*
202 overlap with the baseline model. This suggests that the bootstrapping procedure used for model
203 DYN does not yet fully account for autocorrelation.

204 Overall, however, the Group A confidence intervals for most explanatory variables overlap well
205 and the scatter between the point estimates of different models tends to be bounded by the DK, BT,
206 and Ebisuzaki error bars. Notable exceptions include the NPIs for public spaces, some NPIs for
207 the cultural sector, the hotel and restaurant industry and sports (CHRS), vaccination (with the DYN
208 and 2WFE models as outliers) and the sine component of the seasonal modulation. In the case of
209 the vaccine effect, the relative uncertainty due to the between-model scatter is modest compared
210 to the large effect size, however.

211 Following *StopptCOVID*, the vaccine effect intends to represent the result of halving the frac-
212 tion of unvaccinated individuals. As the original implementation leads to pathological behaviour
213 in the limit of a high vaccination fraction, we also considered a modification of model DYN that
214 correctly implements the vaccine effect (Supplementary Discussion S1.3). The corrected model
215 yields a nominal vaccine efficacy of about 75% against infection, and the other effect estimates
216 remain within the “safe” BT, DK and Ebisuzaki error bars. The bulk of first-dose vaccination oc-
217 curred about three months before the end of the study period, i.e., the waning of vaccine efficacy
218 played a lesser role. The estimated vaccine efficacy is therefore roughly consistent with the high
219 short-term efficacy against infection inferred by the clinical trials and cohort studies [11–13, 26],
220 especially bearing in mind that some conflation of the effects of the first- and second-dose may be
221 implicit in the *StopptCOVID* model.

222 *Model Group B*

223 The models with explicit shrinkage and the RF1 estimates from random forest regression all
224 yield significantly narrower confidence intervals than non-regularised regression with bootstrap
225 errors (which is the purpose of shrinkage in the first place). The RF1 estimates from random
226 forest regression and principal component regression tend to shrink the confidence intervals even
227 more strongly. The RF0 estimates do not show much shrinkage and sometimes result in wider
228 confidence intervals than linear regression with bootstrap errors.

229 With regard to point estimates, there is a rough tendency of elastic net and principal component
230 regression to shrink strong effects in the baseline model to modest or small effects, and to magnify
231 a few small effect estimates. As in the recovery-injection test, random forest regression tends to
232 yield very small effects for many NPIs if the actual NPI activation is used as reference state (case
233 RF1), but some very large effects for single NPIs without any other concurrent interventions (case
234 RF0).

235 The prominent cases where the models with significant shrinkage of the confidence intervals
236 yield *smaller* effects outside the BT error bars are the NPIs for public spaces – as the ARMA(1,11)
237 and 2WFE models in the previous subsection – and vaccination. Exactly as in the injection-
238 recovery tests, the models with shrinkage prefer to attribute a greater effect to physical distancing,
239 and to some extent to policy COVID tests L3.

240 **Ranking of Effects**

241 Despite considerably wider error bars than in *StopptCOVID*, some statistically significant asso-
242 ciations of variations in \mathcal{R} with interventions or environmental factors can be detected. Visually,
243 the effects of vaccination and seasonality emerge most clearly. For vaccination, only model 2WFE
244 has a confidence interval that overlaps with zero, which we consider an outlier for reasons de-
245 scribed above. For the sine component of seasonality, only one confidence interval marginally
246 overlaps with the null, and all Group A models shows a significant cosine modulation. This is re-
247 inforced by the regularised regression models (elastic net and PCR), which also show a significant
248 cosine component.

249 For a more quantitative identification of the NPIs that may be associated with lower $\mathcal{R}(t)$, we
250 define two different scores (Tables 2 and 3) to quantify how confidently a null effect can be ex-

Explanatory variable	Disagreement on sign	Overlap with null	Sum (integer score)
Seasonality (cos)	0	0	0
Vaccination (1st dose)	0	1	1
Easter & Christmas	0	1	1
Seasonality (sin)	0	1	1
Public spaces L2	1	1	2
Public spaces L4	1	1	2
Public spaces L5	1	1	2
Public spaces L3	1	1	2
Service sector L4	0	3	3
Service sector L5	0	3	3
Night life L2	0	3	3
Child care facilities L4	0	3	3
Child care facilities L3	0	3	3
Masks L3	0	3	3
CHRS, 3 at highest Lvl	1	2	3

TABLE 2. Integer score for ranking the likelihood of a real association of covariates with $\mathcal{R}(t)$ based on model Group A. Lower scores are better. Only the top-14 NPIs are shown. Scores of up to 2 are rated as potentially indicative of a real association with variations in $\mathcal{R}(t)$ (see text).

Explanatory variable	False positive risk score
Seasonality (cos)	0.002
Seasonality (sin)	0.008
Easter & Christmas	0.028
Vaccination (1st dose)	0.029
Service sector L4	0.041
Night life L3	0.052
Night life L2	0.077
Child care facilities L4	0.082
Public spaces L2	0.100
CHRS, 2 at 2nd-highest Lvl	0.102
Service sector L5	0.103
Public spaces L3	0.103
Service sector L2	0.104
Child care facilities L3	0.104

TABLE 3. False positive risk score for the 14 most highly-ranked explanatory variables based on model Group A. Scores lower than 0.1 are rated as indicative of a real association with variations in $\mathcal{R}(t)$ (see text).

251 cluded based on within-model error bars and between-model consistency. The first score is the
252 number of models that disagree with the sign of the median effect estimate across models plus the
253 number of models with confidence intervals that overlap with the null.

254 The second score (“false positive risk score”) is the average of the false-positive probabilities
255 Q in the models computed from the cumulative t -distribution,

$$Q = \int_{-\infty}^{\beta_i} t_{\nu} \left(\frac{x}{\sqrt{\text{var}\beta_i}} \right) dx, \quad (1)$$

256 where the number of degrees of freedom ν is the number of observations minus the number of
257 federal states and explanatory variables (although $\nu = \infty$ for practical purposes). In line with
258 guidelines for multi-model comparisons [27], one should avoid interpreting the resulting metric
259 as a probability. Both scores are merely heuristic scales that penalise lack of significance and
260 between-model variation. Recognising the bias in Group B models, these scores are computed
261 only for Group A, and the baseline model is also excluded because of its unrealistic error model.

262 The explanatory variables with the top-14 scores are shown in Tables 2 and 3, respectively. As
263 a cut-off for a potentially indicative association of an NPI with lower $\mathcal{R}(t)$, we tentatively suggest
264 that there should *either* be no more than two vetos (integer score of two or less), *or* that the
265 average false-positive risk score should not exceed 0.1. This allows us to accept effects as likely
266 even when a model is a clear outlier, or when a few models do not find a significant effect, but
267 all point estimates clearly cluster on one side. Along with vaccination, seasonality and Easter &
268 Christmas, this leaves NPIs for public spaces, the service sector (L4), night life and child care
269 facilities (L4) as the best candidates for associations with lower $\mathcal{R}(t)$. Among these, we rate the
270 effect of restrictions in public spaces as most statistically robust because only model 2WFE fails
271 to find a significant association with lower $\mathcal{R}(t)$ and is responsible for the lower ranking of these
272 NPIs on the second score.

273 The models with shrinkage serve to temper the scores in Tables 2 and 3, however. They indicate
274 that the effects of NPIs for public spaces and night life may be weaker and hard to distinguish
275 from the effects of physical distancing. Rather than blindly accepting the scores in Tables 2 and
276 3, further research is required to better distinguish the effects of classes of NPIs that came into
277 force early during the pandemic and are therefore tend to be assigned large effects by the models.

278 DISCUSSION

279 Our model ensemble demonstrates that the confidence intervals from the original evaluation
280 of German NPIs by *StopptCOVID* are substantially too narrow. They are neither consistent with
281 a more rigorous error analysis for linear regression, nor with the between-model variation in the
282 ensemble. *The data for $\mathcal{R}(t)$ and NPIs are insufficient for confidently assigning effects to most*
283 *NPIs.* Moreover, the confidence intervals are generally so wide that it is impossible to confidently
284 detect trends with increased stringency, contrary to claims by *StopptCOVID*. Especially for NPIs in
285 public spaces, even the point estimates do not suggest additional benefits from higher stringency
286 levels, similar to other recent claims in the literature [28].

287 It is important to point out that any inferred effects on $\mathcal{R}(t)$ still need to be translated into
288 relevant public health outcomes (e.g., total or peak hospitalisations, years or quality-adjusted years
289 of lives saved) for a satisfactory assessment of interventions. This complex task requires additional
290 information *beyond* the effects of NPIs on $\mathcal{R}(t)$. Moreover, estimation of NPI effects may have even
291 more error and uncertainty than what we estimated here, if data are unreliable, a common feature
292 in the chaotic circumstances of the COVID-19 pandemic.

293 Furthermore, the epidemic model considered in this study cannot determine feedback effects
294 and non-linear effects can impact both disease spread and the risk factors; hence constructing
295 counterfactual trajectories of $\mathcal{R}(t)$ can be misleading. For example, self-regulating behaviour in
296 the population [29, 30] or population heterogeneity [22, 24] may lead to a slow-down of disease
297 spread independent of NPIs. Time dependence, in particular waning of immunity after vaccina-
298 tion or infection is another relevant complication that can lead to counterintuitive results such as
299 “immunity debt” [33]. In the context of vaccination against seasonal respiratory diseases, higher
300 vaccine coverage can under certain conditions increase infection peaks [34–39]. Similarly unin-
301 tuitive disease dynamics – which obviously cannot be translated directly to COVID-19 – has also
302 long been studied for rubella [40, 41] and measles [42]. The implications of temporary vaccine-
303 induced immunity should only be assessed based on a sufficient understanding of the principles of
304 time-dependent disease dynamics.

305 We are forced to conclude that the approach taken by the German Federal Ministry of Health
306 is insufficient for ascertaining public health outcomes of NPIs given the revealed statistical limita-
307 tions. To enable robust evaluation and adjustments of NPIs to ensure proportionality and balance
308 benefits and harms, future public health interventions should be designed with a sufficient pre- and

309 post-intervention observation times to permit a meaningful determination of effects. This should
310 be coupled with careful epidemiological considerations and plans for relevant data collection [43].
311 Where a preliminary cost-benefit analysis suggests rough equipoise, interventions should include
312 a control group. State-of-the art time series analysis should be used to inform the required design,
313 similar to the use of sample size calculations to ensure sufficient power for other epidemiological
314 studies. Pandemic research programs should be coordinated to systematically identify knowledge
315 gaps and ensure that data for interventions effectiveness in relevant settings (hospitals, nursing
316 homes, etc.) are obtained to complement studies of population-wide disease spread.

317 Most prior literature on time-series based NPI studies does not adequately address critical sta-
318 tistical problems like autocorrelation and multicollinearity. Therefore, underestimation of NPI
319 effect uncertainties is likely a broader problem in the literature. We recommend that key results
320 be subjected to a similar reanalysis to provide reliable information for pandemic planning to pol-
321 icymakers. In future, policymakers and funders should more broadly support validation research,
322 and actively seek validation for critical policy-relevant research.

323 MATERIALS AND METHODS

324 *StopptCOVID* baseline model

325 For estimating the effects of NPIs in Germany, *StopptCovid* uses a linear regression model for
326 the logarithm of the time-dependent reproduction number $\mathcal{R}_j(t)$ for each federal state (“Bundes-
327 land”) j . NPIs are included as a set of $N_{\text{NPI}} = 51$ explanatory variables X_i on a scale from 0 to
328 1 (see below for details). In addition, the model includes two trigonometric terms for a harmonic
329 seasonal modulation of $\mathcal{R}_j(t)$ with arbitrary phase, a dependence on the fraction of vaccinated indi-
330 viduals (at least one dose), and fixed effects α_j for state j . The model assumes that $\ln \mathcal{R}(t)$ increases
331 by 0.3 and 0.6 times the share ν_α and ν_δ of the α - and δ -variant, respectively. Effects of variants
332 are not estimated but imposed manually as fixed parameters. The effects of vaccination and of
333 NPIs are assumed to occur with lags τ_{vac} and τ_{NPI} with respect to the corresponding explanatory
334 variables. Depletion of susceptibles by infections is neglected.

In terms of these explanatory variables and their regression coefficients β_i , the model for $\ln \mathcal{R}_j(t)$

in state j reads,

$$\begin{aligned} \ln \mathcal{R}_j(t) = & \alpha_j + 0.3\nu_\alpha(t) + 0.6\nu_\delta(t) + \beta_0 \cos \frac{2\pi t}{365 \text{ d}} + \beta_1 \sin \frac{2\pi t}{365 \text{ d}} - \beta_2 \log_2[1 - V(t - \tau_{\text{vac}})] \\ & + \sum_{i=3}^{N_{\text{NPI}}+2} \beta_i X_{j,i}(t - \tau_{\text{NPI}}) + \epsilon_t, \text{ with } \epsilon_t \sim \mathcal{N}(0, \sigma^2). \end{aligned} \quad (2)$$

Errors are modelled as normally distributed and uncorrelated across time and states, so that for the observed reproduction number $\ln \mathcal{R}_{j,t}^{\text{obs}}$ at discrete time indices t as response variable,

$$\ln \mathcal{R}_{j,t}^{\text{obs}} \sim \mathcal{N}(\ln \mathcal{R}_{j,t}, \sigma^2). \quad (3)$$

Weighted least-squares (WLS) regression is used for the baseline model, with weights given by 7-day averages of case numbers. This choice of weights can be justified as reducing heteroskedasticity in the observational errors.

Note that the dependence on the vaccination fraction V in this model is taken to be non-linear. This particular form of the vaccine effect is problematic, as explained in Supplementary Discussion S1.3.

StopptCOVID determines the delay between interventions and their effect by optimising the model fit based on the Akaike information criterion (AIC; 112). The optimum delay is found to be negative ($\tau_{\text{NPI}} = -1$ d) for NPIs; and τ_{vac} is found to be 5 d relative to the time of the first dose. The negative delay is problematic (Supplementary Discussion S1.4), but we accept the delays inferred by *StopptCOVID* as fixed parameters *not subject to errors* for the purpose of our statistical analysis, which is an assumption favourable to the original *StopptCOVID* model.

The study considers the time period from 1 March 2020 until 31 August 2021.

Data Sources

The effective reproduction number $\mathcal{R}_{j,t}$ is calculated from smoothed, 7-day average case data. $\mathcal{R}_{j,t}$ is expressed in terms of the incident daily cases \mathcal{I} ,

$$\mathcal{R}_{j,t} = \frac{\sum_{\tau=0}^6 \mathcal{I}_{j,t-\tau}}{\sum_{\tau=4}^{10} \mathcal{I}_{j,t-\tau}}, \quad (4)$$

assuming a generation time of $\tau_{\text{gen}} = 4$ d. Incident daily case data are not taken directly from case reports, but based on a reconstruction of symptom onset.

355 The explanatory NPI variables are constructed from a detailed, county-level data set of NPIs
356 compiled by *infas* (Institut für angewandte Sozialwissenschaft, Institute for Applied Social Sci-
357 ence), which are available online from www.healthcare-datenplattform.de/. The *infas*
358 dataset codes a number of subcategories for 23 main categories of NPIs (e.g., for contacts in
359 private settings, primary and secondary schools, masking). *StopptCOVID* uses a subset of these
360 subcategories to assign a level to each main category of NPIs. Up to 6 levels for NPI settings are
361 distinguished, with Level 2 (L2) representing the least stringent form of restrictions. The NPI level
362 is determined by the most stringent active restriction. For a detailed breakdown of restrictions at
363 each level, see Supplementary Table S1. Due to strong correlations, NPIs for the cultural sector,
364 the hotel and restaurant industry and sports (CHRS) are included in combined categories, depend-
365 ing on how many of these sectors were subject to the highest level of NPI stringency, or failing
366 that, on the second-highest level of stringency.

367 The different levels are treated as binary variables at the county level. Gaps in the NPI dataset
368 were filled by imputation (last observation carried forward). State-level NPI variables on a contin-
369 uous scale from 0 to 1 are then constructed as population-weighted averages of the county-level
370 NPI variables, and a lag by τ_{NPI} is applied before these are fed into the linear regression model.

371 For maximum consistency with *StopptCOVID*, the input data for \mathcal{R} , case numbers NPIs are
372 read out from their publicly available R scripts [12]. In line with our strict focus on the statistical
373 analysis, this eliminates the danger of divergent results due to potential misunderstandings about
374 the coding and imputation of the explanatory variables. We highlight, however, that a superficial
375 examination of the data revealed some anomalies. For example, the coded NPI variables do not
376 show any health restrictions in child care facilities in the state of Mecklenburg-Vorpommern in
377 2020. This contradicts information by the state government [45] and is evidently wrong.

378 However, as the *StopptCOVID* dataset does not include the response variable and the explana-
379 tory variables for every day of the period of interest. Data are not provided for short periods
380 without cases in individual states. During these phases, we impute data for all variables by linear
381 interpolation to permit the application of certain analysis methods for time series that cannot easily
382 deal with data gaps.

383 **Procedures for IV&V exercise**

384 To select methods for an ensemble-based IV&V exercise, we adapted the procedures outlined
385 by den Boon et al. [27]. To identify suitable methods for the reevaluation, we applied predefined
386 general selection criteria, namely,

- 387 • use in prior NPI studies and widespread use for inference and regression problems in other
388 fields,
- 389 • rigorous derivation from first principles,
- 390 • the extent by which key problems (e.g., autocorrelation and multicollinearity) were ad-
391 dressed,
- 392 • sufficient differentiation from other approaches included in the comparison (to avoid acci-
393 dentally obtaining similar results by construction),
- 394 • and for final inclusion the demonstration of superior/competitive sensitivity and precision
395 when compared to the RKI approach.

396 To survey methods commonly used in the evaluation of NPI effectiveness, we screened all studies
397 cited in the Royal Society's recent NPI review [30] examining the effectiveness of NPIs for SARS-
398 CoV-2, unless the study type was deemed not relevant in the context of time series analyses. The
399 following study types were excluded: Case report (study) or series, (prospective or retrospective)
400 cohort study, contact survey, randomised control trial. All others were included, even if the use
401 of time series was not made explicit (e.g., ecological studies). In an initial round of screening,
402 we reviewed the subset of studies that considered $\mathcal{R}(t)$ as outcome, and whose quality of evidence
403 was not rated as *very low* by [30] according to their GRADE assessment [46, 47]. Based on this
404 initial review, we defined various relevant dimensions to broadly categorise all studies according
405 to their different methodological approaches. A description of those categories and dimensions
406 can be found in Supplementary Methods S3 (taxonomy of models for NPI effect estimation), S4
407 (methods for error analysis), and S5 (methods for addressing multicollinearity).

408 The results were presented to the project working group. Based mainly on the above selection
409 criteria, a specific subset of methods representing standard approaches for the different categories
410 were chosen (Table 1). Additional reasons for inclusion or rejection are outlined and further
411 discussed in Supplementary Methods S3. The chosen methods were implemented using the dataset
412 employed by *StopptCOVID* without altering the epidemiological model assumptions. Finally, a

413 framework for determining effect sizes and uncertainties of summary measures/ensemble values
414 from all models was agreed upon and documented. For more details on the literature review and
415 decision process, see Supplementary Methods S2.

416 An open call for participation was sent to a number of scientific societies at the beginning
417 of the project: Association of the Scientific Medical Societies in Germany (AWMF), German
418 Society for Epidemiology, German Association for Medical Informatics, Biometry and Epidemi-
419 ology (GMDS), Deutsche Arbeitsgemeinschaft Statistik (DAGStat), German Statistical Society
420 (DStatG), Verein für Socialpolitik e.V., Deutsche Mathematiker-Vereinigung (DMV), Deutsche
421 Physikalische Gesellschaft, Deutsches Klima-Konsortium (DKK), German Reproducibility Net-
422 work. The German Network for Evidence-Based Medicine and the German Society for Epidemi-
423 ology kindly disseminated the call, and the German Reproducibility Network provided contact
424 details of member institutions for further distribution.

425 **Model Ensemble**

426 Among the methods selected for the model ensemble (Table 1), models DK, Ebisuzaki, and
427 BT merely use different methods for calculating confidence intervals from the residuals in the
428 case of autocorrelated errors, and are implemented on top of the baseline WLS model. **Model**
429 **DK** employs the Driscoll-Kraay estimator [73, 74], which uses an estimate of the error covariance
430 matrix up to a specified temporal lag and across entities to compute the variances of the regression
431 coefficients. **Model Ebisuzaki** adapts a frequency-domain method [93] that takes autocorrelation
432 into account by decomposing the residuals into Fourier components and computes confidence
433 intervals based on the power spectrum of the residuals. **Model BT** computes errors using a time
434 series bootstrap [88] that randomly resamples chunks of the time series of residuals such as to
435 preserve their autocorrelation structure. Such a bootstrap is also used for models 2WFE, DYN, RF,
436 Elastic Net and PCR. **Model 2WFE** uses fixed effects both for entities and time [122] to subtract
437 unmodelled temporal dynamics common to all Federal states. Effect estimates for seasonality
438 are obtained in a hierarchical approach by regressing the fixed effects in terms of the seasonal
439 variables. **Model ARMA**(p, q) models regression errors $n_{j,t}$ as an autoregressive moving-average
440 process of order (p, q),

$$n_{j,t} - \sum_{\tau=1}^p \phi_{\tau} n_{j,t-\tau} = \epsilon_{j,t} + \sum_{\tau=1}^q \theta_{\tau} \epsilon_{j,t-\tau}, \quad (5)$$

441 and estimates the regression coefficients and the autoregression coefficients ϕ_τ and θ_τ by maximum
442 likelihood estimation using a state-space formulation [34]. The optimal choice $(p, q) = (1, 11)$ is
443 obtained by minimising the Bayesian information criterion [133]. **Model DYN** combines Equa-
444 tions (2) and (4) into a renewal equation for the 7-day average case data, and is fitted directly to
445 the case data instead of $\mathcal{R}(t)$. **Model RF** is an implementation of random forest regression [137],
446 which constructs an ensemble of decision trees in the explanatory variables from random samples
447 of the data and then averages the results. The trees are fitted to minimise the squared error. The
448 number and depth of trees and the features considered for the tree splits are optimised by cross
449 validation with a time series split. **Models Elastic Net** and **PCR** use linear regression with reg-
450 ularisation as a possible remedy for multicollinearity. Elastic net regression [105] adds penalty
451 terms to the likelihood for more stable estimates of regression coefficients in exchange for some
452 bias. Model PCR uses truncated singular value decomposition [97, 98] (non-centred principal
453 component analysis) to filter out patterns in the explanatory variables that contribute to unstable
454 estimates. The regularisation parameters are again determined by cross validation.

455 Except for model RF, the fitted models immediately yield estimates of linear effects on $\ln \mathcal{R}(t)$
456 that have exactly the same interpretation as in the baseline model. For model RF, linear effect sizes
457 are extracted as the weighted average difference in $\ln \mathcal{R}(t)$ between two counterfactual scenarios
458 when an intervention is switched on or off completely, while the other NPIs have their actual acti-
459 vation patterns (case RF1). In addition, we also consider this average difference for the case when
460 only seasonal effects and holidays are switched on in the model (case RF0). Crudely speaking,
461 cases RF1 and RF0 give estimates for the effect of an NPI in conjunction with all others, or as a
462 single intervention.

463 The ensemble was implemented in PYTHON using `statsmodels` [56] and `sklearn` [138]. For
464 a detailed technical description, we refer to Supplementary Methods S6.

465 DATA AVAILABILITY

466 Our PYTHON code is freely available on GITHUB (https://github.com/bjmueller/npi_ivv_code). The code utilises processed data from the original *StopptCOVID* project [12], which
467 `vv_code`). The code utilises processed data from the original *StopptCOVID* project [12], which
468 in turn uses NPI data from *infas* (www.healthcare-datenplattform.de/), which are freely
469 available after registration. Our code repository contains instructions to download the required
470 third-party data and code to generate the requisite input data for our analysis.

-
- 471 [1] UK Covid-19 Inquiry. Module 1 report: The resilience and preparedness of the United Kingdom
472 (2024). URL [https://covid19.public-inquiry.uk/reports/module-1-report-the-re-](https://covid19.public-inquiry.uk/reports/module-1-report-the-resilience-and-preparedness-of-the-united-kingdom/)
473 [silience-and-preparedness-of-the-united-kingdom/](https://covid19.public-inquiry.uk/reports/module-1-report-the-resilience-and-preparedness-of-the-united-kingdom/).
- 474 [2] Murphy, C. *et al.* Effectiveness of social distancing measures and lockdowns for reducing transmis-
475 sion of COVID-19 in non-healthcare, community-based settings. *Philos. Trans. A Math. Phys. Eng.*
476 *Sci.* **381**, 20230132 (2023).
- 477 [3] Funk, M. *et al.* Wirksamkeit von Corona-Massnahmen in der Schweiz (2022). URL [https://www.](https://www.seco.admin.ch/seco/de/home/Publikationen_Dienstleistungen/Publikationen_und_Formulare/Strukturwandel_Wachstum/Wachstum/wirksamkeit-corona-massnahmen-schweiz.html)
478 [seco.admin.ch/seco/de/home/Publikationen_Dienstleistungen/Publikationen_und](https://www.seco.admin.ch/seco/de/home/Publikationen_Dienstleistungen/Publikationen_und_Formulare/Strukturwandel_Wachstum/Wachstum/wirksamkeit-corona-massnahmen-schweiz.html)
479 [_Formulare/Strukturwandel_Wachstum/Wachstum/wirksamkeit-corona-massnahmen-s](https://www.seco.admin.ch/seco/de/home/Publikationen_Dienstleistungen/Publikationen_und_Formulare/Strukturwandel_Wachstum/Wachstum/wirksamkeit-corona-massnahmen-schweiz.html)
480 [chweiz.html](https://www.seco.admin.ch/seco/de/home/Publikationen_Dienstleistungen/Publikationen_und_Formulare/Strukturwandel_Wachstum/Wachstum/wirksamkeit-corona-massnahmen-schweiz.html).
- 481 [4] Bundesministerium für Gesundheit. Evaluation der Rechtsgrundlagen und Maßnahmen der Pan-
482 demiepolitik. Bericht des Sachverständigenausschusses nach §5 Abs. 9 IfSG (2022). URL [https:](https://www.bundesgesundheitsministerium.de/fileadmin/Dateien/3_Downloads/S/Sachv_erstaendigenausschuss/BER_lfSG-BMG.pdf)
483 [//www.bundesgesundheitsministerium.de/fileadmin/Dateien/3_Downloads/S/Sachv](https://www.bundesgesundheitsministerium.de/fileadmin/Dateien/3_Downloads/S/Sachv_erstaendigenausschuss/BER_lfSG-BMG.pdf)
484 [erstaendigenausschuss/BER_lfSG-BMG.pdf](https://www.bundesgesundheitsministerium.de/fileadmin/Dateien/3_Downloads/S/Sachv_erstaendigenausschuss/BER_lfSG-BMG.pdf).
- 485 [5] Iezadi, S. *et al.* Effectiveness of non-pharmaceutical public health interventions against COVID-19:
486 A systematic review and meta-analysis. *PLoS One* **16**, e0260371 (2021).
- 487 [6] Raftery, A. E., Madigan, D. & Hoeting, J. A. Bayesian Model Averaging for Linear Regression
488 Models. *Journal of the American Statistical Association* **92**, 179–191 (1997). URL [https://doi.](https://doi.org/10.1080/01621459.1997.10473615)
489 [org/10.1080/01621459.1997.10473615](https://doi.org/10.1080/01621459.1997.10473615).
- 490 [7] Edeling, W. *et al.* The impact of uncertainty on predictions of the CovidSim epidemiological code.
491 *Nat. Comput. Sci.* **1**, 128–135 (2021).
- 492 [8] Chin, V., Ioannidis, J. P. A., Tanner, M. A. & Cripps, S. Effect estimates of COVID-19 non-
493 pharmaceutical interventions are non-robust and highly model-dependent. *J. Clin. Epidemiol.* **136**,
494 96–132 (2021).
- 495 [9] an der Heiden, M., Hicketier, A. & Bremer, V. Wirksamkeit und Wirkung von anti-epidemischen
496 Maßnahmen auf die COVID-19-Pandemie in Deutschland (StopptCOVID-Studie) (2023). URL [ht](https://www.rki.de/DE/Content/InfAZ/N/Neuartiges_Coronavirus/Projekte_RKI/StopptCOVID_studie.html)
497 [tps://www.rki.de/DE/Content/InfAZ/N/Neuartiges_Coronavirus/Projekte_RKI/Sto](https://www.rki.de/DE/Content/InfAZ/N/Neuartiges_Coronavirus/Projekte_RKI/StopptCOVID_studie.html)
498 [pptCOVID_studie.html](https://www.rki.de/DE/Content/InfAZ/N/Neuartiges_Coronavirus/Projekte_RKI/StopptCOVID_studie.html).

- 499 [10] Meyer, G., Mühlhauser, I., Brinks, R. & Müller, B. Wirksame Kontrollmaßnahmen in der SARS-
500 CoV-2-Pandemie? *KVH Journal* **10/2023** (2023).
- 501 [11] Baumgarten, W., Beige, O., Haake, D., Merkl, J. & Wieland, T. Was die „StopptCOVID“-Studie des
502 RKI sagt - und was nicht. *Cicero Online* (2023). URL [https://www.cicero.de/innenpolitik](https://www.cicero.de/innenpolitik/corona-pandemie-robertkochinstitut-studie)
503 [/corona-pandemie-robertkochinstitut-studie](https://www.cicero.de/innenpolitik/corona-pandemie-robertkochinstitut-studie).
- 504 [12] Hicketier, A. & an der Heiden, M. StopptCOVID-Studie - Daten, Analyse und Ergebnisse (2024).
505 URL <https://zenodo.org/records/10888033>.
- 506 [13] Bodderas, E. Kanzleramt drängt Lauterbach zur Offenlegung von Corona-Studie. *WELT*, 8.2.2024
507 (2024). URL [https://www.welt.de/politik/deutschland/plus250799234/Pandemie-A](https://www.welt.de/politik/deutschland/plus250799234/Pandemie-Aufarbeitung-Kanzleramt-draengt-Lauterbach-zur-Offenlegung-von-Corona-Studie.html)
508 [ufarbeitung-Kanzleramt-draengt-Lauterbach-zur-Offenlegung-von-Corona-Studi](https://www.welt.de/politik/deutschland/plus250799234/Pandemie-Aufarbeitung-Kanzleramt-draengt-Lauterbach-zur-Offenlegung-von-Corona-Studie.html)
509 [e.html](https://www.welt.de/politik/deutschland/plus250799234/Pandemie-Aufarbeitung-Kanzleramt-draengt-Lauterbach-zur-Offenlegung-von-Corona-Studie.html).
- 510 [14] Stevens, G. A. *et al.* Guidelines for accurate and transparent health estimates reporting: The
511 GATHER statement. *Lancet* **388**, e19–e23 (2016).
- 512 [15] Durbin, J. & Watson, G. S. Testing for serial correlation in least squares regression. i. *Biometrika*
513 **37**, 409–428 (1950). URL <https://doi.org/10.1093/biomet/37.3-4.409>.
- 514 [16] Greene, W. H. *Econometric Analysis* (Pearson Education, 2003), fifth edn. URL [http://pages.](http://pages.stern.nyu.edu/~wgreene/Text/econometricanalysis.htm)
515 [stern.nyu.edu/~wgreene/Text/econometricanalysis.htm](http://pages.stern.nyu.edu/~wgreene/Text/econometricanalysis.htm).
- 516 [17] Johnston, R., Jones, K. & Manley, D. Confounding and collinearity in regression analysis: a cau-
517 tionary tale and an alternative procedure, illustrated by studies of British voting behaviour. *Quality*
518 *& Quantity* **52**, 1957–1976 (2018). URL <https://doi.org/10.1007/s11135-017-0584-6>.
- 519 [18] Driscoll, J. C. & Kraay, A. C. Consistent Covariance Matrix Estimation With Spatially Dependent
520 Panel Data. *The Review of Economics and Statistics* **80**, 549–560 (1998). URL [https://ideas.re](https://ideas.repec.org/a/tpr/restat/v80y1998i4p549-560.html)
521 [pec.org/a/tpr/restat/v80y1998i4p549-560.html](https://ideas.repec.org/a/tpr/restat/v80y1998i4p549-560.html).
- 522 [19] Ebisuzaki, W. A Method to Estimate the Statistical Significance of a Correlation When the Data Are
523 Serially Correlated. *Journal of Climate* **10**, 2147–2153 (1997).
- 524 [20] Politis, D. N. & Romano, J. P. The stationary bootstrap. *Journal of the American Statistical Associ-*
525 *ation* **89**, 1303–1313 (1994). URL <https://doi.org/10.1080/01621459.1994.10476870>.
- 526 [21] Abbott, B. *et al.* Analysis of first LIGO science data for stochastic gravitational waves. *Phys. Rev. D*
527 **69**, 122004 (2004). [gr-qc/0312088](https://arxiv.org/abs/gr-qc/0312088).
- 528 [22] Rizzuto, A. C. *et al.* Zodiacal Exoplanets in Time (ZEIT). VIII. A Two-planet System in Praesepe
529 from K2 Campaign 16. *Astron. J.* **156**, 195 (2018). 1808.07068.

- 530 [23] Polack, F. P. *et al.* Safety and efficacy of the BNT162b2 mRNA covid-19 vaccine. *N. Engl. J. Med.*
531 **383**, 2603–2615 (2020).
- 532 [24] Baden, L. R. *et al.* Efficacy and safety of the mRNA-1273 SARS-CoV-2 vaccine. *N. Engl. J. Med.*
533 **384**, 403–416 (2021).
- 534 [25] Chemaitelly, H. *et al.* mRNA-1273 COVID-19 vaccine effectiveness against the B.1.1.7 and B.1.351
535 variants and severe COVID-19 disease in Qatar. *Nat. Med.* **27**, 1614–1621 (2021).
- 536 [26] Shah, A. S. *et al.* Effect of Vaccination on Transmission of SARS-CoV-2. *New England Journal of*
537 *Medicine* **385**, 1718–1720 (2021). URL [https://www.nejm.org/doi/full/10.1056/NEJMc2](https://www.nejm.org/doi/full/10.1056/NEJMc2106757)
538 [106757](https://www.nejm.org/doi/pdf/10.1056/NEJMc2106757). <https://www.nejm.org/doi/pdf/10.1056/NEJMc2106757>.
- 539 [27] den Boon, S. *et al.* Guidelines for multi-model comparisons of the impact of infectious disease
540 interventions. *BMC Medicine* **17**, 163 (2019).
- 541 [28] Spiliopoulos, L. On the effectiveness of COVID-19 restrictions and lockdowns: Pan metron ariston.
542 *BMC Public Health* **22**, 1842 (2022).
- 543 [29] Perra, N., Balcan, D., Gonçalves, B. & Vespignani, A. Towards a characterization of behavior-disease
544 models. *PLoS One* **6**, e23084 (2011).
- 545 [30] Agaba, G., Kyrychko, Y. & Blyuss, K. Mathematical model for the impact of awareness on the
546 dynamics of infectious diseases. *Mathematical Biosciences* **286**, 22–30 (2017). URL [https://ww](https://www.sciencedirect.com/science/article/pii/S0025556417300433)
547 [w.sciencedirect.com/science/article/pii/S0025556417300433](https://www.sciencedirect.com/science/article/pii/S0025556417300433).
- 548 [31] Neipel, J., Bauermann, J., Bo, S., Harmon, T. & Jülicher, F. Power-law population heterogeneity
549 governs epidemic waves. *PLoS ONE* **15**, e0239678 (2020). 2008.00471.
- 550 [32] Gomes, M. G. M. *et al.* Individual variation in susceptibility or exposure to SARS-CoV-2 lowers
551 the herd immunity threshold. *Journal of Theoretical Biology* **540**, 111063 (2022). URL [https:](https://www.sciencedirect.com/science/article/pii/S0022519322000613)
552 [//www.sciencedirect.com/science/article/pii/S0022519322000613](https://www.sciencedirect.com/science/article/pii/S0022519322000613).
- 553 [33] Munro, A. P. S. & House, T. Cycles of susceptibility: Immunity debt explains altered infectious
554 disease dynamics post-pandemic. *Clinical Infectious Diseases* ciae493 (2024). URL [https://do](https://doi.org/10.1093/cid/ciae493)
555 [i.org/10.1093/cid/ciae493](https://doi.org/10.1093/cid/ciae493).
- 556 [34] Feng, Z., Towers, S. & Yang, Y. Modeling the effects of vaccination and treatment on pandemic
557 influenza. *AAPS J.* **13**, 427–437 (2011).
- 558 [35] Worby, C. J., Wallinga, J., Lipsitch, M. & Goldstein, E. Population effect of influenza vaccination
559 under co-circulation of non-vaccine variants and the case for a bivalent a/h3n2 vaccine component.
560 *Epidemics* **19**, 74–82 (2017). URL <https://www.sciencedirect.com/science/article/pii>

561 i/S1755436517300208.

- 562 [36] Backer, J., van Boven, M., van der Hoek, W. & Wallinga, J. Vaccinating children against influenza
563 increases variability in epidemic size. *Epidemics* **26**, 95–103 (2019). URL [https://doi.org/10](https://doi.org/10.1016/j.epidem.2018.10.003)
564 [.1016/j.epidem.2018.10.003](https://doi.org/10.1016/j.epidem.2018.10.003).
- 565 [37] de Boer, P. T., Backer, J. A., van Hoek, A. J. & Wallinga, J. Vaccinating children against influenza:
566 overall cost-effective with potential for undesirable outcomes. *BMC Med.* **18**, 11 (2020).
- 567 [38] Dürr, H.-P. & Eichner, M. Corona-Pandemie: Zukunfts-Überlegungen aus der Sicht epidemiologis-
568 cher Modellierung. *Monitor Versorgungsforschung* **02/22**, 57 (2022).
- 569 [39] Castioni, P., Gómez, S., Granell, C. & Arenas, A. Rebound in epidemic control: how misaligned
570 vaccination timing amplifies infection peaks. *npj Complexity* **1**, 20 (2024).
- 571 [40] Anderson, R. M. & May, R. M. Vaccination against rubella and measles: quantitative investigations
572 of different policies. *J. Hyg. (Lond.)* **90**, 259–325 (1983).
- 573 [41] Cohen, T. & Lipsitch, M. Too little of a good thing: a paradox of moderate infection control.
574 *Epidemiology* **19**, 588–589 (2008).
- 575 [42] Heffernan, J. M. & Keeling, M. J. Implications of vaccination and waning immunity. *Proc. Biol. Sci.*
576 **276**, 2071–2080 (2009).
- 577 [43] Chiolero, A., Tancredi, S. & Ioannidis, J. P. A. Slow data public health. *Eur. J. Epidemiol.* **38**,
578 1219–1225 (2023).
- 579 [44] Akaike, H. *Information Theory and an Extension of the Maximum Likelihood Principle*, 199–213
580 (Springer New York, New York, NY, 1998). URL [https://doi.org/10.1007/978-1-4612-1](https://doi.org/10.1007/978-1-4612-1694-0_15)
581 [694-0_15](https://doi.org/10.1007/978-1-4612-1694-0_15).
- 582 [45] State Government of Mecklenburg-Vorpommern. Kitas und Kindertagespflege ab Montag flächen-
583 deckend geschlossen (2020). URL [https://www.regierung-mv.de/Landesregierung/sm/A](https://www.regierung-mv.de/Landesregierung/sm/Aktuell/?id=158498&processor=processor.sa.pressemitteilung)
584 [ktuell/?id=158498&processor=processor.sa.pressemitteilung](https://www.regierung-mv.de/Landesregierung/sm/Aktuell/?id=158498&processor=processor.sa.pressemitteilung).
- 585 [46] Guyatt, G. *et al.* GRADE guidelines: 1. Introduction-GRADE evidence profiles and summary of
586 findings tables. *J. Clin. Epidemiol.* **64**, 383–394 (2011).
- 587 [47] Balshem, H. *et al.* GRADE guidelines: 3. rating the quality of evidence. *J. Clin. Epidemiol.* **64**,
588 401–406 (2011).
- 589 [48] Hoechle, D. Robust standard errors for panel regressions with cross-sectional dependence. *Stata*
590 *Journal* **7**, 281–312 (2007). URL [https://ideas.repec.org/a/tsj/stataj/v7y2007i3p281](https://ideas.repec.org/a/tsj/stataj/v7y2007i3p281-312.html)
591 [-312.html](https://ideas.repec.org/a/tsj/stataj/v7y2007i3p281-312.html).

- 592 [49] Wooldridge, J. M. Two-way fixed effects, the two-way mundlak regression, and difference-in-
593 differences estimators (2021). URL <https://ssrn.com/abstract=3906345>.
- 594 [50] Harvey, A. *State space models*, 269–275 (Palgrave Macmillan UK, London, 2010). URL https://doi.org/10.1057/9780230280830_30.
- 595
- 596 [51] Schwarz, G. Estimating the Dimension of a Model. *The Annals of Statistics* **6**, 461–464 (1978). URL
597 <http://www.jstor.org/stable/2958889>.
- 598 [52] Breiman, L. Random forests. *Machine Learning* **45**, 5–32 (2001).
- 599 [53] Zou, H. & Hastie, T. Regularization and variable selection via the elastic net. *Journal of the Royal*
600 *Statistical Society Series B: Statistical Methodology* **67**, 301–320 (2005). URL [https://doi.org/](https://doi.org/10.1111/j.1467-9868.2005.00503.x)
601 [10.1111/j.1467-9868.2005.00503.x](https://doi.org/10.1111/j.1467-9868.2005.00503.x).
- 602 [54] Hansen, P. C. Truncated Singular Value Decomposition Solutions to Discrete Ill-Posed Problems
603 with Ill-Determined Numerical Rank. *SIAM Journal on Scientific and Statistical Computing* **11**,
604 503–518 (1990). URL <https://doi.org/10.1137/0911028>. <https://doi.org/10.1137/0911028>.
- 605 [55] Sekii, T. Two-Dimensional Inversion for Solar Internal Rotation. *PASJ* **43**, 381–411 (1991).
- 606 [56] Seabold, S. & Perktold, J. statsmodels: Econometric and statistical modeling with python. In *9th*
607 *Python in Science Conference* (2010).
- 608 [57] Pedregosa, F. *et al.* Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*
609 **12**, 2825–2830 (2011).

610 ACKNOWLEDGEMENTS

611 We thank M. an der Heiden and V. Bremer for answering technical questions on the origi-
612 nal *StopptCOVID* project. We acknowledge helpful discussions with W. Baumgarten, O. Beige,
613 G. Meyer, I. Mühlhauser, D. Schuricht, and T. Wieland. We are grateful to the German Network for
614 Evidence-Based Medicine, the German Society for Epidemiology and the German Reproducibility
615 Network for help in distributing the project call.

616 AUTHOR CONTRIBUTIONS

617 B.M.: Conceptualisation, Project Management, Software, Formal analysis, Visualisation, Writ-
618 ing – Original draft, Review and Editing. I.P.: Conceptualisation, Procedures, Literature Review,
619 Writing – Original draft, Review and Editing. ML: Conceptualisation, Code checks (models DK

620 and Ebisuzaki), Writing – Review. R.B.: Conceptualisation, Writing – Review and Editing. S.C.:
621 Conceptualisation, Writing – Review and Editing. M.G.M.M.: Conceptualisation, Writing – Re-
622 view. D.H.: Conceptualisation, Writing – Review and Editing. J.P.A.I.: Conceptualisation, Writ-
623 ing – Review and Editing.

624 **FUNDING**

625 There is no funding to report.

626 **ETHICS DECLARATION**

627 This study did not involve research on humans or animals, and only used publicly available,
628 non-personal data sets.

629 **COMPETING INTERESTS**

630 B.M. I.P. M.L., and R.B. are signatories of a call for a non-partisan pandemic review in Ger-
631 many (<https://pandemieaufarbeitung.net>). B.M. has been engaged in discussion and con-
632 sultation of pandemic policy and science policy with members of several German parties (CDU,
633 CSU, FDP, SPD, BSW, Greens), but is not receiving remuneration. S.C., M.G.M.G., D.H., and
634 J.P.A.I. have no competing interests to declare.