

1 **Title:** From Sequences to Strategies: Early Detection of New SARS-CoV-2 Variants via Genetic
2 Distance to Reduce Hospitalizations

3
4 **Authors:** Marika D'Avanzo^{1,2,*}, Aung Pone Myint^{3,*}, Giacomo Cacciapaglia^{4,5}, Stefan
5 Hohenegger⁶, Francesco Conventi^{2,7}, Marta Nunes³

6
7 **Affiliations**

8 1. PhD National Programme in One Health approaches to infectious diseases and life science
9 research, Department of Public Health, Experimental and Forensic Medicine, University of Pavia,
10 Pavia, 27100, Italy.

11 2. INFN Sezione di Napoli, Complesso Universitario di Monte S. Angelo Edificio 6, Via Cintia,
12 80126 Naples, Italy.

13 3. Center of Excellence in Respiratory Pathogens (CERP), Hospices Civils de Lyon (HCL) and
14 Centre International de Recherche en Infectiologie (CIRI), Équipe Santé Publique, Épidémiologie
15 et Écologie Évolutive des Maladies Infectieuses (PHE3ID), Inserm U1111, CNRS UMR5308,
16 ENS de Lyon, Université Claude Bernard Lyon 1, Lyon, France.

17 4. Laboratoire de Physique Théorique et Hautes Energies (LPTHE), UMR 7589, Sorbonne
18 Université & CNRS, 4 place Jussieu, 75252 Paris Cedex 05, France.

19 5. Quantum Theory Center (QTC) at IMADA & D-IAS, Southern Denmark Univ., Campusvej
20 55, 5230 Odense M, Denmark.

21 6. Université Claude Bernard Lyon 1, CNRS/IN2P3, IP2I UMR 5822, 4 rue Enrico Fermi, F-
22 69100 Villeurbanne, France.

23 7. Dipartimento di Ingegneria, Università degli studi di Napoli Parthenope, Centro Direzionale di
24 Napoli, Isola C 4, lato Sud, 80143 Naples, Italy.

25 * These authors contributed equally to this article and share first authorship.

26
27
28 **Keywords:** SARS-CoV-2 variants, Genomic surveillance, Genetic distance, Machine learning,
29 Containment strategies, Early detection, Hospitalization rates

30

31 **Abstract:** The COVID-19 pandemic highlighted the critical need for robust methods to monitor
32 viral evolution and detect emerging variants of concern (VOCs). Traditional genomic
33 surveillance often lacks predictive power. This study expanded an unsupervised machine learning
34 clustering algorithm, based on SARS-CoV-2 Spike protein Levenshtein distance, to track and
35 predict variant predominance across six European countries from 2020 to January 2024. We also
36 investigated the influence of genetic distances and containment strategies on hospitalization rates.
37 Sequences were transformed into temporal chains, and growth parameters were extracted via
38 sigmoid fitting. A deep neural network (DNN) was trained to classify emerging chains as likely
39 dominant, while a CatBoost model assessed variable importance for predicting weekly
40 hospitalizations in Denmark. Simulations explored modifying vaccine genetic distance,
41 containment measures, and VCR.
42 Approximately 5,000 sequences per week enabled early chain detection within four weeks. The
43 DNN achieved near-perfect classification of chain predominance within 3-4 weeks of
44 appearance. Genetic distances within consecutive chains and with vaccine strains were significant
45 predictors of hospitalizations. Simulations suggest that better-matched vaccines or stricter
46 containment measures could reduce hospitalizations. Doubling vaccination coverage alone had
47 minimal effect but showed additional reductions when combined with strict containment.
48 This integrated framework demonstrates the utility of combining unsupervised and supervised
49 machine learning for real-time tracking and prediction of SARS-CoV-2 variant dynamics and
50 their impact on public health. Our findings underscore the critical role of genetic distances and
51 effective public health interventions in mitigating the burden of emerging variants, supporting
52 timely genomic surveillance and adaptive public health strategies.
53

54 From Sequences to Strategies: Early Detection of New SARS- 55 CoV-2 Variants via Genetic Distance to Reduce Hospitalizations

56 Introduction

57 Severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2) has evolved continuously
58 producing variants with different levels of transmissibility, immune escape, and disease severity.
59 Its spike protein is central to host cell entry and continuous accumulation of mutations in the
60 spike protein have led to the emergence of phylogenetically distinct variants[1] which can evade
61 immune responses from prior infections or vaccination[2]. Four of these variants, Alpha, Beta,
62 Delta, and Omicron, have been classified as variants of concern (VOCs) by the World Health
63 Organization (WHO), while others have been designated as variants of interest (VOIs) or variants
64 under monitoring (VUMs)[3]. Since December 2021, Omicron and its sub-lineages have globally
65 dominated the epidemiological dynamics. On May 5, 2023, the WHO declared the end of the
66 COVID-19 public health emergency of international concern while acknowledging the ongoing
67 risks posed by future SARS-CoV-2 evolution[4]. This highlights the importance of continuous
68 viral genomic surveillance to promptly detect and assess new variants that could impact public
69 health, particularly in the context of diverse immunity profiles resulting from varied exposure due
70 to irregular vaccines coverage and prior infections.

71 To address this need, de Hoffer et al.[5] developed an unsupervised machine learning (ML)
72 algorithm to define new variants by clustering the amino acid sequences of the spike protein of
73 SARS-CoV-2 based on the Levenshtein distance, and a time-binned hierarchical clustering with
74 Ward's method. Clusters across consecutive time bins are then linked as chains when they
75 contain the same dominant spike sequence. These chains were empirically found to correspond to
76 persistent and potentially epidemiologically significant variants. The approach effectively
77 predicted the Alpha and Delta variants in the United Kingdom.

78 Levi et al.[6] analyzed data from 30 countries to identify variants associated with over 1,000
79 cases per million population within a three-month period using Jaccard distance. However, the
80 method's reliance on retrospective metrics, the maximum weekly case count observed over the
81 full variant duration, limits its applicability for real-time prediction. Nicora et al.[7] used k-mer

82 counts and a one-class Support Vector Machine (SVM) to flag anomalous sequences, though it
83 suffered from a high false-positive rate. Rancati et al.[8] also employed spike k-mer
84 representations with an autoencoder to predict lineages that would exceed 10% of total
85 sequences, but its performance was inconsistent in countries with relatively lower sequencing
86 volumes, such as France and Denmark. More recently, Feng et al.[9] developed a transformer-
87 based model designed to forecast future lineage frequencies up to two months in advance,
88 however, it relies on prior lineage designation and focuses exclusively on frequency forecasting.

89 While COVID-19 vaccines have been shown to protect against infections and
90 hospitalizations[10], the Omicron variant with extensive spike mutations has caused substantial
91 hospitalizations[11] despite lower intrinsic virulence[12]. Current approaches for assessing
92 vaccine effectiveness against novel variants (primarily in vitro neutralization assays[13,14]) have
93 limitations in their timeliness and predictive capacity. Studies such as Cao et al.[15] have
94 explored the relationship between genetic divergence and vaccine effectiveness, suggesting the
95 genetic distance between circulating and vaccine strains as an indicator for immune escape
96 potential. However, a direct, quantitative link between genetic divergence and real-world clinical
97 outcomes has not been established.

98 In this study, we extend the unsupervised ML algorithm by de Hoffer et al.[5] to analyze spike
99 protein sequences from six European countries. We assess the robustness of the algorithm across
100 geographically and temporally diverse datasets, identify the optimal sequence volume required
101 for early variant detection and evaluate predictive parameters for variant predominance.
102 Furthermore, we incorporate a deep learning classifier to improve early predictions of variant
103 dominance based on early prevalence trends. Additionally, we investigate impact of genetic
104 distance between variants and distances with vaccine strains on observed hospitalizations by
105 leveraging comprehensive Danish health data, together with different containment strategies. By
106 integrating the identification of emerging variants through genetic surveillance with a novel
107 analysis of the direct link between spike protein evolution and hospitalization, this study offers a
108 valuable framework for predicting and responding to future epidemics caused by novel viruses.

109 Methods

110 Data sources

111 We analyzed SARS-CoV-2 Spike protein sequences from Germany, Italy, Sweden, Denmark,
112 France, and Spain sourced from the GISAID database[16]. The Oxford Covid-19 Government
113 Response Tracker provided data on Containment and Health Index (CHI) through the end of
114 2022[17]. Weekly hospitalization numbers, reinfection percentages and weekly vaccination
115 coverage (VCR) in Denmark were obtained from the Denmark Statens Serum Institut
116 (SSI)[18] and complemented by the European Centre for Disease Prevention and Control
117 (ECDC)[19]. Variant-specific hospitalization risks were obtained from a retrospective study in
118 Washington State covering December 2020 to January 2022[20].

119 Cross-Country Variant Surveillance

120 Our analysis began with the algorithm developed by A. de Hoffer et al.[5], to identify viral chains
121 with the data from spike protein sequences collected between January 23, 2020 and January 15,
122 2024. Chains that were observed at least four weeks and had a peak prevalence exceeding 20%
123 were selected for analysis. Chains with the same or similar dominant sequence that represented
124 the same variant were combined. These chains were associated with the respective variants using
125 spike protein sequences available from NCBI (National Center for Biotechnology Information,
126 National Library of Medicine) and the CoVsurver app on GISAID.

127 We then fitted the frequencies of selected chains (P_c) using a mathematical function of the time x
128 to model their growth, with parameters L , b , and a , and decline, with L_2 , b_2 , and a_2 .

$$P_c(x, a, b, L, a_2, b_2, L_2) = L \left(\frac{1}{1 + e^{-\frac{(x-b)}{a}}} \right) - L_2 \left(\frac{1}{1 + e^{-\frac{(x-b_2)}{a_2}}} \right)$$

129 Key parameters derived from this fitting included (a) the growth rate, (b) the inflection point
130 (point of change in the trend), and (L) the upper asymptote (maximum expected prevalence).

131 When fitting only a limited number of weeks k , we labelled the parameters with a corresponding
132 subscript (i.e. a_k)

133 We also derived a parameter t_0 that measures the time taken to isolate the chain ($t_0 = b -$
134 $6.91 \times a$), and studied its dependence upon the average number of sequences available per week.
135 A schematic representation of the main parameters used in this work has been reported in
136 Supplementary Figure S1).

137 Deep Learning-Based Early detection of predominant variant

138 The dataset was constructed by identifying viral lineages that were observed for at least four
139 consecutive weeks of sequence detection and could be used to derive the parameters needed for
140 the model. Real-world chains from six countries were labeled as the predominant chains if they
141 eventually reached $\geq 50\%$ prevalence, and the rest as the “transient chains”.

142 Since real-world data for early-stage variant used to derive the parameters is limited, we applied a
143 data augmentation strategy to generate synthetic chains. For predominant chains, we fitted the
144 key parameters with different distributions from real chains (Supplementary Figure S2), and
145 generated new chains by sampling from these fitted distributions, adding a small noise using an
146 exponential error model. The prevalence at week 1 was computed and the growth curve was
147 recalculated iteratively.

148 For transient chains, we classified them into two transient chain groups (Supplementary Figure
149 S3): above-threshold (showing some growth tendency) and below-threshold (typically vanishing
150 early). Weekly prevalence differences, expressed as the percentage of chains ($W_{i+1} - W_i$), were
151 computed for both groups and modeled using Gaussian distributions. Chains were generated
152 using an iterative procedure sampling weekly difference, ensuring non-negativity and mimicking
153 early extinction events when the values dropped to zero.

154 We next applied a deep neural network (DNN) model consisting of an input layer with 12
155 epidemiological features: prevalence at week 1,2, and 3, early-stage fit parameters
156 $(a_3, b_3, L_3, a_4, b_4, L_4)$, and their first-order derivatives (a_4', b_4', L_4') . It was followed by two
157 hidden layers with 64 and 32 neurons respectively, both Rectified Linear Unit (ReLU) activation.
158 The output layer utilized a sigmoid function for binary classification. Model training was carried
159 out using the Adam optimization algorithm and binary cross-entropy as the loss function. The
160 model was trained for 50 epochs using a batch size of 128. The training dataset was randomly
161 split into 80% training and 20% validation.

162 To evaluate the impact of the initial prevalence of the chains on the classification performance,
163 we generated three different datasets: case 1 - all chains, case 2 – only predominant chains with
164 15% of prevalence at week 1, and case 3 – only predominant chains with 2% prevalence at week
165 1. Each dataset was generated to get 10,000 predominant chains and 10,000 transient chains
166 chosen to allow a reasonable training of the models. DNN models were constructed for each
167 dataset using different feature scenarios like using all twelve input features, prevalence values +
168 a_3, b_3, L_3 , or prevalence values only.

169 Model performance was evaluated using the Receiver Operating Characteristic (ROC) curve, the
170 False Positive Rate (FPR) at high True Positive Rate (TPR) thresholds, the confusion matrices,
171 and the signal-to-noise score distribution patterns.

172 Analysis of hospitalizations by variants in Denmark

173 To investigate the relationship between SARS-CoV-2 genetic variation and hospitalization rates,
174 we focused on Denmark due to the availability of comprehensive public data. We selected the
175 chains with >3 week prevalence and $\geq 50\%$ of prevalence at least at one time point to reduce noise
176 from transient and low-prevalent chains. Chains with the same dominating sequences were
177 combined, even when gaps in their detection weeks led the model to classify them as separate.
178 We calculated the Levenshtein distances between dominating sequences of the consecutive
179 chains (LD1). We computed the distance between dominant sequence and the spike sequence
180 used in the corresponding vaccine being employed at the time (LD_vac). We used the Wuhan
181 strain (EPI_ISL_402124) sequence for the 2020 to 2021-22 season, the BA.5 strain
182 (EPI_ISL_14026118) for 2022-2023 season as both BA.1 and BA.4/5 bivalent vaccines were
183 introduced in that season[21,22], and the XBB.1.5 strain (EPI_ISL_16134259) for the 2023/24
184 season[23]. To estimate hospitalizations attributable to each variant chain, we assumed that the
185 proportion of sequences belonging to a given chain in any week reflected its proportion of
186 hospitalizations, which were then log-transformed. Other variables included in the analysis were
187 2-week lagged VCR if primary or booster doses were reported, 2-week lagged CHI which were
188 projected to be 21.43 beyond 2022 as this magnitude had been sustained since April 2022,
189 interaction of the lagged CHI with the detection week (CH_week), and hospitalization risk due to
190 variants in which sub-lineages of Omicron were considered as similar[24–27].

191 We used a CatBoost regression model to predict weekly hospitalizations per chain. The dataset
192 was randomly split into 80% for training and 20% for testing. We used 5-fold cross-validation
193 and fine-tuned the model hyperparameters with 2000 iterations. Finally, the model's predictive
194 performance was assessed by comparing predicted hospitalizations to observed values, examining
195 residual plots, and calculating R^2 and mean absolute error (MAE). Feature importance scores,
196 SHAP (SHapley Additive exPlanations) plot, and partial dependence plots were generated to
197 understand each variable's contribution.

198 Finally, we used the trained model to explore the impact of different scenarios: employing better-
199 matched vaccine (distance at 5 for those with ≥ 15), implementing sustained strong containment
200 measures (maximum observed CHI at 68.57) or temporarily for 8 weeks, doubling VCR, and
201 combinations of high CHI and doubling VCR. All scenarios except the first one were simulated
202 to take actions at week 4 and impact observed at week 6 due to 2-week lag.

203 We assessed the simulated impact of each scenario by comparing the predicted outcomes with
204 observed hospitalizations. We used 1000 bootstrap simulations to generate 95% confidence
205 intervals (CIs). The results were reported as weekly, cumulative, and percentage reductions.

206 Results

207 Cross-Country Variant Surveillance

208 A total of 2862331 sequences for the six countries were analyzed, resulting into 109 chains for
209 Germany, 142 chains for Italy, 95 chains for Sweden, 82 chains for Denmark, 89 chains for
210 France and 98 chains for Spain. A total of 10 chains from Germany, 11 chains from Italy, 8
211 chains from Sweden, 8 chains from Denmark, 7 chains from France and 12 chains from Spain
212 were selected and successfully matched with the circulating variants.

213 Comparison of the fit results for the prevalence time evolution for the six considered European
214 countries is reported in Figure 1. The direct comparison of the six main variants across various
215 countries highlighted that the time evolution dynamics of each variant was largely independent of
216 the geographic location (Supplementary Figure S4).

217 The growth rate parameters obtained from the fit at weeks 3 and 4 (a_3 and a_4) differed between
218 dominant and transient chains, and may serve as early indicators of a chain's long-term behavior.
219 In dominant chains such as France chain 221, the parameters a , b , and L showed stability after a
220 few weeks in comparison with non-dominant chains like Denmark chain 169 with ongoing
221 fluctuations (Supplementary Figure S5).

222 The fitted curve provided the magnitude of sequencing effort required to minimize t_0 and
223 enhance early variant detection (Figure 2), and approximately 5000 sequences per week were
224 needed for the benchmark, $t_0 = 4$.

225 Deep Learning-Based Early detection of predominant variant

226 The initial dataset for the DNN model was composed of 37 predominant chains and 78 transient
227 chains. The models trained on the final dataset including simulated chains showed a very low
228 FPR value even at a very high value (99%) of TPR for all the three scenarios using the full set of
229 parameters (Supplementary Table S1). A simplified network obtained adding only growth-related
230 parameters (a_3 , b_3 , L_3) also significantly improve the prevalence-only model accuracy.

231 ROC curves confirmed strong overall classification performance, with growth-related parameters
232 offering clear advantages in sensitivity, especially in the more restrictive Case 3 dataset
233 (Supplementary Figure S6). However, relying solely on prevalence values implies a performance
234 drop.

235 The DNN score distribution further highlighted the effective discriminative power of the deep
236 learning model across all datasets except in the model with prevalence data only (Figure 3). A
237 summary of the classification performance of the deep learning model across all scenarios,
238 considering both different feature subsets and increasing levels of data restriction, can be found in
239 the supplementary Figure S7.

240 Analysis of hospitalizations by variants in Denmark

241 14 chains from Denmark were included, and 56,345 (80% of total) hospitalizations were
242 observed by these chains. Weekly hospitalizations by each chain were shown in Figure 4. High
243 distances were observed during the emergence of major variants such as Alpha (chain 53), Delta
244 (105), and the Omicron sublineages BA.2 (151), XBB.1.5 (303), and JN.1 (418) (Figure 4). The

245 reinfection rate, CHI, VCR at each week, and distribution of CH_week can be found in
246 Supplementary Figures S8 and S9.

247 The CatBoost model demonstrated high predictive accuracy (Supplementary Figure S10),
248 achieving an overall R^2 score of 0.98 and 0.85 on the unseen test dataset. The MAE was 0.10 for
249 the training set and 0.31 for the test set. Despite the higher error on the test data, the magnitude
250 remained small when compared with the outcome with a minimum of 0.52 and an average of
251 4.90. Residual diagnostics showed no evidence of heteroscedasticity or non-linearity
252 (Supplementary Figure S10).

253 The detection week of each chain was the most impactful variable, contributing 26.5% to the
254 model's explanatory power, followed by LD_vac (15.5%), VCR (14.6%), and LD1 (10.8%)
255 (Figure 5). The SHAP plot (Supplementary Figure S11) indicated that the week counter showed
256 the highest impact on the outcome with top position, and exhibiting bell-shaped association in the
257 partial dependence plot (Supplementary Figure S12). High LD1 and LD_vac values indicated a
258 positive relationship with hospitalizations and high impact on the outcome, while negative
259 association was observed for CHI_lagged. Finally, hospitalizations also showed seasonal trends,
260 increasing during December and January and decreasing during July and August (Supplementary
261 Figure S13).

262 While simulating with better-matched vaccines, significant reductions in hospitalizations (26-
263 55%) were observed for chains 151, 216, and 418 while nonsignificant reductions observed for
264 chain 303 and no change for chain 378. Temporary or sustained CHI had minimal impact on
265 early waves chain until 123 observed before 2022, since CHI levels were already high during
266 these periods (average CHI >50). Doubling VCR alone produced no statistically significant
267 changes in hospitalization trends across any of the analyzed chains. Additional reductions in
268 hospitalizations were observed in high CHI simulations when combined with doubling VCR.

269 Details of the simulated impact on the reduction in the percentages from total hospitalizations are
270 presented in Figure 6, and detail simulations of weekly hospitalizations and cumulative
271 hospitalizations can be found in Supplementary Figures S14,S15, and S16.

272 Discussion

273 This study highlights several critical insights into the optimization of SARS-CoV-2 surveillance
274 strategies and enabling proactive variant detection and response. It was found that approximately
275 5000 sequences per week for detection within 4 weeks target enables timely detection and
276 facilitates early intervention measures, and provides a universal guideline for sequencing
277 requirements. Currently, reporting of SARS-CoV-2 sequences in GISAID has been diminishing
278 sharply in all countries, and reached lowest point in July 2025 with 7,802 sequences[16].

279 Unlike traditional surveillance, our DNN model rapidly identified dominant variants using only
280 the first few weeks of prevalence data. It improved this assessment by integrating sigmoid growth
281 parameters and their derivatives. Our findings with low initial-prevalence chains only even
282 showed similar (or even better) classification performance, supporting the idea that early-stage
283 variant detection is feasible even when initial chain prevalence is quite low, which is critical for
284 timely public health responses.

285 These insights highlight a key advantage of deep learning over rule-based classification or
286 threshold-based methods. Instead of relying on arbitrary cutoffs for identifying concerning
287 variants, our model learned from observed epidemiological dynamics, allowing for data-driven
288 decision-making. Compared to previous studies our model demonstrates improved
289 performance[7,8], usability in real-time monitoring[6], and the adaptability and reliability in
290 scenarios with limited sequencing data[8].

291 Our study also reinforces the importance of vaccine matching and genetic monitoring in
292 mitigating variant impact. The distances between consecutive variants and distance from the
293 vaccine strain were among the most significant predictors of hospitalization rates, which was
294 supported by reduced vaccine effectiveness for new variants[28,29].

295 Our simulations indicated a clear link between genetic distance and clinical outcomes at a
296 population level, except in one case, chain 378, which was first reported just before the
297 introduction of the matched XBB.1.5 vaccine. Building upon this, our finding highlights that
298 tracking genetic distance alone can provide actionable insights for hospital capacity planning,
299 even without complete epidemiological data.

300 Strict non-pharmaceutical interventions showed a significant impact in reduction of
301 hospitalizations, and combining with doubling VCR showed additive benefits, which suggests
302 that a multifaceted approach could yield even greater benefits in the absence of matched
303 vaccines. The observed reductions, particularly in chains with lower initial containment levels or
304 delayed vaccine strain matching, emphasize the importance of timely public health measures in
305 the emergence of genetically distinct variants like chains 151 (BA.2), 216 (BA.5), and 418
306 (JN.1).

307 These insights provide a compelling case for sustained investment in genomic surveillance,
308 vaccine development, and integrated public health responses to safeguard against the evolving
309 threat of SARS-CoV-2 and other emerging pathogens.

310 Several limitations should be taken into account in our study. First, the performance of both the
311 clustering algorithm and the deep learning model was inherently dependent on the availability of
312 genomic data, and may introduce biases due to sequencing delays, particularly during the early
313 emergence of new variants. Second, the predictive power of the DNN relied on features extracted
314 during the first few weeks after a variant's initial detection from the clustering algorithm, and
315 variants with delayed or irregular dynamics may challenge the model's classification ability.
316 Thirdly, the data augmentation strategy might not fully capture the diversity of real-world
317 evolutionary behaviors, particularly for variants exhibiting novel or outlier dynamics. Fourthly,
318 its ability to predict the trajectory of newly emerging variants, especially different viruses, in
319 real-time remains to be systematically validated. Future research should explore integrating real-
320 time data streams, improving robustness against sequencing gaps, and evaluating transferability
321 to other viruses. Another limitation is that the proportion of variants among hospitalized patients
322 might be different from that in the general infected population. However, focusing on major
323 variants in the model may partially address this issue, as these variants are more likely to be
324 associated with higher hospitalization rates. Additionally, we were unable to incorporate detailed
325 population structures of vaccine recipients or hospitalized patients due to data availability
326 constraints. Furthermore, the long-term impact of vaccinations from previous seasons on
327 hospitalization rates was not included in our analysis. While this omission may not significantly
328 alter outcomes, it remains a factor worth exploring in future studies to ensure a more
329 comprehensive understanding. Finally, modifying key variables within the CatBoost model

330 introduced wide CIs resulting many scenarios' reduction insignificant. Future research should
331 aim to precisely quantify the relationship between the magnitude of variable modification and the
332 resulting predictive accuracy, thereby refining the ability to generate nuanced and robust hospital
333 burden forecasts. This would not only improve the interpretation of scenario-based predictions
334 but also support the development of more resilient and informative forecasting methodologies.

335 Conclusion

336 The integration of optimized sequencing rates, early-stage classification with deep learning, and
337 genetic-based hospitalization risk assessment offers a powerful, data-driven framework for
338 SARS-CoV-2 surveillance. We show that adequate sequencing volume (impacting on the time
339 needed to detect chains) together with classification models based on growth dynamics and
340 leveraging genetic distance as a key predictor can enhance the ability of public health systems to
341 promptly identify and mitigate the impact of new variants. We demonstrate that early variant
342 classification is feasible using only a few weeks of prevalence data and validating the importance
343 of sigmoid-based features in improving classification accuracy even with low initial-prevalence
344 chains. We prove evidence that genetic distance predicts hospitalization risk, supporting real-time
345 response strategies including vaccine adaptation programs and we highlight the need for
346 integrated deep learning, genetic monitoring, and vaccination strategies for pandemic
347 preparedness.

348 These findings contribute to a growing framework for real-time, adaptive surveillance systems
349 that can rapidly respond to emerging epidemiological threats. Future work will focus on
350 enhancing the robustness of the DNN model, exploring alternative architectures (e.g.,
351 transformer-based models), and refining transient chain modeling techniques to further improve
352 predictive accuracy.

353

354 References

- 355 1. Jackson CB, Farzan M, Chen B, Choe H. Mechanisms of SARS-CoV-2 entry into cells. *Nat*
356 *Rev Mol Cell Biol.* 2022 Jan;23(1):3–20.
- 357 2. Markov PV, Ghafari M, Beer M, Lythgoe K, Simmonds P, Stilianakis NI, et al. The
358 evolution of SARS-CoV-2. *Nat Rev Microbiol.* 2023 Jun;21(6):361–79.
- 359 3. WHO. Updated working definitions and primary actions for SARSCoV2 variants [Internet].
360 2023 [cited 2025 Jul 23]. Available from:
361 [https://www.who.int/publications/m/item/updated-working-definitions-and-primary-actions-](https://www.who.int/publications/m/item/updated-working-definitions-and-primary-actions-for--sars-cov-2-variants)
362 [for--sars-cov-2-variants](https://www.who.int/publications/m/item/updated-working-definitions-and-primary-actions-for--sars-cov-2-variants)
- 363 4. WHO. Statement on the fifteenth meeting of the IHR (2005) Emergency Committee on the
364 COVID-19 pandemic [Internet]. 2023 [cited 2024 Aug 26]. Available from:
365 [https://www.who.int/news/item/05-05-2023-statement-on-the-fifteenth-meeting-of-the-](https://www.who.int/news/item/05-05-2023-statement-on-the-fifteenth-meeting-of-the-international-health-regulations-(2005)-emergency-committee-regarding-the-coronavirus-disease-(covid-19)-pandemic)
366 [international-health-regulations-\(2005\)-emergency-committee-regarding-the-coronavirus-](https://www.who.int/news/item/05-05-2023-statement-on-the-fifteenth-meeting-of-the-international-health-regulations-(2005)-emergency-committee-regarding-the-coronavirus-disease-(covid-19)-pandemic)
367 [disease-\(covid-19\)-pandemic](https://www.who.int/news/item/05-05-2023-statement-on-the-fifteenth-meeting-of-the-international-health-regulations-(2005)-emergency-committee-regarding-the-coronavirus-disease-(covid-19)-pandemic)
- 368 5. de Hoffer A, Vatani S, Cot C, Cacciapaglia G, Chiusano ML, Cimorelli A, et al. Variant-
369 driven early warning via unsupervised machine learning analysis of spike protein mutations
370 for COVID-19. *Sci Rep.* 2022 Jun 3;12(1):9275.
- 371 6. Levi R, Zerhouni EG, Altuvia S. Predicting the spread of SARS-CoV-2 variants: An
372 artificial intelligence enabled early detection. *PNAS Nexus.* 2024 Jan 1;3(1):pgad424.
- 373 7. Nicora G, Salemi M, Marini S, Bellazzi R. Predicting emerging SARS-CoV-2 variants of
374 concern through a One Class dynamic anomaly detection algorithm. *BMJ Health Care*
375 *Inform.* 2022 Dec 9;29(1):e100643.
- 376 8. Rancati S, Nicora G, Prospero M, Bellazzi R, Salemi M, Marini S. Forecasting dominance of
377 SARS-CoV-2 lineages by anomaly detection using deep AutoEncoders. *Briefings in*
378 *Bioinformatics.* 2024 Nov 1;25(6):bbae535.
- 379 9. Feng Y, Goldberg EE, Kupperman M, Zhang X, Lin Y, Ke R. CovTransformer: A
380 transformer model for SARS-CoV-2 lineage frequency forecasting. *Virus Evolution.* 2024
381 Oct 17;10(1):veae086.
- 382 10. Graña C, Ghosn L, Evrenoglou T, Jarde A, Minozzi S, Bergman H, et al. Efficacy and
383 safety of COVID-19 vaccines - Graña, C - 2022 | Cochrane Library. [cited 2025 Jul 23];
384 Available from:
385 <https://www.cochranelibrary.com/cdsr/doi/10.1002/14651858.CD015477/full>
- 386 11. WHO. WHO COVID-19 dashboard [Internet]. datadot. [cited 2025 Apr 18]. Available
387 from: <https://data.who.int/dashboards/covid19/hospitalizations>

- 388 12. Carabelli AM, Peacock TP, Thorne LG, Harvey WT, Hughes J, de Silva TI, et al. SARS-
389 CoV-2 variant biology: immune escape, transmission and fitness. *Nat Rev Microbiol.* 2023
390 Mar;21(3):162–77.
- 391 13. Cromer D, Steain M, Reynaldi A, Schlub TE, Khan SR, Sasson SC, et al. Predicting vaccine
392 effectiveness against severe COVID-19 over time and against variants: a meta-analysis. *Nat*
393 *Commun.* 2023 Mar 24;14(1):1633.
- 394 14. Khoury DS, Docken SS, Subbarao K, Kent SJ, Davenport MP, Cromer D. Predicting the
395 efficacy of variant-modified COVID-19 vaccine boosters. *Nat Med.* 2023 Mar;29(3):574–8.
- 396 15. Cao L, Lou J, Chan SY, Zheng H, Liu C, Zhao S, et al. Rapid evaluation of COVID-19
397 vaccine effectiveness against symptomatic infection with SARS-CoV-2 variants by analysis
398 of genetic distance. *Nat Med.* 2022 Aug;28(8):1715–22.
- 399 16. Khare S, Gurry C, Freitas L, Schultz MB, Bach G, Diallo A, et al. GISAID’s Role in
400 Pandemic Response. *CCDCW.* 2021 Dec 3;3(49):1049–51.
- 401 17. Hale T, Angrist N, Goldszmidt R, Kira B, Petherick A, Phillips T, et al. A global panel
402 database of pandemic policies (Oxford COVID-19 Government Response Tracker). *Nat*
403 *Hum Behav.* 2021 Apr;5(4):529–38.
- 404 18. SSI. SSI’s interaktive dashboards [Internet]. Denmark Statens Serum Institut. [cited 2024
405 Apr 27]. Available from:
406 <https://experience.arcgis.com/template/099eb5c9acea4e18b411997815be2f98>
- 407 19. ECDC. Data on COVID-19 vaccination in the EU/EEA [Internet]. 2020 [cited 2024 Apr
408 27]. Available from: <https://www.ecdc.europa.eu/en/covid-19/data>
- 409 20. Paredes MI, Lunn SM, Famulare M, Frisbie LA, Painter I, Burstein R, et al. Associations
410 between SARS-CoV-2 variants and risk of COVID-19 hospitalization among confirmed
411 cases in Washington State: a retrospective cohort study. *medRxiv.* 2022 Feb
412 16;2021.09.29.21264272.
- 413 21. SSI. Vaccination against COVID-19, influenza and pneumococcal disease [Internet]. 2022
414 [cited 2025 Jun 2]. Available from: <https://en.ssi.dk/news/epi-news/2022/no-35---2022>
- 415 22. SSI. This autumn’s influenza and COVID-19 vaccination programme, 2023/2024 [Internet].
416 2023 [cited 2025 Jun 2]. Available from: <https://en.ssi.dk/news/epi-news/2023/no-39---2023>
- 417 23. SSI. Retningslinje for vaccination mod covid-19 og influenza - historisk [Internet]. 2023
418 [cited 2025 Jul 23]. Available from: [http://www.sst.dk/da/udgivelser/2023/Retningslinje-](http://www.sst.dk/da/udgivelser/2023/Retningslinje-for-vaccination-mod-covid-19-og-influenza)
419 [for-vaccination-mod-covid-19-og-influenza](http://www.sst.dk/da/udgivelser/2023/Retningslinje-for-vaccination-mod-covid-19-og-influenza)
- 420 24. Harman K, Nash SG, Webster HH, Groves N, Hardstaff J, Bridgen J, et al. Comparison of
421 the risk of hospitalisation among BA.1 and BA.2 COVID-19 cases treated with sotrovimab
422 in the community in England. *Influenza and Other Respiratory Viruses.* 2023;17(5):e13150.

- 423 25. Aziz NA, Nash SG, Zaidi A, Nyberg T, Groves N, Hope R, et al. Risk of severe outcomes
424 among SARS-CoV-2 Omicron BA.4 and BA.5 cases compared to BA.2 cases in England.
425 *Journal of Infection*. 2023 Jul 1;87(1):e8–11.
- 426 26. WHO. XBB.1.5 Updated Rapid Risk Assessment [Internet]. 2023 Jan. Available from:
427 [https://www.who.int/docs/default-](https://www.who.int/docs/default-source/coronaviruse/25012023xbb.1.pdf?sfvrsn=c3956081_1)
428 [source/coronaviruse/25012023xbb.1.pdf?sfvrsn=c3956081_1](https://www.who.int/docs/default-source/coronaviruse/25012023xbb.1.pdf?sfvrsn=c3956081_1)
- 429 27. WHO. JN.1 variant update and risk evaluation [Internet]. 2024 Apr. Available from:
430 [https://www.who.int/docs/default-](https://www.who.int/docs/default-source/coronaviruse/15042024_jn1_ure.pdf?sfvrsn=8bd19a5c_7)
431 [source/coronaviruse/15042024_jn1_ure.pdf?sfvrsn=8bd19a5c_7](https://www.who.int/docs/default-source/coronaviruse/15042024_jn1_ure.pdf?sfvrsn=8bd19a5c_7)
- 432 28. Gram MA, Emborg HD, Schelde AB, Friis NU, Nielsen KF, Moustsen-Helms IR, et al.
433 Vaccine effectiveness against SARS-CoV-2 infection or COVID-19 hospitalization with the
434 Alpha, Delta, or Omicron SARS-CoV-2 variant: A nationwide Danish cohort study. *PLOS*
435 *Medicine*. 2022 Sep 1;19(9):e1003992.
- 436 29. Moustsen-Helms IR, Bager P, Larsen TG, Møller FT, Vestergaard LS, Rasmussen M, et al.
437 Relative vaccine protection, disease severity, and symptoms associated with the SARS-
438 CoV-2 omicron subvariant BA.2.86 and descendant JN.1 in Denmark: a nationwide
439 observational study. *The Lancet Infectious Diseases*. 2024 Sep 1;24(9):964–73.

440

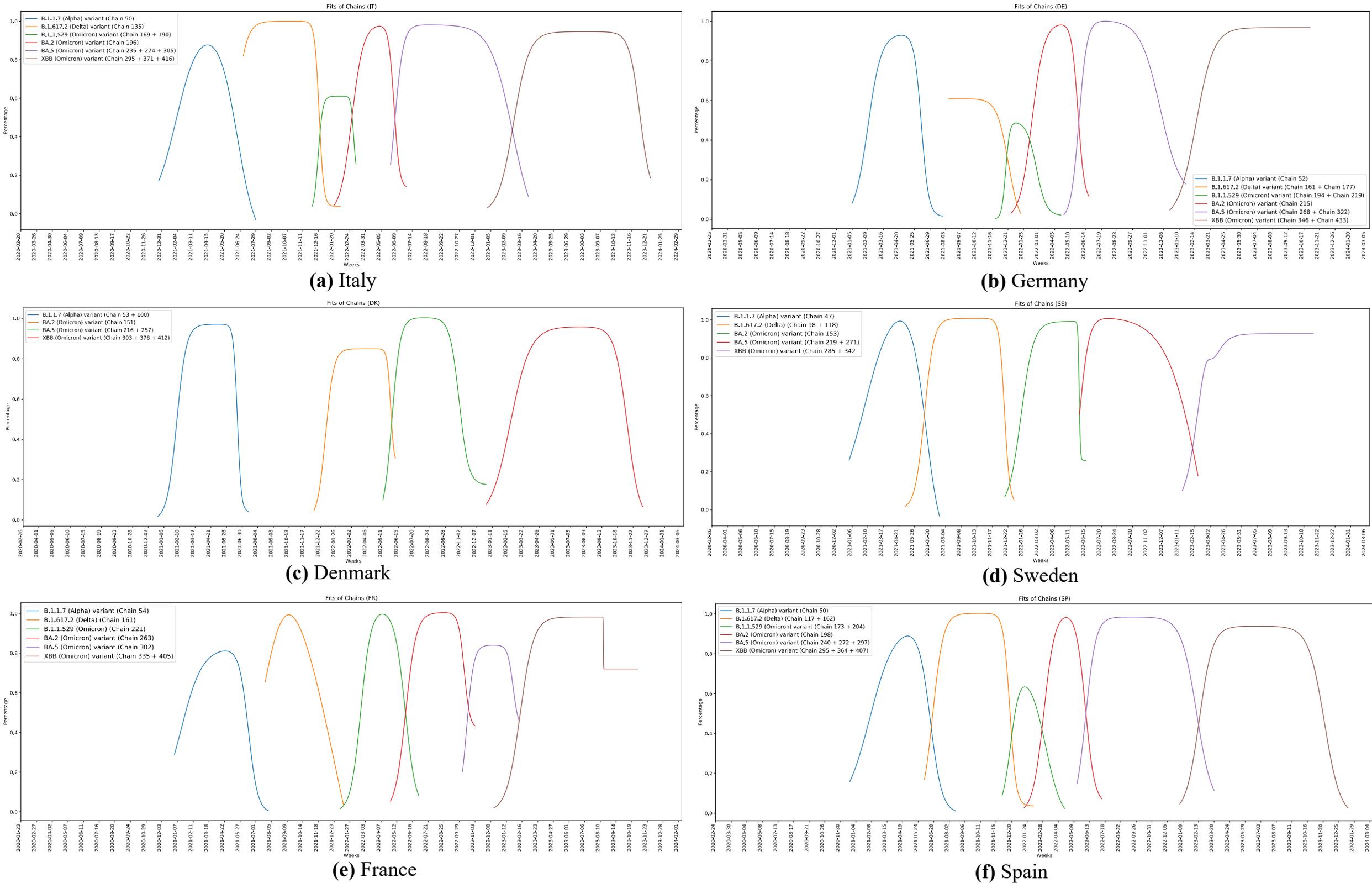


Figure 1. Fitted chains for Italy (a), Germany (b), Denmark (c), Sweden (d), France (e), and Spain (f) from 2020 to 2023. For each country, chains representing the evolution of a variant were fitted using a combination of two sigmoid functions to capture the increase and decrease phase of variant prevalence. On the x-axis, time in weeks is plotted, while on the y-axis the percentage of prevalence (from 0 to 1) of the variant of the chain with respect to all the sequences is considered. Chains are numbered by the algorithm and numbers are assigned according to the number of the first cluster of that chain [5].

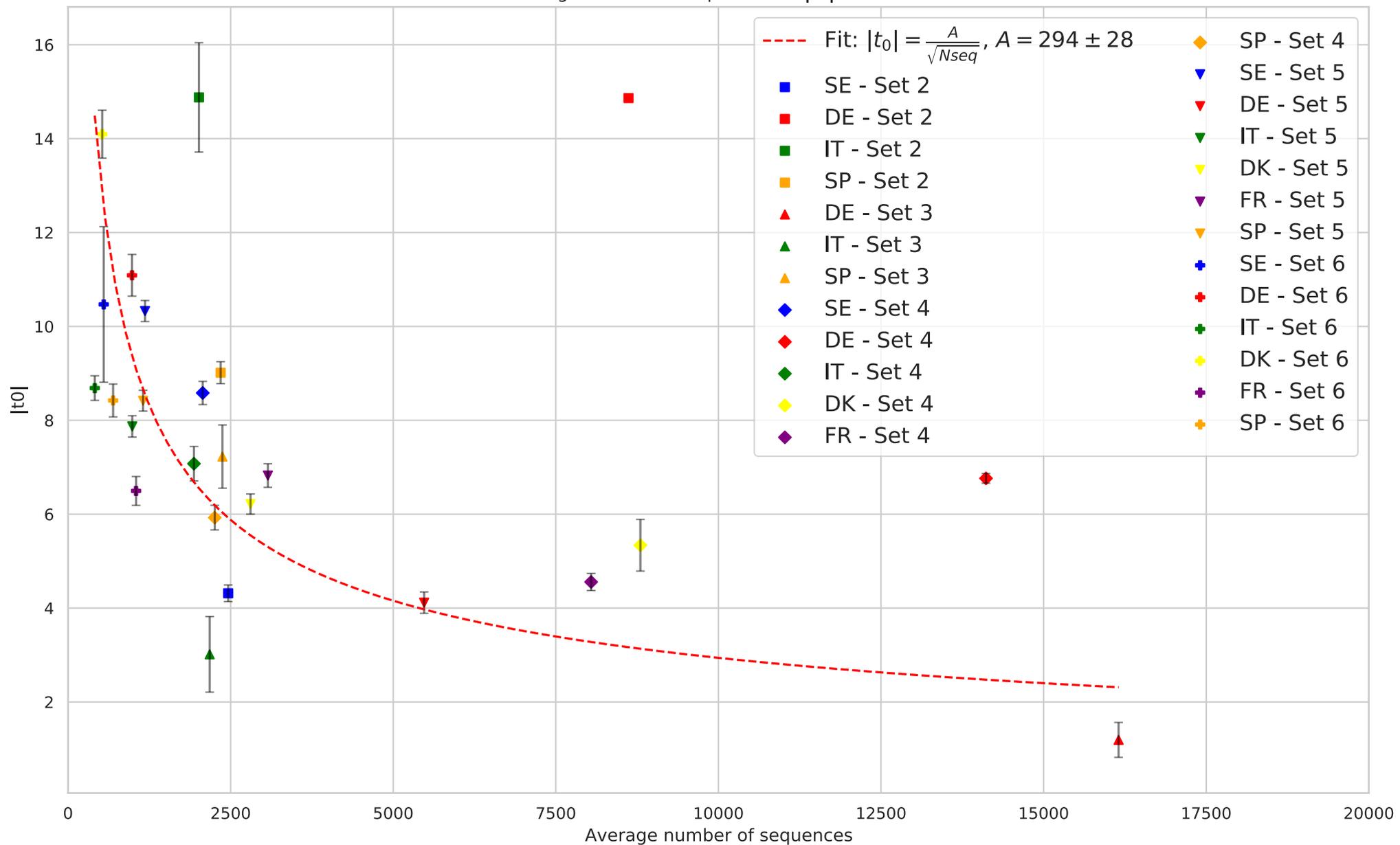
Average number of sequences vs $|t_0|$ for all sets

Figure 2. Calibration curve. The average number of sequences available per week is plotted against the absolute value of the first detection time, $|t_0|$, for variant chains across different country sets. Each point corresponds to a variant detected in a specific country, with marker shape indicating the set and color indicating the country. Error bars represent uncertainty on $|t_0|$ from the fitting procedure. A red dashed curve shows the best fit according to an inverse power law model, $|t_0| = \frac{A}{\sqrt{x}}$, where x is the average weekly number of sequences. The fitted parameter is $A = 294 \pm 28$.

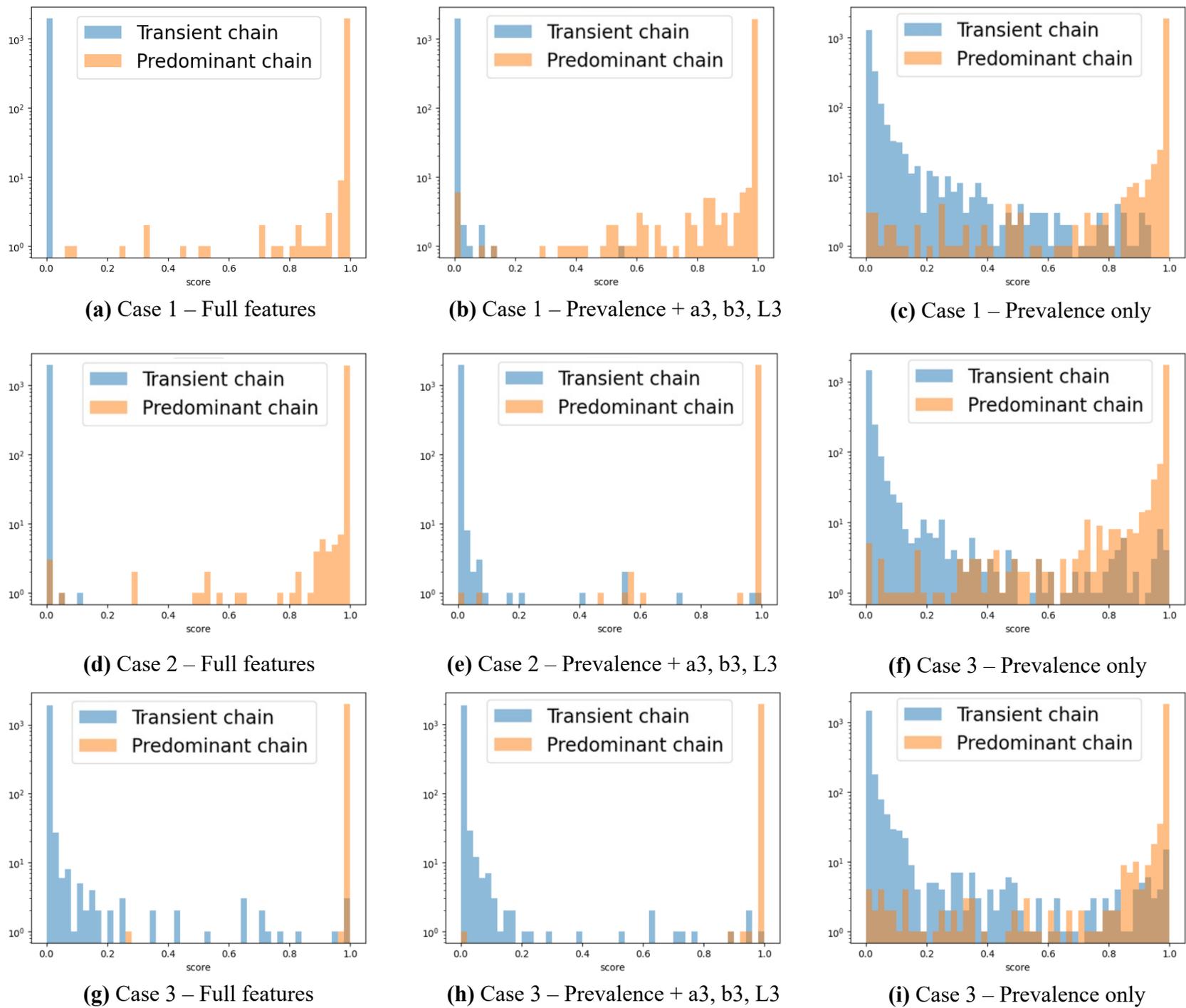
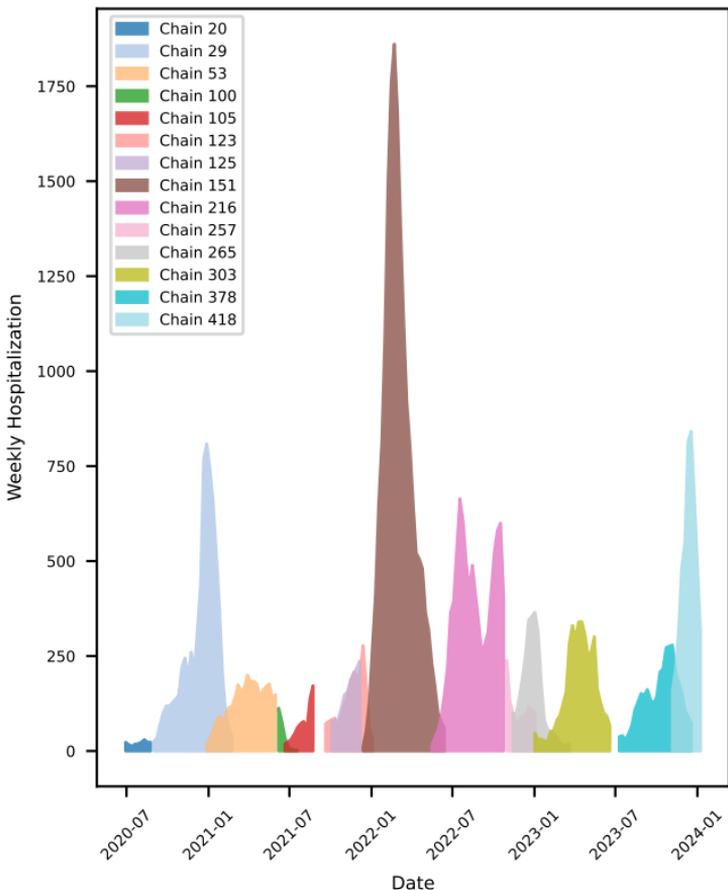


Figure 3. DNN score distribution comparison across dataset configurations and feature sets. The signal-to-noise ratio (i.e. Predominant chains to Transient chains ratio) quantifies the separation between true positive and false positive predictions for different classifier scores. Including growth-related parameters consistently improves the signal-to-noise ratio, particularly for challenging early-stage chains.

Weekly Hospitalization by Chain



Genetic Distances (LD1 & LD_vac)

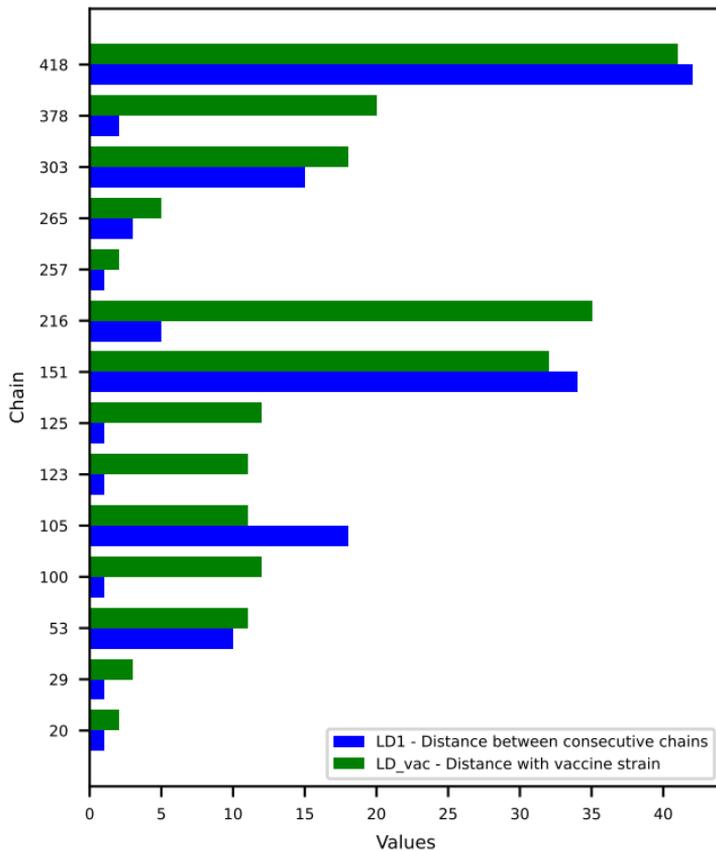


Figure 4. Weekly hospitalizations and Levenshtein distances between consecutive chains and their respective vaccine strain distances for each chain, observed in Denmark from July 2020 to January 2024. Variant classifications : Chain 53 - Alpha, 100 - Delta, 151 - Omicron BA.2, 216 - BA.5, 303 - XBB.1.5, 418 - JN.1.

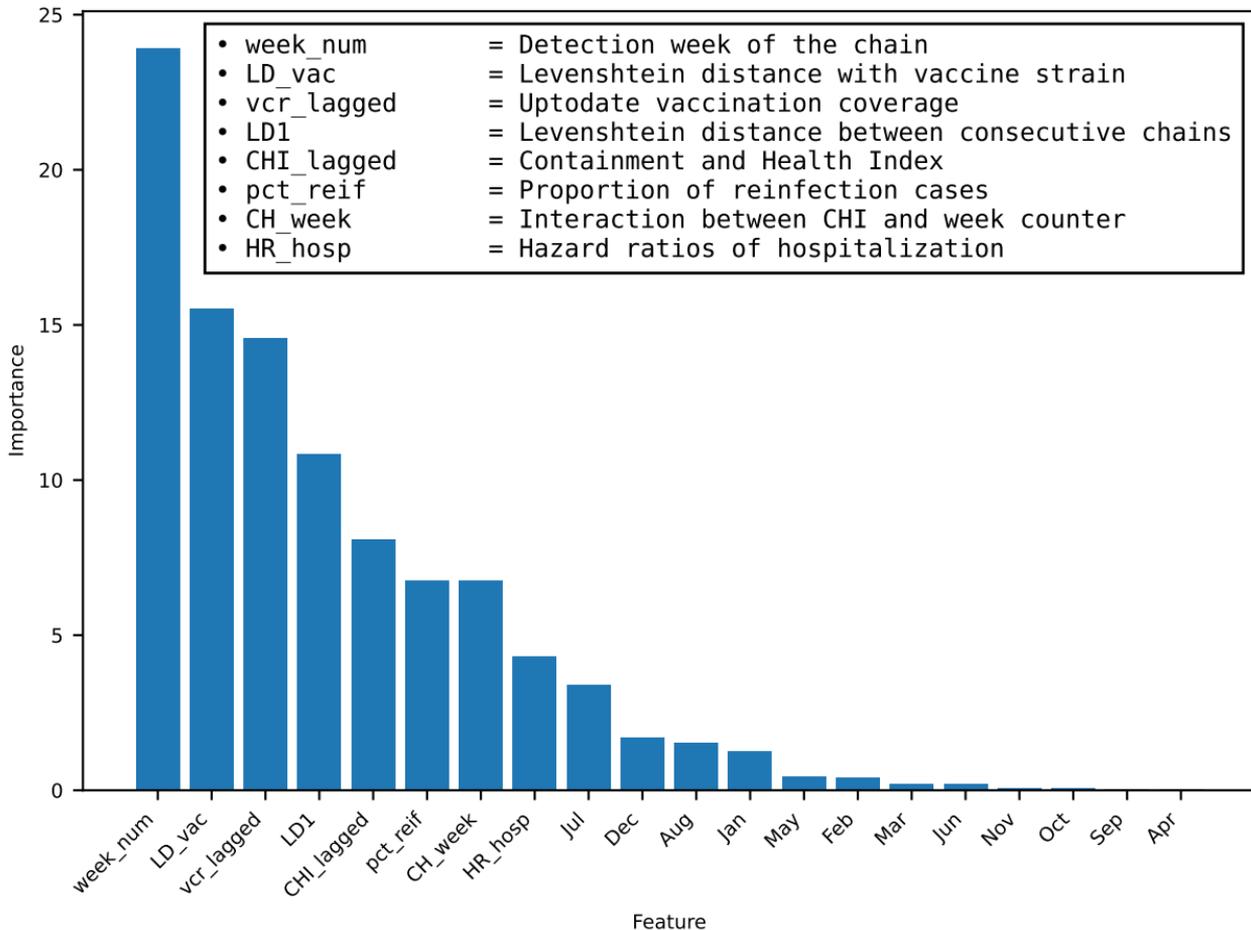
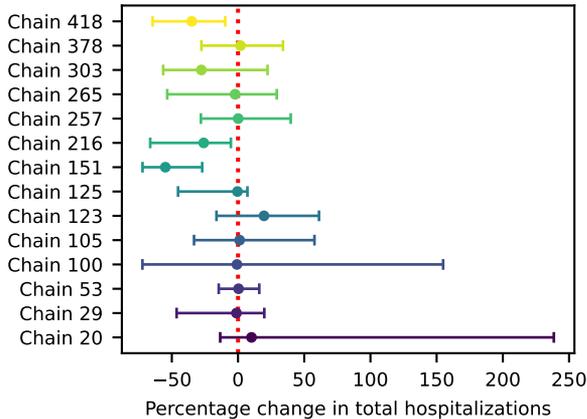
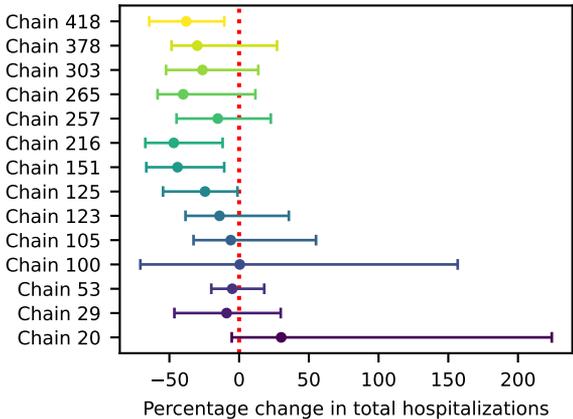


Figure 5. Feature importance of the CatBoost model for prediction of weekly hospitalizations by chains in Denmark (Jul 2020 - Jan 2024).

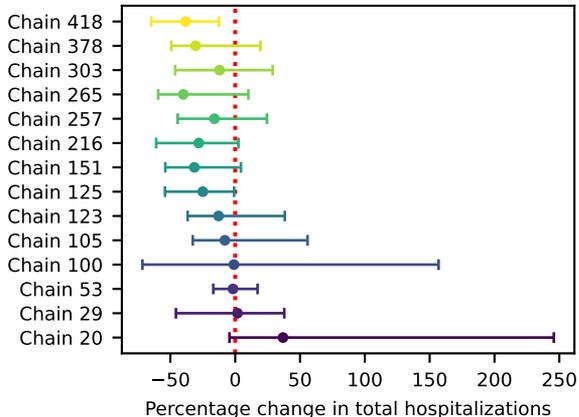
Simulation 1 - better-matched vaccine strain



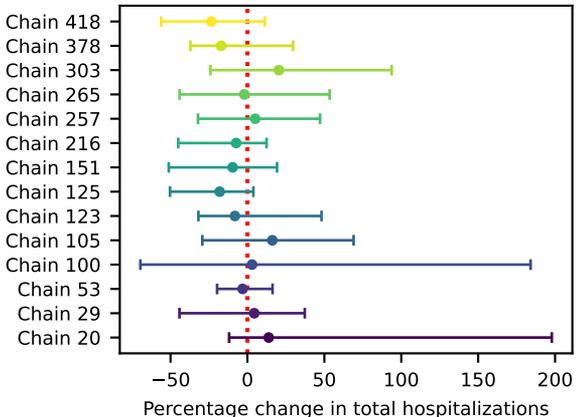
Simulation 2 - sustained high CHI



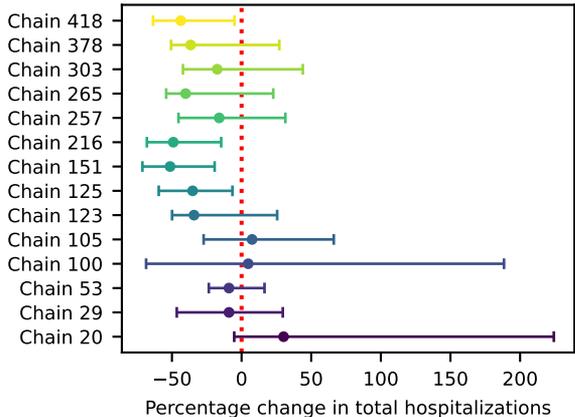
Simulation 3 - temporary high CHI



Simulation 4 - doubling Vaccination Coverage (VCR)



Simulation 5 - sustained high CHI + doubling VCR



Simulation 6 - temporary high CHI + doubling VCR

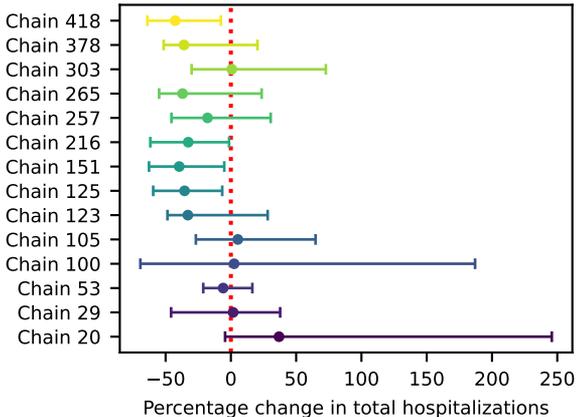


Figure 6. Changes in the proportion of total hospitalizations per chain in Denmark, across different scenarios, with respective 95 % confidence intervals.