

1 **Application of ARIMA, hybrid ARIMA and Artificial Neural Network Models in**
2 **predicting and forecasting tuberculosis incidences among children in Homa Bay and**
3 **Turkana Counties, Kenya**

4
5 Siamba Stephen^{1*}, Otieno Argwings¹, Koech Julius¹

6
7 ¹University of Eldoret, School of Science, Department of Mathematics and Computer Science

8
9 *Corresponding author

10 Email: stephen.siamba@gmail.com

26 **Abstract**

27 **Background**

28 Tuberculosis (TB) infections among children (below 15 years) is a growing concern, particularly
29 in resource-limited settings. However, the TB burden among children is relatively unknown in
30 Kenya where two-thirds of estimated TB cases are undiagnosed annually. Very few studies have
31 used Autoregressive Integrated Moving Average (ARIMA), hybrid ARIMA, and Artificial
32 Neural Networks (ANNs) models to model infectious diseases globally. We applied ARIMA,
33 hybrid ARIMA, and Artificial Neural Network models to predict and forecast TB incidences
34 among children in Homa bay and Turkana Counties in Kenya.

35 **Methods**

36 The ARIMA, ANN, and hybrid models were used to predict and forecast monthly TB cases
37 reported in the Treatment Information from Basic Unit (TIBU) system for Homa bay and
38 Turkana Counties between 2012 and 2021. The data were split into training data, for model
39 development, and testing data, for model validation using an 80:20 split ratio respectively.

40 **Results**

41 The hybrid ARIMA model (ARIMA-ANN) produced better predictive and forecast accuracy
42 compared to the ARIMA (0,0,1,1,0,1,12) and NNAR (1,1,2) [12] models. Furthermore, using the
43 Diebold-Mariano (DM) test, the predictive accuracy of NNAR (1,1,2) [12] versus ARIMA-
44 ANN, and ARIMA-ANN versus ARIMA (0,0,1,1,0,1,12) models were significantly different,
45 $p < 0.001$, respectively. The 12-month forecasts showed a TB prevalence of 175 to 198 cases per
46 100,000 children in Homa bay and Turkana Counties in 2022.

47 **Conclusion**

48 The hybrid (ARIMA-ANN) model produces better predictive and forecast accuracy compared to
49 the single ARIMA and ANN models. The findings show evidence that the prevalence of TB
50 among children below 15 years in Homa bay and Turkana Counties is significantly under-
51 reported and is potentially higher than the national average.

52

53

54

55

56

57

58

59

60

61

62

63

64

65

66

67 **Introduction**

68 **Background**

69 Tuberculosis (TB) is a highly infectious disease ranked among the top ten most lethal
70 causes of mortality. Approximately 33% of the global population has been plague-ridden with
71 TB, particularly in developing countries [1]. In 2016, over 10 million new TB cases were
72 reported globally with children below 15 years of age accounting for about 7% of those cases
73 [3]. In 2018, about 1 million TB disease cases and over 230,000 TB-related deaths occurred
74 among children below 15 years with about 55% of these reported TB cases going undiagnosed
75 and/or unreported [2].

76 Pediatric TB is usually overlooked [4] amid diagnosis and treatment challenges.
77 Developing countries account for over 85% of new cases of TB globally with Asian and African
78 countries contributing 61% and 25% of global new TB cases [2] respectively. In 2016,
79 approximately 7 countries, globally, accounted for close to 65% of all new TB cases [2] while in
80 2019, 30 high TB burdened countries accounted for 87% of all new TB cases while only 8
81 countries accounted for approximately 67% of the total new TB cases [5].

82 The TB burden in Sub-Saharan Africa (SSA) is far much greater and is exacerbated by
83 poverty, political strife, and weak health systems which have curtailed the implementation of TB
84 control interventions. Consequently, TB has become an enormous burden to health systems that
85 are already overstretched [6].

86 Tuberculosis is a disease of major concern in Kenya and is among the top five causes of
87 mortality. Kenya is listed among the top 30 TB high burdened countries [7]. Kenya is also
88 among 14 countries globally that suffer from the TB, TB-HIV co-infection, and Multi-Drug
89 Resistant TB [8] triple burden. The TB prevalence for Kenya in 2015 was 233 per 100,000

90 population with a mortality of 20 per 100,000. In Kenya, the TB case notification increased from
91 11,000 to 116,723 between 1990 and 2007 [9] occasioned by the HIV epidemic and improved
92 case detection due to improved diagnostic capacity.

93 The use of mathematical models in the modeling of epidemic interactions within
94 populations has been detailed extensively. While existing interventions to control TB have been
95 partially successful, within the context of resource constraints, mathematical modeling can
96 increase understanding and result in better policies toward the implementation of effective
97 strategies that would compound better health and economic benefits [10]. In addition,
98 mathematical models are essential in guiding policymakers in resource allocation toward the
99 prevention and control of diseases.

100 Several studies have utilized ARIMA, Seasonal ARIMA (SARIMA), neural network, and
101 hybrid ARIMA models to model TB incidences in other countries [11, 13] and in these studies,
102 the hybrid models were demonstrated to offer better predictive and forecast accuracy. In Africa,
103 *Azeez et al.* compared the predictive capabilities of the SARIMA and the hybrid SARIMA
104 neural network auto-regression (SARIMA-NNAR) models in modeling TB incidences in South
105 Africa and the SARIMA-NNAR model was found to have better goodness-of-fit [12]. In
106 addition, *Li et al.* compared the the predictive power of the ARIMA and ARIMA-generalized
107 regression neural network (GRNN) hybrid models in forecasting TB incidences in China and
108 concluded that the hybrid model was superior to the single ARIMA model [13].

109 The ARIMA and Neural Networks models have also been applied in modeling other
110 infectious diseases. *Zeming and Yanning* predicted HIV-AIDS incidences in China in 2017 using
111 the back propagation (BP) artificial neural network and the ARIMA models and concluded that

112 the hybrid (BP-ANN) model offered better predictive power compared to the single ARIMA
113 model [14]. Zhou *et al.* modeled the prevalence of schistosomiasis in Qianjiang city in China
114 using a hybrid model of ARIMA-NARNN (Nonlinear Autoregressive Neural Network) and
115 concluded that the hybrid ARIMA-NARNN model produced high-quality prediction accuracy
116 [15]. Yu *et al.* used the hybrid seasonal ARIMA and NARNN model to forecast incidences of the
117 Coxsackie viral infection in Shenzhen China, and concluded the hybrid seasonal ARIMA-
118 NARNN was the best model [16].

119 While hybrid ARIMA models have been applied in forecasting both the short-term and
120 long-term incidences of infectious diseases in other countries, there has been little to no
121 application of these cutting-edge methods in African countries with the majority of the models
122 limited to only ARIMA models [17]. In Kenya, while ARIMA models have been applied in
123 forecasting disease incidence [18], very little has been done in the application of hybrid ARIMA
124 models in predicting disease incidence except in non-public health settings such as agriculture
125 and economics.

126 The popularity of ARIMA models stems from their flexibility to represent varieties of
127 time series with simplicity but with a profound limitation stemming from their linear
128 assumptions which in many cases is usually impractical [19] since real-world applications
129 mainly involve data exhibiting non-linear patterns. Consequently, to overcome this disadvantage,
130 non-linear stochastic models such as the ANN models have been proposed [20]. Despite this, a
131 single ANN model is not able to incorporate both linear and non-linear patterns and this has led
132 to the adoption of hybrid models that can address this challenge [21]. To attain a higher degree of

133 predictive and forecasting accuracy, theoretical and empirical findings show that combining
134 different models can be effective [22].

135 To better understand the status of TB infection among children in Kenya, it is important
136 to assess the trend and forecast these incidences using available surveillance data and novel
137 models to elicit a better understanding and innovative interventions to curtail the spread of
138 pediatric TB in Kenya. This study compares linear-based ARIMA, non-linear-based ANN, and
139 hybrid ARIMA in modeling TB incidences among children below years in Homa bay and
140 Turkana counties in Kenya.

141 **Materials and methods**

142 **Study design and setting**

143 This was a retrospective study that utilized aggregated monthly TB cases among children
144 data reported by health facilities located in Homa Bay and Turkana Counties to the National
145 Tuberculosis, Leprosy and Lung Disease Program (NTLLDP) in the Treatment Information from
146 Basic Unit (TIBU) electronic system between January 2012 to December 2021. The data
147 comprised 120 data points. The study utilized data reported by health facilities in Homa bay and
148 Turkana Counties which are among the top 10 TB endemic Counties in Kenya [50].

149 **Data collection and analysis**

150 Tuberculosis case data were abstracted and aggregated for each month between January
151 2012 to December 2021 for health facilities located in Homa bay and Turkana Counties in
152 Kenya. In 2012, the Kenya Ministry of Health (MoH) through the Division of Leprosy,
153 Tuberculosis and Lung Disease transitioned the reporting of TB cases from paper-based to the
154 TIBU system [23]. The TIBU system is a national case-based surveillance system used in the

155 storage of individual cases of TB that are reported to the national TB program monthly with
156 nationwide coverage [24]. This study did not collect or utilize patient-level data.

157 The 120 TB cases data points in the data used in this study were split using the 80:20
158 ratio with 80% (2012 to 2019) of the data assigned to the training set and 20% (2020 to 2021)
159 assigned to the testing set. Splitting the data into training and testing data using the 80:20 ratio,
160 in this case, was based on the available chronological data points as this has been proved to yield
161 test error rate estimates with reasonably low bias and variance [25]. Furthermore, when the
162 training data is large enough, the model learns well and this gives predicted values that are much
163 closer to the actual values [26]. The models were built on the training set and their performance
164 was evaluated on the testing set; this method is the holdout-validation method for model
165 performance evaluation.

166 Data analysis was performed using R statistical software [27] together with applicable
167 packages for analyzing time-series data. The results were summarized using tables and figures.

168 **The Time Series concept**

169 A time series is a sequential set of data points measured over time and is typically
170 composed of the trend, cyclical, seasonal, and irregular (random) components.

171 An autoregressive (AR) model is a type of random process used to describe certain time-
172 varying processes within a time series [28]. The basic idea of AR models is that the present value
173 of a series Y_t can be linearly explained by a function of p past values, that is,
174 $Y_{t-1}, Y_{t-2}, \dots, Y_{t-p}$.

175 Assuming $E(Y_t) = 0$, an AR process of order p can be written as;

176
$$Y_t = \phi_1 Y_{t-1} + \phi_2 Y_{t-2} + \dots + \phi_p Y_{t-p} + \varepsilon_t, \quad (1)$$

177 Where ε_t is white noise (WN), and is uncorrelated with Y_s for all $s < t$

178 However, if the mean is $E(Y_t) = \mu \neq 0$, then Y_t is replaced by $Y_t - \mu$ to obtain;

179
$$Y_t - \mu = \phi_1 (Y_{t-1} - \mu) + \phi_2 (Y_{t-2} - \mu) + \dots + \phi_p (Y_{t-p} - \mu) + \varepsilon_t \quad (2)$$

180 Equation 2 can also be written as;

181
$$Y_t = \alpha + \phi_1 Y_{t-1} + \phi_2 Y_{t-2} + \dots + \phi_p Y_{t-p} + \varepsilon_t \quad (3)$$

182 Where;

183
$$\alpha = \mu(1 - \phi_1 - \dots - \phi_p)$$

184 Furthermore, equation (1) can be written in the form;

185
$$Y_t - \phi_1 Y_{t-1} - \phi_2 Y_{t-2} - \dots - \phi_p Y_{t-p} = \varepsilon_t \quad (4)$$

186 However, by applying the backshift operator $BY_t = Y_{t-1}$, we get;

187
$$(1 - \phi_1 B - \phi_2 B^2 - \dots - \phi_p B^p)Y_t = \varepsilon_t$$

188 Or using notation, we can write;

189
$$\phi(B)Y_t = \varepsilon_t, \quad (5)$$

190 Where $\phi(B)$ denotes the autoregressive (AR) operator;

191
$$\phi(B) = 1 - \phi_1 B - \phi_2 B^2 - \dots - \phi_p B^p \quad (6)$$

192 A moving average (MA) model uses the dependency between an observed value and the
193 residual error from a moving average model applied to the lagged observations. This implies that
194 the output variable is linearly dependent on the current and past values of a stochastic term [28].

195 Consequently, Y_t is a moving average process of order q if

196
$$Y_t = \varepsilon_t + \theta_1 \varepsilon_{t-1} + \theta_2 \varepsilon_{t-2} + \dots + \theta_q \varepsilon_{t-q}, \quad (7)$$

197 Where ε_t is WN and $\theta_1, \dots, \theta_q$ are constants

198 On the other hand, a MA(q) process can also be written in the form;

$$199 \quad Y_t = \theta(B)\varepsilon_t, \quad (8)$$

200 Where the moving average operator;

$$201 \quad \theta(B) = 1 + \theta_1 B + \theta_2 B^2 + \dots + \theta_q B^q = 1 + \sum_{j=1}^q \theta_j B^j \quad (9)$$

202 defines a linear combination of values in the shift operator $B^k \varepsilon_t = \varepsilon_{t-k}$

203 **Autoregressive Integrated Moving Average models (ARIMA) models**

204 An ARMA (p, q) model is a class of stochastic processes whose auto-covariance
205 functions depend on a finite number of unknown parameters. Generally, an ARMA process of
206 orders p and q can be represented mathematically [29] as;

$$207 \quad Y_t = \sum_{j=1}^p \phi_j Y_{t-j} + \sum_{j=1}^q \theta_j \varepsilon_{t-j} + \varepsilon_t \quad \forall t \in \mathbb{Z} \quad (10)$$

208 In the lag operator notation, the ARMA (p, q) process is given by

$$209 \quad \phi(B)Y_t = \theta(B)\varepsilon_t \quad \forall t \in \mathbb{Z} \quad (11)$$

210 Box and Jenkins introduced the ARIMA model in 1960 [30]. The ARIMA model requires
211 only historical time series data on the variable under forecasting. Most importantly, ARIMA
212 models are represented as ARIMA (p, d, q) where p is the number of AR terms, d is the number
213 of non-seasonal differences, and q is the number of lagged forecast errors [31]. The ARIMA
214 model assumes that the residuals are independent and normally distributed with homogeneity of
215 variance and zero mean value.

216 **Seasonal Autoregressive Integrated Moving Average models (SARIMA) models**

217 The SARIMA model is made up of non-seasonal and seasonal components in a multiplicative
218 model. A SARIMA model can be written as $ARIMA(p, d, q)(P, D, Q)_S$ where p is the non-
219 seasonal AR order, d is the non-seasonal differencing, q is the non-seasonal MA order, P is the
220 seasonal AR order, D is the seasonal differencing, Q is the seasonal MA order and S is the period
221 of repeating seasonal pattern. Generally, $S=12$ for monthly data.

222 Let the backshift operator be presented as $BY_t = Y_{t-1}$

223 Without differencing, a SARIMA model can be written formally as;

$$224 \vartheta(B^S)\phi(B)(Y_t - \mu) = \theta(B)\Theta(B^S)\varepsilon_t, \forall t \in \mathbb{Z} \quad (12)$$

225 Where on the left of equation 12, the seasonal and non-seasonal AR processes multiply each
226 other, and on the right of equation 12, the seasonal and non-seasonal MA processes multiply
227 each other. Also, in this study, $S=12$ since monthly TB case data are used.

228 **Artificial Neural Networks (ANNs) models**

229 Artificial Neural Networks have been suggested as alternative and better modeling
230 approaches to time series forecasting [32]. The main goal of ANNs is to construct a model that
231 mimics the human brain intelligence into a machine [33] and is biologically motivated [34]. The
232 most common ANNs are multi-layer perceptrons (MLPs) that utilize a single hidden layer feed-
233 forward network (FNN) [35] made up of the input layer, the hidden layer, and the output layer
234 connected by acyclic links [36]. A neuron is a data processing unit while the nodes in the various
235 layers of ANNs are the processing elements. In this study, the inputs were the lagged TB case
236 observations and the outputs were the predicted or fitted TB cases from the model. According to
237 Zhang [37], the ANN can be written as:

238
$$Y_t = \gamma_0 + \sum_{j=1}^q \gamma_j h(\beta_{0j} + \sum_{i=1}^p \beta_{ij} Y_{t-i}) + \varepsilon_t, \forall t \quad (13)$$

239 Where Y_{t-i} ($i=1, 2, \dots, p$) are the p inputs, Y_t is the output, q are the hidden nodes, γ_j ($j=0, 1, 2,$
240 \dots, q) and β_{ij} ($i=0, 1, 2, \dots, p; j=0, 1, 2, \dots, q$) are the connection weights and ε_t is the random
241 shock; γ_0 and β_{0j} are the bias terms. There is no systematic rule in deciding the choice of q while
242 p , which is the number of neurons is equal to the number of features in the data. The logistic
243 function $h(\cdot)$ is applied as the nonlinear activation function, where:

244
$$h(\cdot) = \frac{1}{1+e^{-x}} \quad (14)$$

245 As a matter of fact, the model in equation 13 performs a nonlinear functional mapping from past
246 observations of a time series to the future value. That is:

247
$$Y_t = f(Y_{t-1}, Y_{t-2}, \dots, Y_{t-p}, v) + \varepsilon_t \quad (15)$$

248 Where v is a vector of all parameters and $f(\cdot)$ is a function determined by the structure of the
249 network and the connection weights.

250 **Hybrid (ARIMA-ANN) models**

251 Generally, a time series can be observed as having linear and nonlinear components as
252 shown in equation 16.

253
$$Y_t = I_t + n_t \quad (16)$$

254 Where I_t and n_t are the linear (from the ARIMA model) and nonlinear (ANN fitted
255 ARIMA model residuals) components respectively. Residuals from the ARIMA model are fitted
256 with the ANN model.

257 **Proposed Methodology**

258 The proposed methodology for this study was based on the combination of the Box-
259 Jenkins methodology for ARIMA modeling, and the ANN and hybrid ARIMA models as shown
260 in Fig 1.

261 **Fig 1: The Proposed methodology**

262 **Model identification and specification**

263 Optimal values of p , d , and q for the ARIMA model were determined by examining the
264 autocorrelation functions and the best model was determined by testing models with different
265 parameters of p , d , and q . The models were estimated using the maximum likelihood estimation
266 (MLE) method and the Akaike Information Criterion (AIC) and Bayesian Information Criterion
267 (BIC) [38] penalty function statistics were used to determine the best model that minimizes AIC
268 or BIC.

269 One assumption of the ARIMA model is that the residuals should be white noise. As
270 such, the Ljung-Box Q test [39] was used to test the hypothesis of independence, constant
271 variance and zero mean of the model residuals.

272 **Accuracy measures**

273 Various accuracy measures have been proposed [40] to determine predictive and forecast
274 performance. This study used the Root Mean Squared Error (RMSE), Mean Absolute Error
275 (MAE), and the Mean Absolute Percent Error (MAPE) to measure the predictive and forecast
276 accuracy of the three models. The lower the values of these accuracy measures the better the
277 model. Furthermore, MAPE values of 10% or below, 10-20%, and 20-50% should be considered
278 as high accuracy, good accuracy, and reasonable accuracy [41].

279 The study also compared the predictive accuracy of the forecasts from the three models
280 using the Diebold-Mariano (DM) test [42]. The test was used to test the null hypothesis that two
281 models have similar predictive accuracy.

282 **Ethical approval and considerations**

283 Authorization for use of the data from the TIBU system was obtained from the National
284 Commission for Science, Technology, and Innovation (NACOSTI) through a research permit
285 and letters of approval from the research unit of the department of health services Homabay
286 County and the department of medical services Turkana County.

287 **Results**

288 **Exploratory data analysis**

289 There was a total of 120 data points in this dataset. The trend of the TB cases among
290 children below 15 years in Homa bay and Turkana counties in the data is shown in Fig 2
291 showing a marginal increase in the TB cases reported between 2018 and 2021.

292 **Fig 2: Monthly TB cases among children below 15 years from Homabay and Turkana** 293 **Counties between 2012 and 2021**

294 The monthly cycle box plot of TB cases is shown in Fig 3 where there is a potential
295 presence of seasonality within the reported TB cases. However, whether or not to account for
296 seasonality in the model depends on whether this would improve model accuracy. This implies
297 that there is need to account for seasonality within the ARIMA model. Furthermore, outliers
298 were detected in some months.

299 **Fig 3: Monthly cycle plot of TB cases**

300

301 **Comparison of model performance in predicting TB cases among children below 15 years**
 302 **in Homabay and Turkana Counties, Kenya**

303 **ARIMA Model estimation and accuracy**

304 The Akaike Information Criterion (AIC) and the Bayesian Information Criterion (BIC)
 305 were used to pick the best parsimonious model based on the least AIC or BIC estimated values.
 306 The best model was ARIMA (0,0,1,1,0,1,12); where $p=0$, $d=0$ and $q=1$ respectively and $P=1$,
 307 $D=0$ and $Q=1$ respectively. The Ljung-Box Q test for the best model showed a p-value of 0.971
 308 implying that the ARIMA (0,0,1,1,0,1,12) model residuals were independently distributed.
 309 The best model was made up of non-differenced seasonal AR (1), non-seasonal MA (1) model
 310 and seasonal MA (1) polynomials and using equation 12, the ARIMA (0,0,1,1,0,1,12) model
 311 equation can be written as:

$$312 (1 - \theta_1 B^{12})(Y_t - \mu) = (1 + \Theta_1 B^{12})(1 + \theta_1 B)\varepsilon_t, \forall t \in \mathbb{Z} \quad (17)$$

313 Where ε_t is WN, and uncorrelated with Y_s for each $s < t$

314 From the model output, the estimated coefficients were (see Table 1); $ma1 = \theta_1 = 0.291$, $sar1 =$
 315 $\theta_1 = 0.997$, $sma1 = \Theta_1 = -0.953$ and $\mu = \text{Intercept} = 50.902$

316 Plugging these estimated coefficients into equation 17, yields the model equation:

$$317 (1 - 0.997B^{12})(Y_t - 50.902) = (1 - 0.953B^{12})(1 + 0.291B)\varepsilon_t, \forall t \in \mathbb{Z} \quad (18)$$

318 **Table 1: Estimated model coefficients**

	Coefficients	Std. Error	Z-value	Pr(> z)
ma1	0.291	0.108	2.701	0.007*
sar1	0.997	0.015	68.167	<0.001**
sma1	-0.953	0.127	-7.512	<0.001**
Intercept	50.902	5.064	10.052	<0.001**

319 Significance codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

320 **ARIMA Model residual diagnostics**

321 The best ARIMA model was assessed for fit using the standard model residual analysis.
322 In model diagnostic checking, 4 plots were used to test the underlying assumptions as shown in
323 Fig 4. The Q-Q plot was relatively normal except for a few outliers at the tails, with model
324 residuals being normally distributed. Inspection of the Autocorrelation Function (ACF) and
325 Partial Autocorrelation Function (PACF) plots to test residual randomness to identify patterns or
326 extreme values showed significant auto-correlations at lag 3.

327 **Fig 4: ARIMA Model residual diagnostics**

328 **ARIMA model performance**

329 The performance of the ARIMA (0,0,1,1,0,1,12) model was carried out by comparing
330 predicted and forecasted TB cases with the actual TB cases reported and presented in Fig 5 and
331 Fig 6 respectively.

332 **Fig 5: Comparison of ARIMA predicted versus actual (training data) TB cases**

333 **Fig 6: Comparison of ARIMA 24-month forecast versus actual (test data) TB cases**

334 **ARIMA (0,0,1,1,1,0,1,12) model accuracy**

335 Table 2 compares the accuracy of parameters/measures of the ARIMA (0,0,1,1,0,1,12)
336 model on the training (2012 to 2019) and testing (2020 to 2021) data. The accuracy measures
337 compared were RMSE, MAE, and MAPE correspondingly.

338 The accuracy measure comparison in table 2 demonstrates that the model performs slightly
339 worse on the testing data an RMSE value of 29.17 compared to 18.74 on the training data.

340 Large RMSE values generally imply that the fitted model fails to account for important
341 information within the underlying data hence the need to account for such information which
342 might be captured within the non-linearities as shown in the additive hybrid model methodology
343 which provides fitted values that are as closer to the actual/observed data.

344 **Table 2: Model (ARIMA (0,0,1,1,0,1,12)) accuracy comparison on training and testing data**

Data	RMSE	MAE	MAPE
Training	18.74	14.39	39.00
Testing	29.17	20.47	33.15

345 **Artificial Neural Network (ANN) model estimation and accuracy**

346 The training dataset was fit using an ANN model using the Neural Network Auto-
347 Regressive (NNAR) function to produce an NNAR (p,P,k) [m] model. The optimal lag
348 parameter, p , and the number of nodes in the hidden layer, k , were automatically selected while
349 $P=1$ by default. In addition, a decay parameter of 0.001 and a maximum iteration of 200 were
350 pre-set for the model to help restrict the weights from becoming too large and ensure that the
351 model can test different models until the optimal model that has the minimal RMSE produced
352 respectively.

353 The optimal NNAR model produced was NNAR (1,1,2) [12] with an average of 20
354 networks each of which was a 2-2-1 network with 9 weights. The plot of the point forecast
355 values from the NNAR (1,1,2) [12] model on the training set against the actual training data were
356 presented in Fig 7 while Fig 8 presents the 24-month forecast based on the NNAR (1,1,2) [12]
357 model. Table 3 shows the accuracy measures from the model. The model seems to perform badly
358 on the testing dataset although the MAPE values below 50% are reasonable.

359 **Fig 7: Comparison of NNAR (1,1,2) [12] predicted TB cases versus actual TB cases**
360 **(training data)**

361 **Fig 8: Comparison of NNAR (1,1,2) [12] 24-month forecast TB cases versus actual TB cases**
362 **(test data)**

363 **Table 3: Model (NNAR (1,1,2) [12]) accuracy comparison between training and testing set**

Data	RMSE	MAE	MAPE
Training	18.56	14.58	29.89
Testing	28.65	21.95	38.86

364 **Hybrid (ARIMA-ANN) model estimation and accuracy**

365 Residuals from the optimal ARIMA (0,0,1,1,0,1,12) model were fit using an ANN model
366 and the accuracy measures and comparison of the forecast and prediction were presented in
367 Table 4, Fig 9 and Fig 10 respectively.

368 **Fig 9: Comparison of ARIMA-ANN predicted TB cases versus actual TB cases (training**
369 **data)**

370 **Fig 10: Comparison of ARIMA-ANN 24-month forecasted TB cases versus actual TB cases**
371 **(testing data)**

372 **Table 4: Hybrid ARIMA-ANN model accuracy**

Data	RMSE	MAE	MAPE
Training	19.08	15.32	42.42
Testing	27.61	19.69	32.89

373 **Comparison of predictive accuracy of the models**

374 The predictive accuracy of the models was compared using the Diebold-Mariano (DM)
375 test with the null hypothesis that the predictive accuracy of the two models compared are the
376 same. The results in table 5 show that the NNAR (1,1,2) [12] and ARIMA-ANN models present
377 significantly different predictive accuracies, similar to ARIMA (0,0,1,1,0,1,12) versus ARIMA-

378 ANN models. In general, the ARIMA-ANN model offers better predictive accuracy compared to
379 the NNAR (1,1,2) [12] and ARIMA (0,0,1,1,0,1,12) models.

380 **Table 5: Comparison of predictive accuracy**

Model	DM statistic	Loss Function Power	P-value
ARIMA (0,0,1,1,0,1,12) Vs NNAR (1,1,2) [12]	0.732	2	0.466
NNAR (1,1,2) [12] Vs ARIMA-ANN	6.260	2	<0.001
ARIMA (0,0,1,1,0,1,12) Vs ARIMA-ANN	8.732	2	<0.001

381 **Comparison of model performance in forecasting temporal trends of TB incidences**

382 The resulting ARIMA (0,0,1,1,0,1,12), NNAR (1,1,2) [12] and ARIMA-ANN models
383 were used to forecast TB cases for the next 12 months (2022). The forecast results were
384 presented in table 6. The 12-month mean forecasted TB cases was 55, 59 and 52 cases per month
385 based on the ARIMA (0,0,1,1,0,1,12), ANN (1,1,2) [12] and ARIMA-ANN respectively for
386 2022 giving a total of 657, 706 and 629 TB cases forecasted for the year 2022 from the ARIMA
387 (0,0,1,1,0,1,12), ANN (1,1,2) [12] and ARIMA-ANN models respectively.

388

389

390

391

392

393

394

395 **Table 6: TB cases 12 month point forecast comparison**

Month	ARIMA (0,0,1,1,0,1,12)	NNAR (1,1,2) [12]	ARIMA-ANN
Jan-22	57	53	44
Feb-22	61	53	45
Mar-22	62	52	60
Apr-22	56	58	43
May-22	46	52	44
Jun-22	54	52	71
Jul-22	63	72	54
Aug-22	46	52	52
Sep-22	53	56	62
Oct-22	63	103	49
Nov-22	48	57	46
Dec-22	48	46	59
Mean	55	59	52
Total	657	706	629

396 **Discussion**

397 Although all three models were able to predict TB cases among children below 15 years,
398 the hybrid ARIMA-ANN model was able to offer better predictive performance compared to
399 single ARIMA and ANN models. These findings compare with those from other studies which
400 applied either hybridized ARIMA or SARIMA in the modeling of TB incidences and other
401 infectious diseases [12, 16, 43, 44] with the overall conclusion that hybrid models have better
402 predictive performance. The majority of infectious disease data are neither purely linear nor non-
403 linear and mostly present with both linear and nonlinear properties. As such, single models are
404 not enough in modeling such kinds of data. Hybrid models are found to be most appropriate for
405 the accurate estimation of such data [45]. The use of hybridized ARIMA models has been
406 proposed in recent years and used extensively with improvements proposed over time.

407 The estimated TB prevalence in Kenya was 259 TB cases per 100,000 population in 2020
408 [4]. This translates to approximately 134,680 TB cases and with children accounting for about

409 20% (26,936) of these cases [46], the prevalence among children below 15 years was
410 approximately 121 TB cases per 100,000 population of children.

411 From the forecasted mean of 52 to 59 TB cases per month in 2022 for Homabay and
412 Turkana Counties, this study estimates that for the year 2022, the number of TB cases reported
413 would be approximately 629 to 706. However, given that these are estimated reported cases, they
414 most likely represent only about 35% of TB cases since up to 65% of pediatric TB cases are
415 potentially missed each year [2]. Taking this into account, the estimated TB cases for 2022 will
416 be approximately 1797 to 2017 for Homa Bay and Turkana Counties among children. The
417 estimated population of children below 15 years in Homabay and Turkana Counties for 2022 is
418 approximately 1,020,795 [47]. As such, the estimated TB prevalence among children in Homa
419 Bay and Turkana Counties in 2022 would be approximately 176 to 198 TB incidences per
420 100,000 population. These TB prevalence estimates for 2022 are way higher than the estimated
421 national average of 121 per 100,000 population of children below 15 years in 2020.

422 The findings of this study show that the estimated TB prevalence among children below
423 15 years is much higher compared to the estimated national average for 2020. These findings are
424 in line with the WHO newsletter that showed that the number of people developing TB and
425 dying from the disease could be much higher in 2021 and 2022 mainly because of the COVID-19
426 pandemic [48]. These findings also confirm those by Oliwa *et al.* who indicated that notification
427 data may underestimate the TB burden among children [49] while Mbithi *et al.* reported a
428 decrease in TB diagnosis in Kenya by an average of 28% in the year 2020 [50].

429 **Conclusion**

430 The hybrid ARIMA model offers better predictive accuracy and forecast performance
431 compared to single ARIMA and ANN models in modeling TB cases among children below 15
432 years in Homabay and Turkana Counties.

433 The findings in this study allude to the under-reporting of TB cases among children
434 below 15 years. As such, there is need to re-look at the TB surveillance framework data more
435 closely to understand existing gaps. There is an urgency to re-align vital resources towards the
436 National TB program to have the TB fight back on track.

437 **Limitations**

438 This study utilized data collected and reported in the TIBU system, as such, the study did
439 not have control over the quality and accuracy of the data.

440 This study utilized data between 2012 to 2021 which comprised 120 data points. Deep
441 learning algorithms such as ANNs usually demand a large amount of data to allow the algorithm
442 to effectively learn. In this study, there were 96 data points for model development and while this
443 represented 80% of the records, it might not have been sufficient enough to allow proper learning
444 of the algorithm. Furthermore, there were only 24 records used for model validation and this
445 might not have been large enough to allow for better learning by the algorithm.

446 This study combined data and analysis for Turkana and Homabay County. However,
447 these two counties might present different scenarios when it comes to pediatric TB.

448 Finally, since the study focused on modeling TB cases among children below 15 years in
449 Homa bay and Turkana counties, the findings might not be generalized to other counties of
450 Kenya.

451 **Acknowledgments**

452 We would like to thank the departments of health in Homa bay and Turkana County for
453 granting permission to access and utilize TB data in the TIBU system, and we hope that this
454 work aids in the realization of the goals of the TB program in these two counties and beyond.

455 **References**

- 456 1. Liao CM, Lin YJ, Cheng YH. Modeling the impact of control measures on tuberculosis
457 infection in senior care facilities. *Building and environment*. 2013 Jan 1; 59:66-75.
- 458 2. WHO *Global tuberculosis report*. 2018
- 459 3. WHO *Global tuberculosis report*. 2017
- 460 4. WHO *Global tuberculosis report*. 2020
- 461 5. Floyd K, Glaziou P, Zumla A, Raviglione M. The global tuberculosis epidemic and progress
462 in care, prevention, and research: an overview in year 3 of the End TB era. *The Lancet*
463 *Respiratory Medicine*. 2018 Apr 1;6(4):299-314.
- 464 6. Zumla A, Petersen E, Nyirenda T, Chakaya J. Tackling the tuberculosis epidemic in sub-
465 Saharan Africa—unique opportunities arising from the second European Developing
466 Countries Clinical Trials Partnership (EDCTP) programme 2015-2024. *International Journal*
467 *of Infectious Diseases*. 2015 Mar 1;32:46-9.
- 468 7. WHO *Global tuberculosis report*. 2016
- 469 8. Kimani E, Muhula S, Kiptai T, Orwa J, Odero T, Gachuno O. Factors influencing TB
470 treatment interruption and treatment outcomes among patients in Kiambu County, 2016-
471 2019. *PloS one*. 2021 Apr 6;16(4):e0248820.

- 472 9. Kipruto H, Mung'atu J, Ogila K, Adem A, Mwalili S, Masini E, Kibuchi E. The
473 epidemiology of tuberculosis in Kenya, a high TB/HIV burden country (2000-2013).
474 International Journal of Public Health and Epidemiology Research. 2015;1(1):2-13.
- 475 10. Houben RM, Dowdy DW, Vassall A, Cohen T, Nicol MP, Granich RM, Shea JE, Eckhoff P,
476 Dye C, Kimerling ME, White RG. TB MAC TB-HIV meeting participants. How can
477 mathematical models advance tuberculosis control in high HIV prevalence settings? Int J
478 Tuberc Lung Dis. 2014;18:509-14
- 479 11. Cao S, Wang F, Tam W, Tse LA, Kim JH, Liu J, Lu Z. A hybrid seasonal prediction model
480 for tuberculosis incidence in China. BMC medical informatics and decision making. 2013
481 Dec;13(1):1-7.
- 482 12. Azeez A, Obaromi D, Odeyemi A, Ndege J, Muntabayi R. Seasonality and trend forecasting
483 of tuberculosis prevalence data in Eastern Cape, South Africa, using a hybrid model.
484 International journal of environmental research and public health. 2016 Aug;13(8):757.
- 485 13. Li Z, Wang Z, Song H, Liu Q, He B, Shi P, Ji Y, Xu D, Wang J. Application of a hybrid
486 model in predicting the incidence of tuberculosis in a Chinese population. Infection and Drug
487 Resistance. 2019;12:1011.
- 488 14. Zeming Li, Yanning Li A. comparative study on the prediction of the BP artificial neural
489 network model and the ARIMA model in the incidence of AIDS. BMC medical informatics
490 and decision making. 2020 Dec;20(1):1-3.
- 491 15. Zhou L, Xia J, Yu L, Wang Y, Shi Y, Cai S, Nie S. Using a hybrid model to forecast the
492 prevalence of schistosomiasis in humans. International journal of environmental research and
493 public health. 2016 Apr;13(4):355.

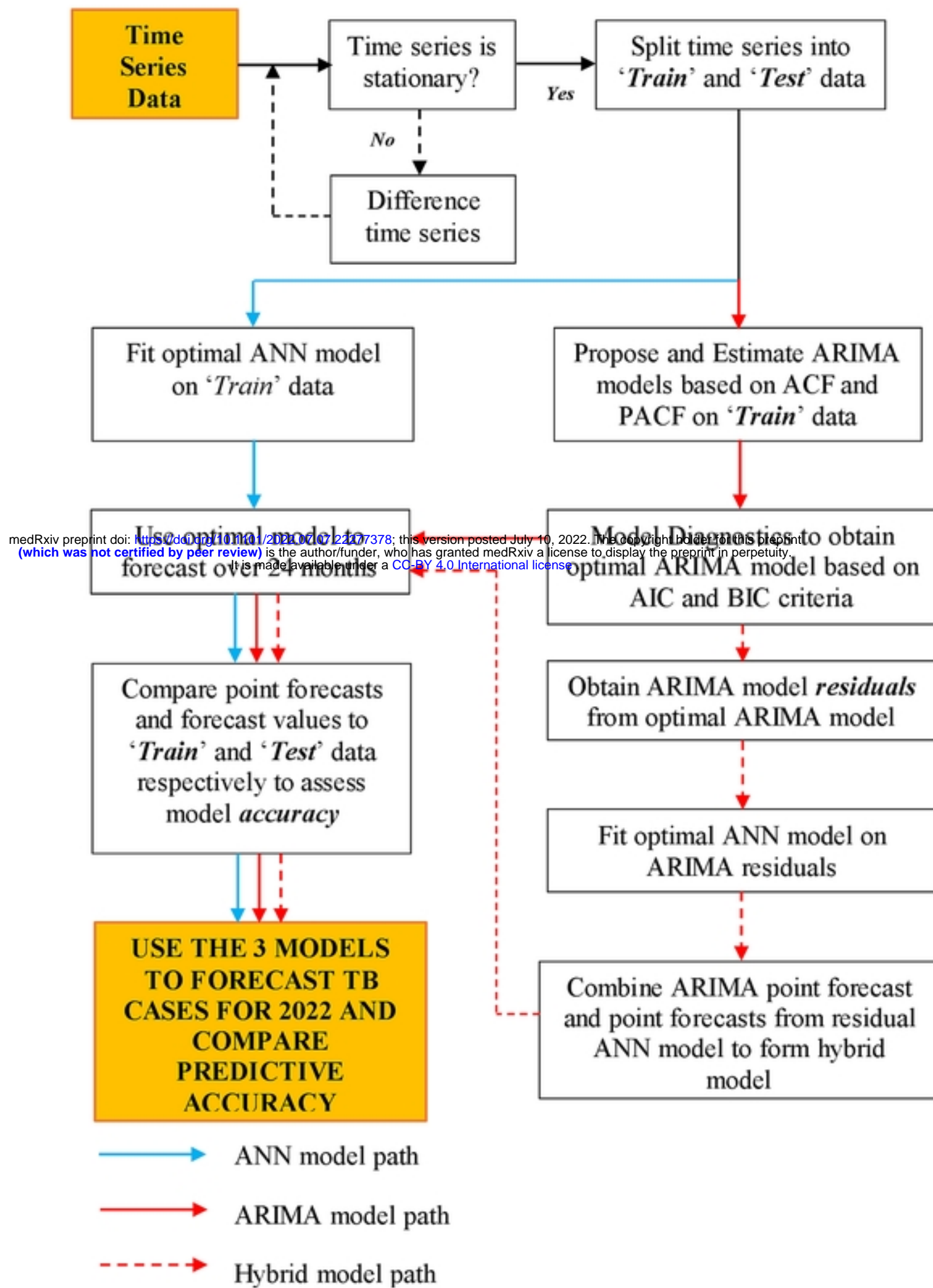
- 494 16. Yu L, Zhou L, Tan L, Jiang H, Wang Y, Wei S, Nie S. Application of a new hybrid model
495 with seasonal auto-regressive integrated moving average (ARIMA) and nonlinear
496 autoregressive neural network (NARNN) in forecasting incidence cases of HFMD in
497 Shenzhen, China. *PloS one*. 2014 Jun 3;9(6):e98241.
- 498 17. Anokye R, Acheampong E, Owusu I, Isaac Obeng E. Time series analysis of malaria in
499 Kumasi: Using ARIMA models to forecast future incidence. *Cogent social sciences*. 2018
500 Jan 1;4(1):1461544.
- 501 18. Achieng E, Otieno V, Mung'atu J. Modeling the trend of reported malaria cases in Kisumu
502 county, Kenya. *F1000Research*. 2020 Jun 12;9:600.
- 503 19. Ren H, Li J, Yuan ZA, Hu JY, Yu Y, Lu YH. The development of a combined mathematical
504 model to forecast the incidence of hepatitis E in Shanghai, China. *BMC infectious diseases*.
505 2013 Dec;13(1):1-6.
- 506 20. Zhang GP. A neural network ensemble method with jittered training data for time series
507 forecasting. *Information Sciences*. 2007 Dec 1;177(23):5329-46.
- 508 21. Yolcu U, Egrioglu E, Aladag CH. A new linear & nonlinear artificial neural network model
509 for time series forecasting. *Decision support systems*. 2013 Feb 1;54(3):1340-7.
- 510 22. Khashei M, Bijari M. A new class of hybrid models for time series forecasting. *Expert*
511 *Systems with Applications*. 2012 Mar 1;39(4):4344-57.
- 512 23. Ministry of Health, Government of Kenya, Division of Leprosy, Tuberculosis, and Lung
513 Disease. *National Monitoring and Evaluation Plan*. 2010
- 514 24. Ministry of Health, Government of Kenya, Division of Leprosy, Tuberculosis, and Lung
515 Disease. *Annual Report*. 2012

- 516 25. James G, Witten D, Hastie T, Tibshirani R. An introduction to statistical learning. New York:
517 Springer; 2013 Jun.
- 518 26. Medar R, Rajpurohit VS, Rashmi B. Impact of training and testing data splits on the accuracy
519 of time series forecasting in machine learning. In 2017 International Conference on
520 Computing, Communication, Control and Automation (ICCUBEA) 2017 Aug 17 (pp. 1-6).
521 IEEE.
- 522 27. Team RC. R: A language and environment for statistical computing. R Foundation for
523 Statistical Computing, Vienna, Austria. <http://www.R-project.org/>. 2021.
- 524 28. Bakar NA, Rosbi S. Data clustering using Autoregressive Integrated Moving Average
525 (ARIMA) model for Islamic country currency: an econometrics method for Islamic financial
526 engineering. The International Journal of Engineering and Science (IJES). 2017;6(6):22-31.
- 527 29. Lee J. Univariate time series modeling and forecasting (Box-Jenkins method), Econ 413,
528 Lecture 4. Department of Economics, University of Illinois. 2018.
- 529 30. GE P. Box, GM Jenkins, Time series analysis: forecasting and control, revised ed.
- 530 31. Shrivastav AK, Ekata D. Applicability of Box Jenkins ARIMA model in crime forecasting:
531 A case study of counterfeiting in Gujarat state. Int J Adv Res Comput Eng Technol.
532 2012;1(4):494-7.
- 533 32. Khashei M, Bijari M. An artificial neural network (p, d, q) model for time-series forecasting.
534 Expert Systems with applications. 2010 Jan 1;37(1):479-89.
- 535 33. Kihoro JM, Otieno RO, Wafula C. Seasonal time series forecasting: a comparative study of
536 Arima and ann models. African Journal of Science and Technology. 2006.

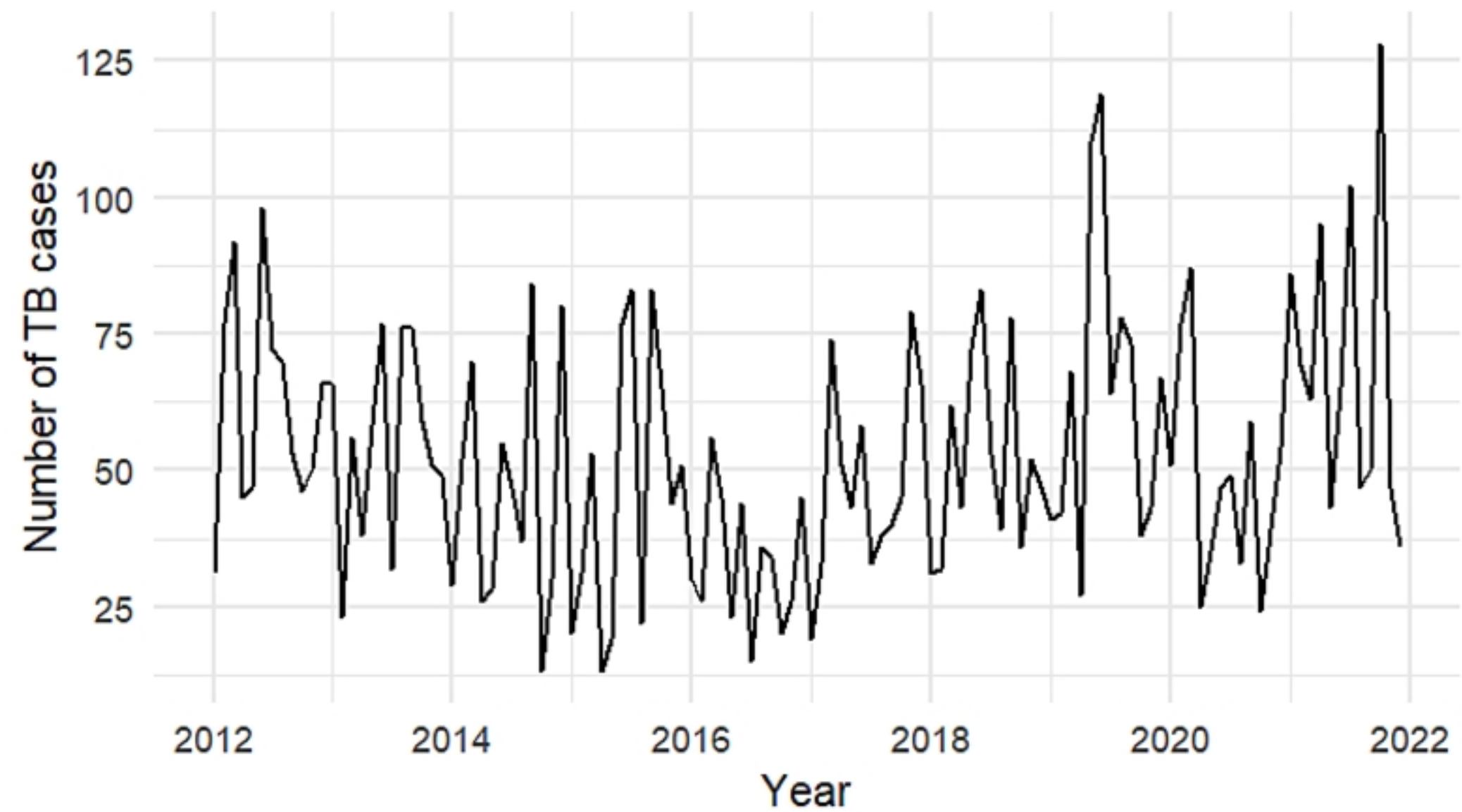
- 537 34. Larie D, An G, Cockrell RC. The Use of Artificial Neural Networks to Forecast the Behavior
538 of Agent-Based Models of Pathophysiology: An Example Utilizing an Agent-Based Model
539 of Sepsis. *Frontiers in Physiology*. 2021;12.
- 540 35. Cinar AC. Training feed-forward multi-layer perceptron artificial neural networks with a
541 tree-seed algorithm. *Arabian Journal for Science and Engineering*. 2020 Dec;45(12):10915-
- 542 36. Darji MP, Dabhi VK, Prajapati HB. Rainfall forecasting using neural network: A survey. In
543 2015 international conference on advances in computer engineering and applications 2015
544 Mar 19 (pp. 706-713). IEEE.
- 545 37. Zhang GP. Time series forecasting using a hybrid ARIMA and neural network model.
546 *Neurocomputing*. 2003 Jan 1;50:159-75.
- 547 38. Chakrabarti A, Ghosh JK. AIC, BIC, and recent advances in model selection. *Philosophy of*
548 *statistics*. 2011 Jan 1:583-605.
- 549 39. Ljung GM, Box GE. On a measure of lack of fit in time series models. *Biometrika*. 1978 Aug
550 1;65(2):297-303.
- 551 40. Hamzaçebi C. Improving artificial neural networks' performance in seasonal time series
552 forecasting. *Information Sciences*. 2008 Dec 1;178(23):4550-9.
- 553 41. Lewis C. *International and Business Forecasting Methods* Butterworths: London. 1982.
- 554 42. DIEBOLD F, MARIANO R. Comparing predictive accuracy. *journal of business and*
555 *Economics Statistics*, v. 13. 1995.
- 556 43. Nyoni SP, Nyoni T. Modeling and Forecasting TB Incidence in Bolivia Using the Multilayer
557 Perceptron Neural Network. *International Research Journal of Innovations in Engineering*
558 *and Technology*. 2021 Mar 1;5(3):301.

- 559 44. Zhou L, Zhao P, Wu D, Cheng C, Huang H. Time series model for forecasting the number of
560 new admission inpatients. *BMC medical informatics and decision making*. 2018 Dec;18(1):1-
561 1.
- 562 45. Chakraborty T, Chakraborty AK, Biswas M, Banerjee S, Bhattacharya S. Unemployment rate
563 forecasting: A hybrid approach. *Computational Economics*. 2021 Jan;57(1):183-201.
- 564 46. Okwara FN, Oyore JP, Were FN, Gwer S. Correlates of isoniazid preventive therapy failure
565 in child household contacts with infectious tuberculosis in high burden settings in Nairobi,
566 Kenya—a cohort study. *BMC infectious diseases*. 2017 Dec;17(1):1-1.
- 567 47. https://www.census.gov/datatools/demo/idb/#!/table?COUNTRY_YEAR=2022&COUNTRY
568 [_YR_ANIM=2022](https://www.census.gov/datatools/demo/idb/#!/table?COUNTRY_YEAR=2022&COUNTRY)
- 569 48. WHO Global Tuberculosis Report. 2021
- 570 49. Oliwa JN, Maina J, Ayieko P, Gathara D, Kathure IA, Masini E, Van't Hoog AH, van
571 Hensbroek MB, English M. Variability in distribution and use of tuberculosis diagnostic tests
572 in Kenya: a cross-sectional survey. *BMC infectious diseases*. 2018 Dec;18(1):1-3.
- 573 50. Mbithi I, Thekkur P, Chakaya JM, Onyango E, Owiti P, Njeri NC, Kumar AM,
574 Satyanarayana S, Shewade HD, Khogali M, Zachariah R. Assessing the real-time impact of
575 COVID-19 on TB and HIV services: the experience and response from selected health
576 facilities in Nairobi, Kenya. *Tropical Medicine and Infectious Disease*. 2021 May 10;6(2):74.
- 577 51. Makori L, Gichana H, Oyugi E, Nyale G, Ransom J. Tuberculosis in an urban hospital
578 setting: Descriptive epidemiology among patients at Kenyatta National Hospital TB clinic,
579 Nairobi, Kenya. *International Journal of Africa Nursing Sciences*. 2021 Jan 1;15:100308.

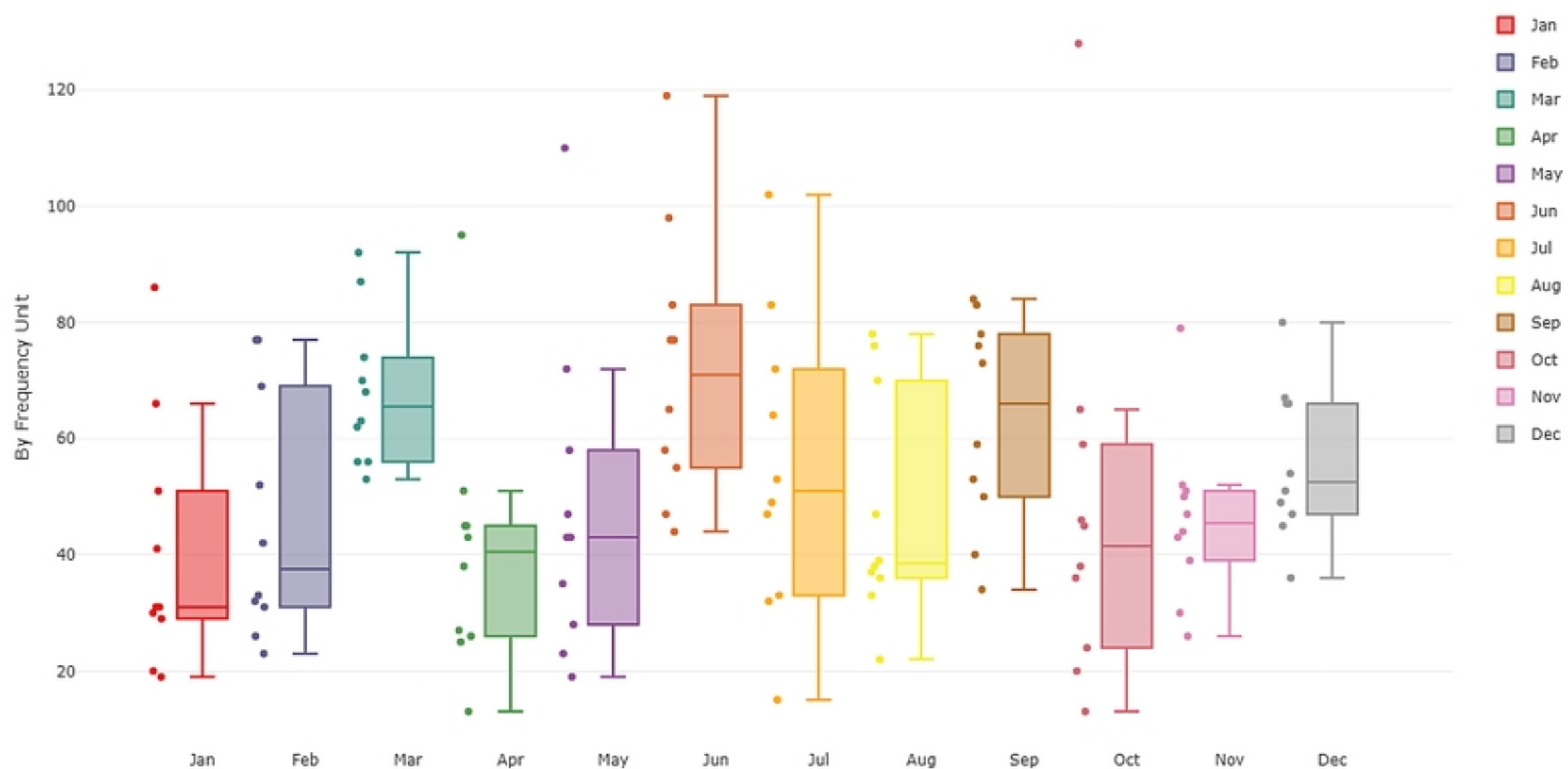
580



The Proposed methodology

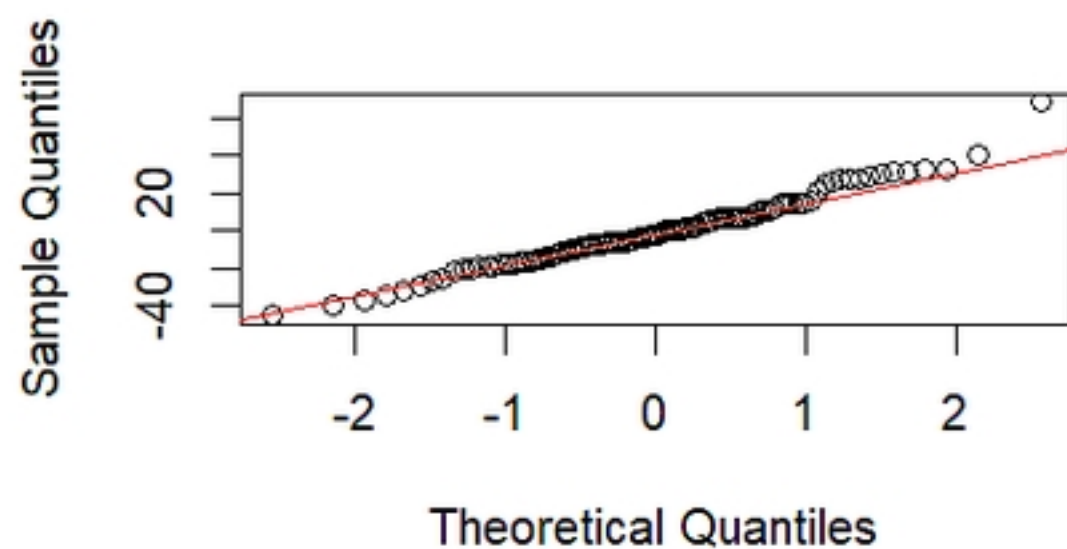
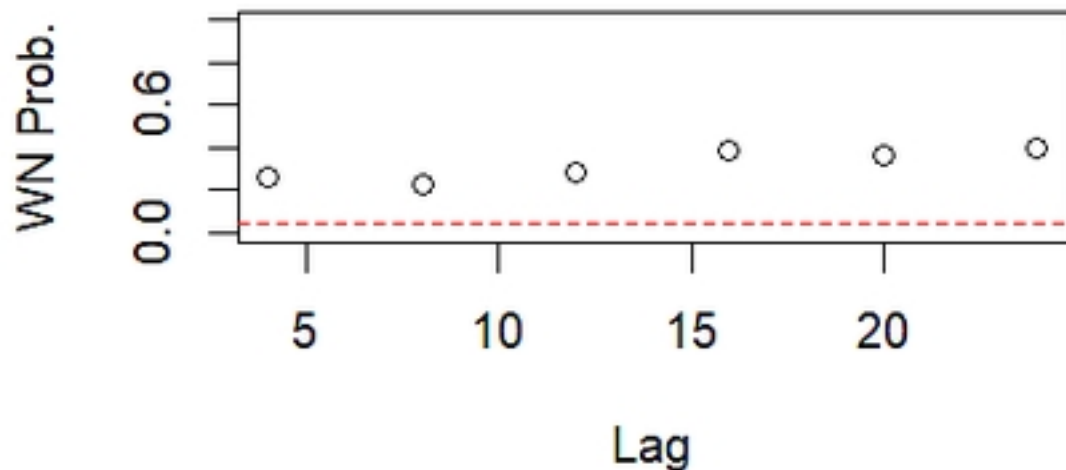
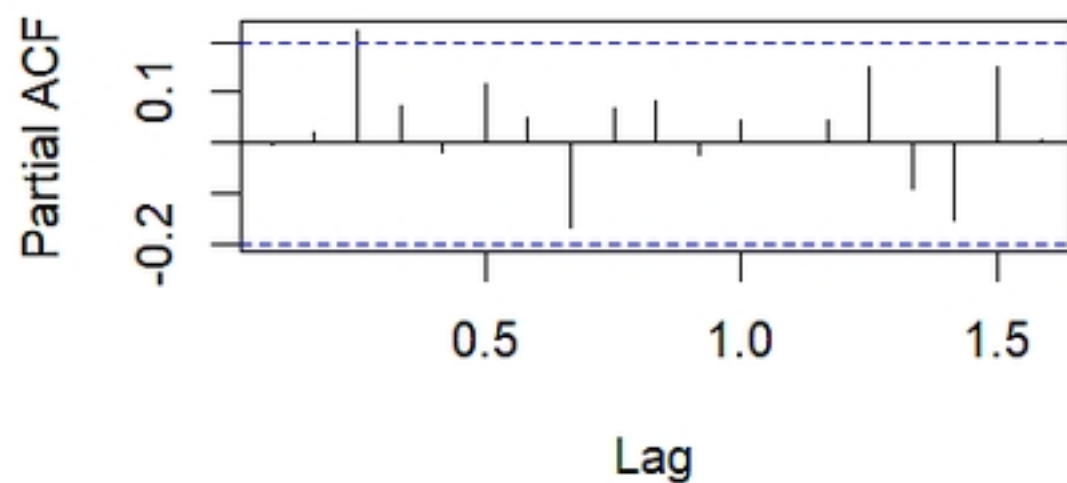
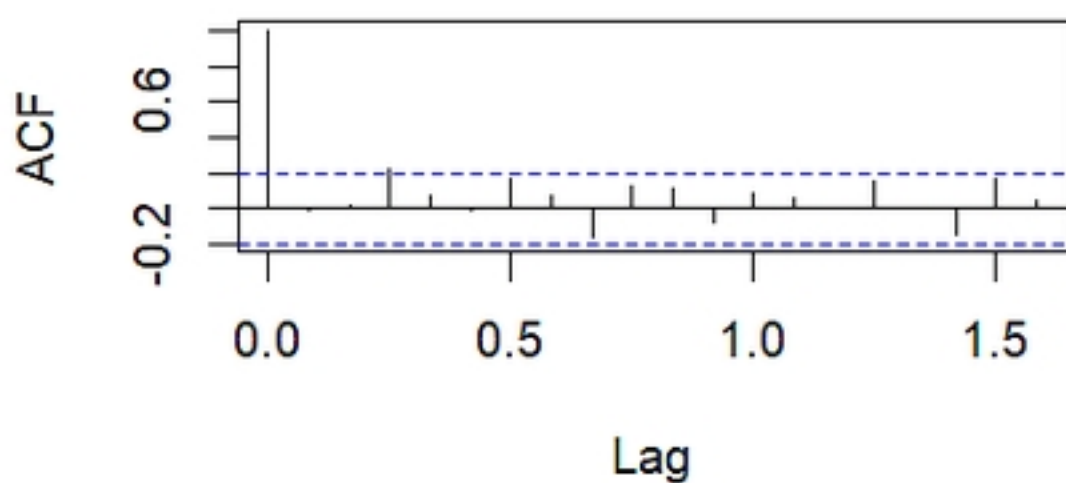


Monthly TB cases among children below 15 years from H

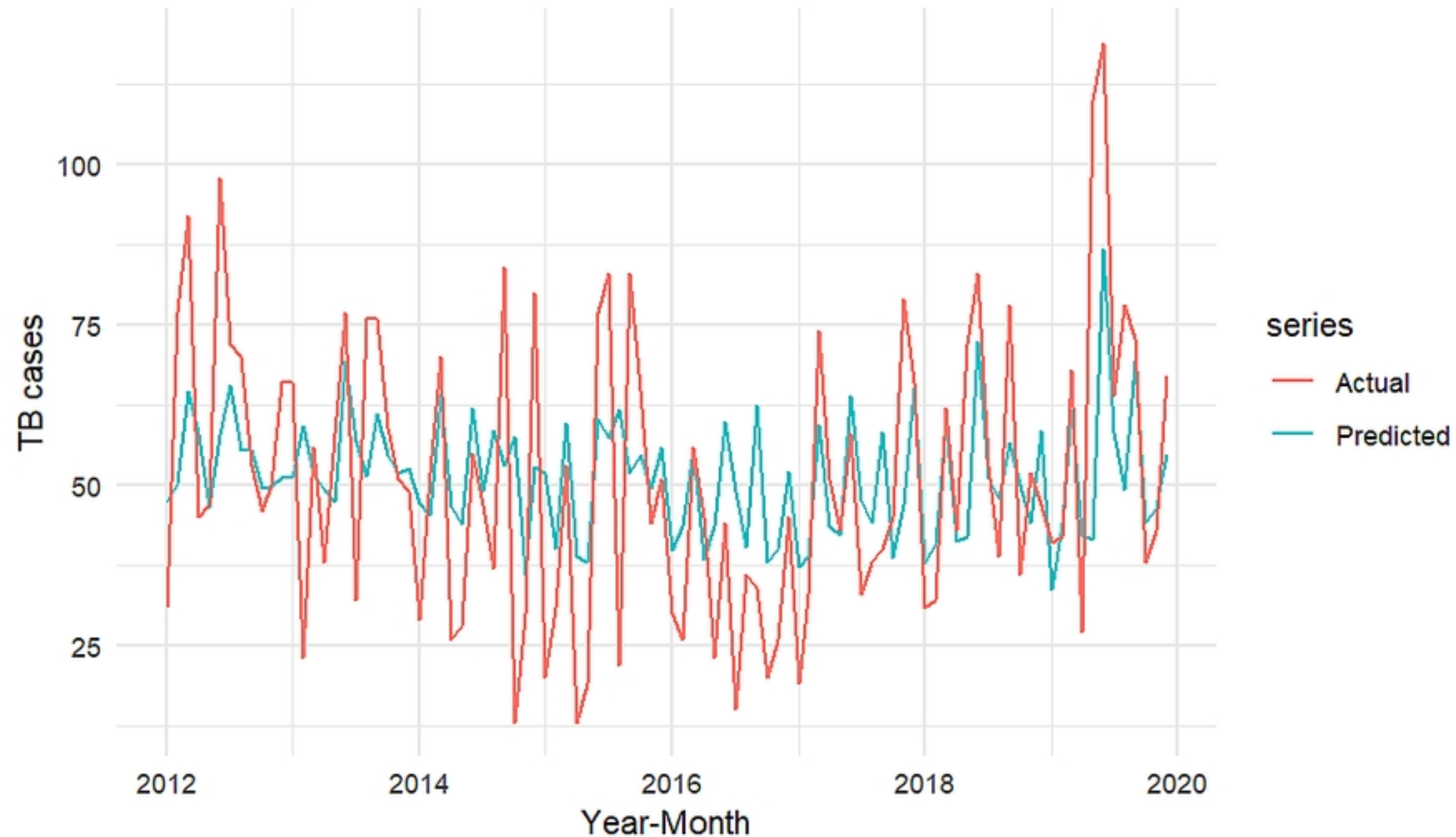


Monthly cycle plot of TB cases

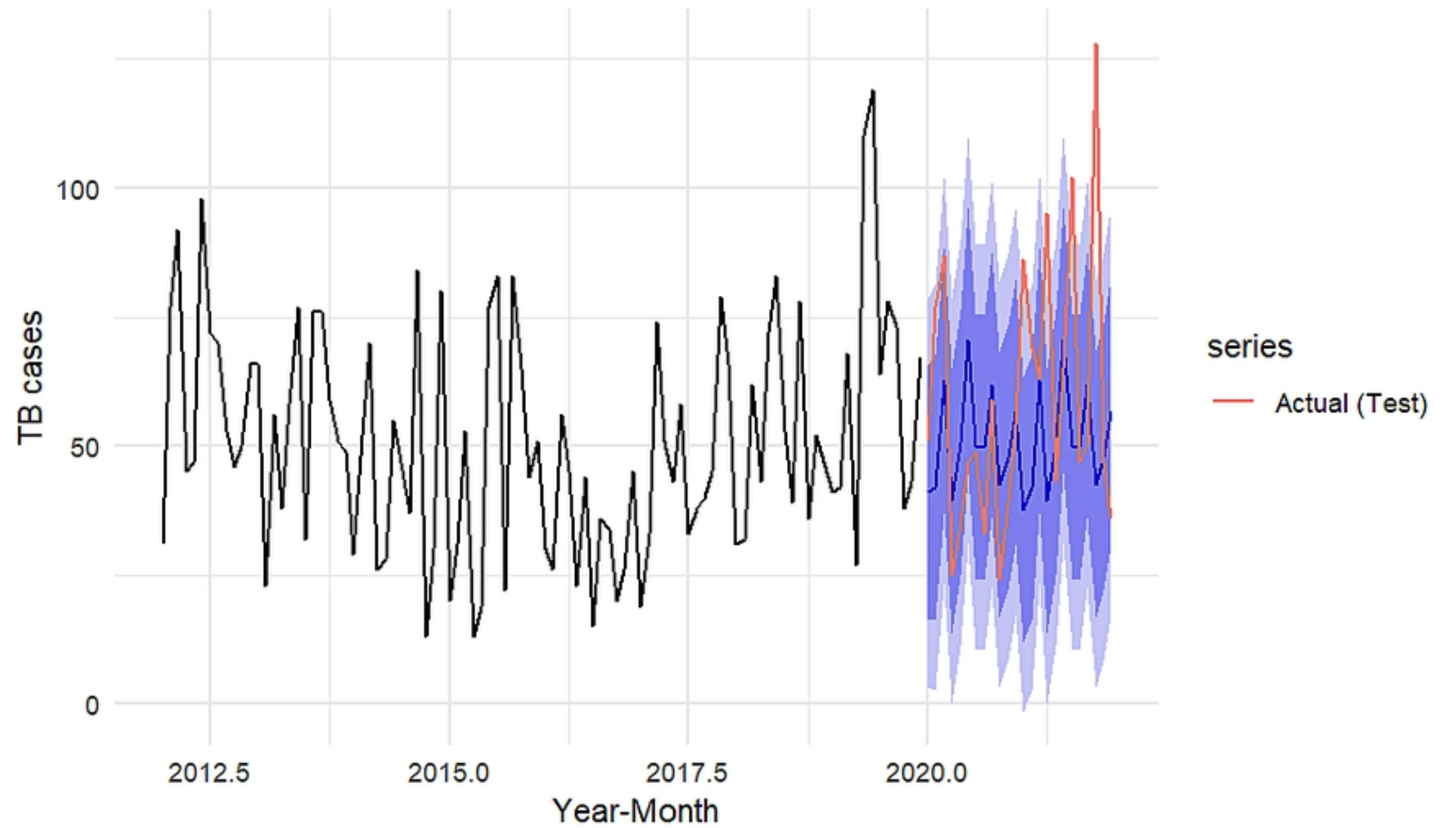
Residual Diagnostics Plots



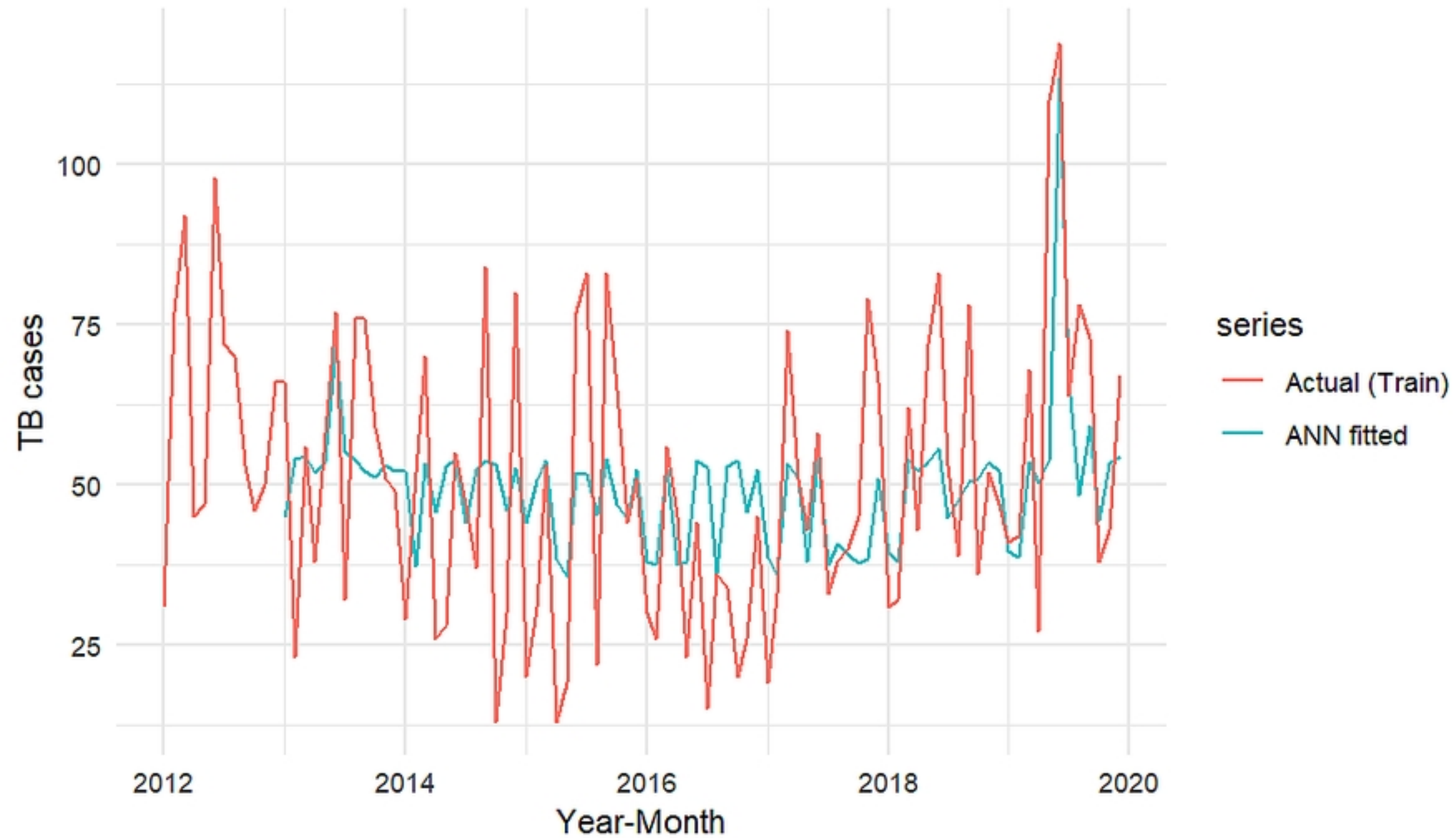
ARIMA Model residual diagnostics



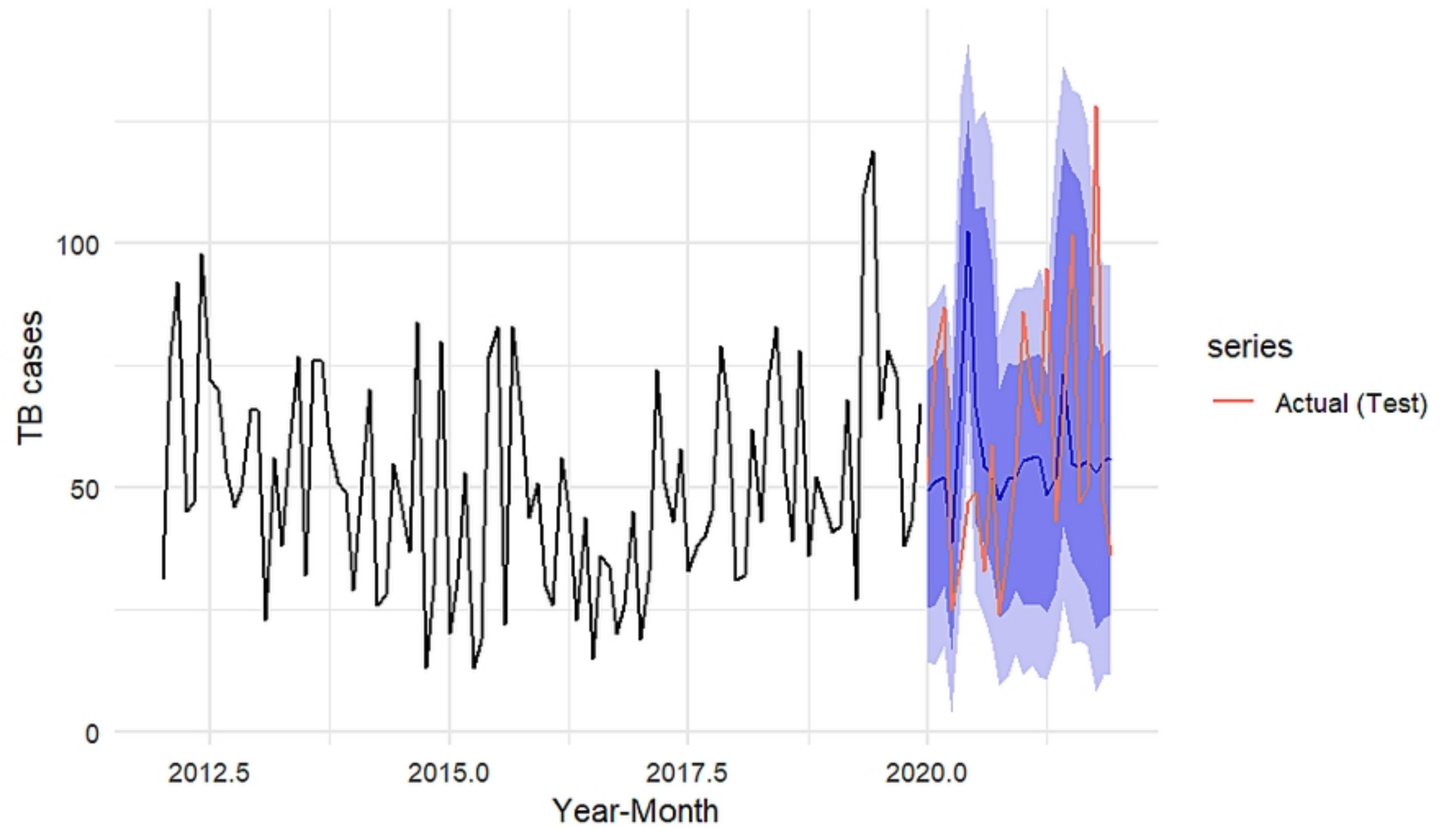
Comparison of ARIMA predicted versus actual (training data) TB



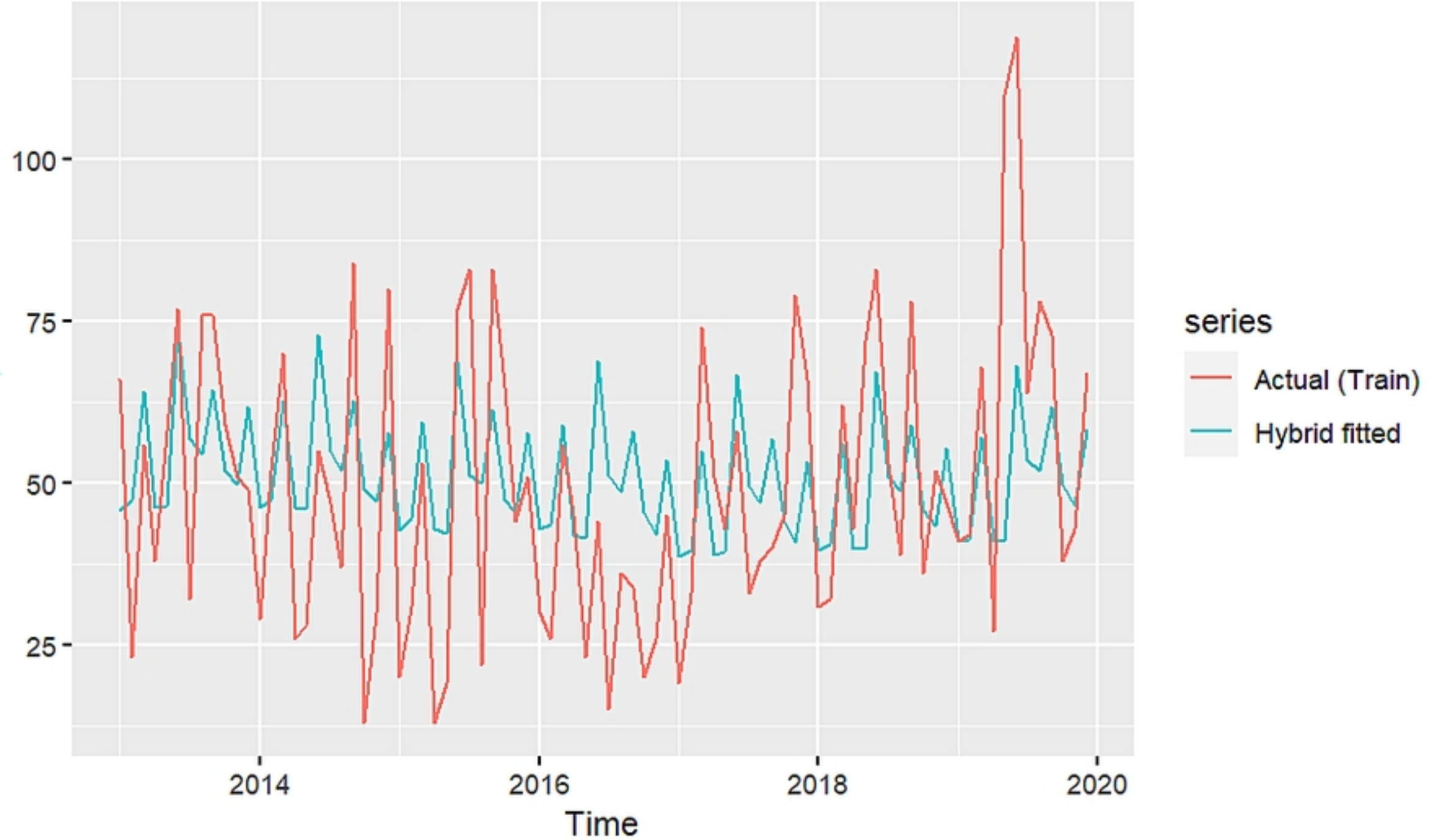
Comparison of ARIMA 24-month forecast versus actual (test data)



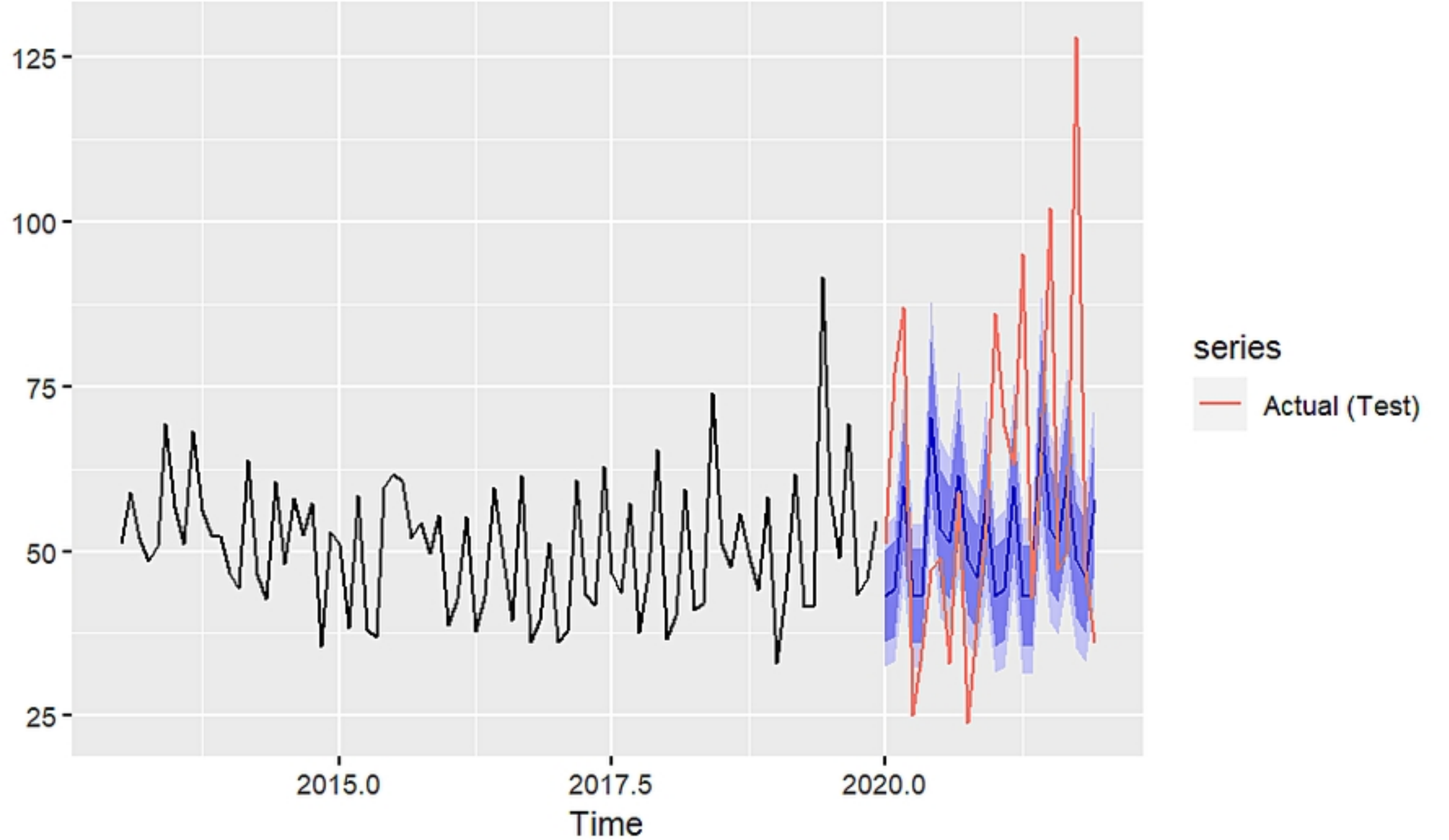
Comparison of NNAR (1,1,2) [12] predicted TB cases versus actual



Comparison of NNAR (1,1,2) [12] 24-month forecast TB cases ver



Comparison of ARIMA-ANN predicted TB cases versus actual TB



Comparison of ARIMA-ANN 24-month forecasted TB cases versus