

Task 3 - Statistical Report – Time Series Analysis of UK Live Births

1. Introduction

Understanding demographic changes is fundamental to national planning, particularly in policy areas such as healthcare, education, migration, taxation, and ageing population strategy. One of the most important demographic indicators is the annual number of live births, which directly affects future population size and age structure. A sustained decline in births can lead to long-term challenges, including labour shortages, increased dependency ratios, and pressure on pension systems, while sudden increases can strain public services such as maternity care, schools, and childcare.

This report focuses on forecasting the **number of live births in the United Kingdom** using historical records found in the *Vital Statistics in the UK* dataset. The dataset spans more than a century of demographic change, making it suitable for time-series modelling and providing an opportunity to analyse both long-term trends and short-term fluctuations.

Task Understanding

This task requires:

1. Conducting **initial exploratory analysis** of the time-series dataset
2. Decomposing the series into components such as trend and seasonal effects
3. Developing, comparing, and validating **multiple forecasting models**
4. Testing assumptions associated with each model
5. Interpreting results clearly with supporting graphs and outputs
6. Producing a final forecast and a concise summary of findings

Three forecasting approaches were used in this analysis:

- A **naïve model**, which acts as a benchmark
- An **ETS (Exponential Smoothing)** model
- An **ARIMA** model, with log transformation where appropriate

The analysis also includes a hypothesis test on UK fertility rates to complement the birth-rate forecast.

2. Data Exploration and Preparation

Before modelling, the dataset required substantial cleaning due to metadata rows, scattered formatting, and numerical columns stored as text.

2.1 Loading and Inspecting the Data

The raw dataset contained descriptive notes, footnotes, and blank rows which needed to be removed. After skipping the top metadata rows, the cleaned dataset consisted of annual records of live births for the entire UK.

```
#install Packages
install.packages(c("readxl", "dplyr", "janitor", "ggplot2", "forecast", "tseries"))

library(readxl) # reading Excel files
library(dplyr) # data wrangling
library(janitor) # clean_names()
library(ggplot2) # plotting
library(forecast) # time series models (ETS, ARIMA, etc.)
library(tseries) # ADF test for stationarity

# 1. LOAD RAW DATA

# Load data
births_raw <- read_excel("vital statistics in the UK.xlsx", sheet = "Birth")

# quick look
str(births_raw)
names(births_raw)
head(births_raw)

# Ignore the top 5 metadata rows
births_raw <- read_excel(
  "vital statistics in the UK.xlsx",
  sheet = "Birth",
  skip = 5      # skip note rows so row 6 becomes header
)

# Inspect structure
str(births_raw)
head(births_raw)

> births_raw <- read_excel("vital statistics in the UK.xlsx", sheet = "Birth")
New names:
 * `..` -> `...2`
 * `..` -> `...3`
 * `..` -> `...4`
 * `..` -> `...5`
 * `..` -> `...6`
 * `..` -> `...7`
 * `..` -> `...8`
 * `..` -> `...9`
 * `..` -> `...10`
 * `..` -> `...11`
 * `..` -> `...12`
 * `..` -> `...13`
 * `..` -> `...14`
 * `..` -> `...15`
 * `..` -> `...16`
 * `..` -> `...17`
 * `..` -> `...18`
 * `..` -> `...19`
 * `..` -> `...20`
 * `..` -> `...21`
 * `..` -> `...22`
 * `..` -> `...23`
 * `..` -> `...24`
 * `..` -> `...25`
 * `..` -> `...26`
 * `..` -> `...27`
 * `..` -> `...28`
 * `..` -> `...29`
 * `..` -> `...30`
 * `..` -> `...31`
 * `..` -> `...32`
 * `..` -> `...33`
 * `..` -> `...34`
 * `..` -> `...35`
 * `..` -> `...36`
 * `..` -> `...37`
 * `..` -> `...38`
 * `..` -> `...39`
```

```

> # quick look
> str(births_raw)
tibble [177 x 39] (s3:tbl_df/tbl/data.frame)
$ Annual data: Births (numbers, rates and standardised mean age of mother): chr [1:177] "[note 1][note 10][note 11][note 12][note 13][note 14][note 15][note 16][note 17][note 18][note 19][note 20]" "This worksheet contains one table. The notes can be found on the Notes sheet." "United Kingdom and constituent countries" "Sources: Office for National Statistics, National Records of Scotland, Northern Ireland Statistics and Research Agency" ...
$ ...2 : chr [1:177] NA NA NA NA ...
$ ...3 : chr [1:177] NA NA NA NA ...
$ ...4 : chr [1:177] NA NA NA NA ...
$ ...5 : chr [1:177] NA NA NA NA ...
$ ...6 : chr [1:177] NA NA NA NA ...
$ ...7 : chr [1:177] NA NA NA NA ...
$ ...8 : chr [1:177] NA NA NA NA ...
$ ...9 : chr [1:177] NA NA NA NA ...
$ ...10 : chr [1:177] NA NA NA NA ...
$ ...11 : chr [1:177] NA NA NA NA ...
$ ...12 : chr [1:177] NA NA NA NA ...
$ ...13 : chr [1:177] NA NA NA NA ...
$ ...14 : logi [1:177] NA NA NA NA NA ...
$ ...15 : logi [1:177] NA NA NA NA NA ...
$ ...16 : logi [1:177] NA NA NA NA NA ...
$ ...17 : logi [1:177] NA NA NA NA NA ...
$ ...18 : num [1:177] NA NA NA NA NA NA NA NA NA ...
$ ...19 : logi [1:177] NA NA NA NA NA NA ...
$ ...20 : logi [1:177] NA NA NA NA NA NA ...
$ ...21 : logi [1:177] NA NA NA NA NA NA ...
$ ...22 : logi [1:177] NA NA NA NA NA NA ...
$ ...23 : logi [1:177] NA NA NA NA NA NA ...
$ ...24 : logi [1:177] NA NA NA NA NA NA ...
$ ...25 : logi [1:177] NA NA NA NA NA NA ...
$ ...26 : logi [1:177] NA NA NA NA NA NA ...
$ ...27 : logi [1:177] NA NA NA NA NA NA ...
$ ...28 : logi [1:177] NA NA NA NA NA NA ...
$ ...29 : logi [1:177] NA NA NA NA NA NA ...
$ ...30 : logi [1:177] NA NA NA NA NA NA ...
$ ...31 : logi [1:177] NA NA NA NA NA NA ...
$ ...32 : logi [1:177] NA NA NA NA NA NA ...
$ ...33 : logi [1:177] NA NA NA NA NA NA ...
$ ...34 : logi [1:177] NA NA NA NA NA NA ...
$ ...35 : logi [1:177] NA NA NA NA NA NA ...
$ ...36 : logi [1:177] NA NA NA NA NA NA ...
$ ...37 : logi [1:177] NA NA NA NA NA NA ...
$ ...38 : logi [1:177] NA NA NA NA NA NA ...
$ ...39 : logi [1:177] NA NA NA NA NA NA ...

> names(births_raw)
[1] "Annual data: Births (numbers, rates and standardised mean age of mother)"
[2] "...2"
[3] "...3"
[4] "...4"
[5] "...5"
[6] "...6"
[7] "...7"
[8] "...8"
[9] "...9"
[10] "...10"
[11] "...11"
[12] "...12"
[13] "...13"
[14] "...14"
[15] "...15"
[16] "...16"
[17] "...17"
[18] "...18"
[19] "...19"
[20] "...20"
[21] "...21"
[22] "...22"
[23] "...23"
[24] "...24"
[25] "...25"
[26] "...26"
[27] "...27"
[28] "...28"
[29] "...29"
[30] "...30"
[31] "...31"
[32] "...32"
[33] "...33"
[34] "...34"
[35] "...35"
[36] "...36"
[37] "...37"
[38] "...38"
[39] "...39"

```

```

> head(births_raw)
# A tibble: 6 × 39
  Annual data: Births (nu...` ...2 ...3 ...4 ...5 ...6 ...7 ...8 ...9 ...10 ...11 ...12
  <chr>           <chr> <chr> <chr> <chr> <chr> <chr> <chr> <chr> <chr> <chr>
1 [note 1][note 10][note 1... NA   NA   NA   NA   NA   NA   NA   NA   NA   NA
2 This worksheet contains ... NA   NA   NA   NA   NA   NA   NA   NA   NA   NA
3 United Kingdom and const... NA   NA   NA   NA   NA   NA   NA   NA   NA   NA
4 Sources: Office for Nati... NA   NA   NA   NA   NA   NA   NA   NA   NA   NA
5 Year               Numb... Numb... Numb... Numb... Numb... Numb... Tota... Tota... Tota... Tota...
6 2021              6946... 6248... 5959... 28781 47786 22071 1.53 1.55 1.55 1.5 1.31
# i abbreviated name:
#   `Annual data: Births (numbers, rates and standardised mean age of mother)`
# i 27 more variables: ...13 <chr>, ...14 <lgl>, ...15 <lgl>, ...16 <lgl>, ...17 <lgl>,
#   ...18 <dbl>, ...19 <lgl>, ...20 <lgl>, ...21 <lgl>, ...22 <lgl>, ...23 <lgl>,
#   ...24 <lgl>, ...25 <lgl>, ...26 <lgl>, ...27 <lgl>, ...28 <lgl>, ...29 <lgl>,
#   ...30 <lgl>, ...31 <lgl>, ...32 <lgl>, ...33 <lgl>, ...34 <lgl>, ...35 <lgl>,
#   ...36 <lgl>, ...37 <lgl>, ...38 <lgl>, ...39 <lgl>
> # Ignore the top 5 metadata rows
> births_raw <- read_excel(
+   "vital statistics in the uk.xlsx",
+   sheet = "Birth",
+   skip = 5      # skip note rows so row 6 becomes header
+ )
New names:
• `..` -> `...14`
• `..` -> `...15`
• `..` -> `...16`
• `..` -> `...17`
• `..` -> `...18`

> # Inspect structure
> str(births_raw)
tibble [172 × 18] (s3:tbl_df/tbl/data.frame)
$ Year          : chr [1:172] "2021" "2020" "2019" "2018" ...
$ Number of live births: United Kingdom : chr [1:172] "694685" "681560" "712680" "731213" ...
$ Number of live births: England and Wales: num [1:172] 624828 613936 640370 657076 679106 ...
$ Number of live births: England       : chr [1:172] "595948" "585195" "610505" "625651" ...
$ Number of live births: Wales        : chr [1:172] "28781" "28638" "29704" "31274" ...
$ Number of live births: Scotland     : chr [1:172] "47786" "46809" "49863" "51308" ...
$ Number of live births: Northern Ireland: chr [1:172] "22071" "20815" "22447" "22829" ...
$ Total fertility rate: United Kingdom: chr [1:172] "1.53" "1.56" "1.63" "1.68" ...
$ Total fertility rate: England and Wales: chr [1:172] "1.55" "1.58" "1.65" "1.7" ...
$ Total fertility rate: England       : chr [1:172] "1.55" "1.59" "1.66" "1.7" ...
$ Total fertility rate: Wales        : chr [1:172] "1.5" "1.47" "1.54" "1.63" ...
$ Total fertility rate: Scotland     : chr [1:172] "1.31" "1.29" "1.37" "1.42" ...
$ Total fertility rate: Northern Ireland: chr [1:172] "1.81" "1.71" "1.82" "1.85" ...
$ ...14          : logi [1:172] NA NA NA NA NA NA ...
$ ...15          : logi [1:172] NA NA NA NA NA NA ...
$ ...16          : logi [1:172] NA NA NA NA NA NA ...
$ ...17          : logi [1:172] NA NA NA NA NA NA ...
$ ...18          : num [1:172] NA NA NA NA NA NA 1 NA NA NA ...

> head(births_raw)
# A tibble: 6 × 18
  Year `Number of live births: United Kingdom` `Number of live births:...` `Number of live birth...` 
  <chr> <chr> <dbl> <chr> <chr>
1 2021 694685 624828 595948 28781
2 2020 681560 613936 585195 28638
3 2019 712680 640370 610505 29704
4 2018 731213 657076 625651 31274
5 2017 755042 679106 646794 32176
6 2016 774835 696271 663157 32936
# i abbreviated names: `Number of live births: England and Wales`, `Number of live births: England`,
#   `Number of live births: Wales`
# i 13 more variables: `Number of live births: Scotland` <chr>, `Number of live births: Northern Ireland` <chr>,
#   `Total fertility rate: United Kingdom` <chr>, `Total fertility rate: England and Wales` <chr>,
#   `Total fertility rate: England` <chr>, `Total fertility rate: Wales` <chr>,
#   `Total fertility rate: Scotland` <chr>, `Total fertility rate: Northern Ireland` <chr>, ...14 <lgl>, ...15 <lgl>,
#   ...16 <lgl>, ...17 <lgl>, ...18 <dbl>

```

2.2 Cleaning Column Names and Removing Redundant Columns

Column names were standardised using `clean_names()`, and extraneous columns (typically footnote markers labelled “x13”, “x14”, etc.) were removed.

```

# 2. CLEAN COLUMN NAMES

# Turn long, messy names into snake_case
births <- births_raw %>%janitor::clean_names()

# Look at the new names
names(births)

# 3. REMOVE JUNK COLUMNS

# Remove any column whose name starts with "x"
births <- births %>%select(-starts_with("x"))

str(births)

> # 2. CLEAN COLUMN NAMES
>
> # Turn long, messy names into snake_case
> births <- births_raw %>%janitor::clean_names()
>
> # Look at the new names
> names(births)
[1] "year"                               "number_of_live_births_united_kingdom"
[3] "number_of_live_births_england_and_wales" "number_of_live_births_england"
[5] "number_of_live_births_wales"           "number_of_live_births_scotland"
[7] "number_of_live_births_northern_ireland" "total_fertility_rate_united_kingdom"
[9] "total_fertility_rate_england_and_wales" "total_fertility_rate_england"
[11] "total_fertility_rate_wales"           "total_fertility_rate_scotland"
[13] "total_fertility_rate_northern_ireland" "x14"
[15] "x15"                                 "x16"
[17] "x17"                                 "x18"
> # 3. REMOVE JUNK COLUMNS
>
> # Remove any column whose name starts with "x"
>
> births <- births %>%select(-starts_with("x"))
>
> str(births)
tibble [172 x 13] (s3:tbl_df/tbl/data.frame)
$ year                           : chr [1:172] "2021" "2020" "2019" "2018" ...
$ number_of_live_births_united_kingdom : chr [1:172] "694685" "681560" "712680" "731213" ...
$ number_of_live_births_england_and_wales: num [1:172] 624828 613936 640370 657076 679106 ...
$ number_of_live_births_england        : chr [1:172] "595948" "585195" "610505" "625651" ...
$ number_of_live_births_wales         : chr [1:172] "28781" "28638" "29704" "31274" ...
$ number_of_live_births_scotland      : chr [1:172] "47786" "46809" "49863" "51308" ...
$ number_of_live_births_northern_ireland: chr [1:172] "22071" "20815" "22447" "22829" ...
$ total_fertility_rate_united_kingdom : chr [1:172] "1.53" "1.56" "1.63" "1.68" ...
$ total_fertility_rate_england_and_wales: chr [1:172] "1.55" "1.58" "1.65" "1.7" ...
$ total_fertility_rate_england        : chr [1:172] "1.55" "1.59" "1.66" "1.7" ...
$ total_fertility_rate_wales         : chr [1:172] "1.5" "1.47" "1.54" "1.63" ...
$ total_fertility_rate_scotland      : chr [1:172] "1.31" "1.29" "1.37" "1.42" ...
$ total_fertility_rate_northern_ireland: chr [1:172] "1.81" "1.71" "1.82" "1.85" ...

```

2.3 Converting Text Values to Numeric

Many numeric values (birth counts, fertility rates) appeared as character strings due to embedded commas or symbols. These were cleaned with regex and converted to numeric values.

```

# 4. FIX DATA TYPES (NUMERIC VS TEXT)

# For each character column:
# 1. Remove everything that is not a digit or decimal point.
# 2. Convert the cleaned string to numeric.

births <- births %>%
  mutate(
    across(
      where(is.character),
      ~ gsub("[^0-9\\.]", "", .) # keep only digits and decimal points
    )
  ) %>%
  mutate(
    across(
      where(is.character),
      ~ as.numeric(.)
    )
  )

# Check structure again
str(births)

# Check how many NAs are present in each column
sapply(births, function(x) sum(is.na(x)))

> # 4. FIX DATA TYPES (NUMERIC VS TEXT)
>
> # For each character column:
> # 1. Remove everything that is not a digit or decimal point.
> # 2. Convert the cleaned string to numeric.
>
> births <- births %>%
+   mutate(
+     across(
+       where(is.character),
+       ~ gsub("[^0-9\\.]", "", .) # keep only digits and decimal points
+     )
+   ) %>%
+   mutate(
+     across(
+       where(is.character),
+       ~ as.numeric(.)
+     )
+   )
>
> # Check structure again
> str(births)
tibble [172 x 13] (s3:tbl_df/tbl/data.frame)
$ year                  : num [1:172] 2021 2020 2019 2018 2017 ...
$ number_of_live_births_united_kingdom : num [1:172] 694685 681560 712680 731213 755042 ...
$ number_of_live_births_england_and_wales: num [1:172] 624828 613936 640370 657076 679106 ...
$ number_of_live_births_england        : num [1:172] 595948 585195 610505 625651 646794 ...
$ number_of_live_births_wales         : num [1:172] 28781 28638 29704 31274 32176 ...
$ number_of_live_births_scotland     : num [1:172] 47786 46809 49863 51308 52861 ...
$ number_of_live_births_northern_ireland: num [1:172] 22071 20815 22447 22829 23075 ...
$ total_fertility_rate_united_kingdom: num [1:172] 1.53 1.56 1.63 1.68 1.74 1.79 1.8 1.82 1.83 1.92 ...
$ total_fertility_rate_england_and_wales: num [1:172] 1.55 1.58 1.65 1.7 1.76 1.81 1.82 1.83 1.85 1.94 ...
$ total_fertility_rate_england        : num [1:172] 1.55 1.59 1.66 1.7 1.76 1.81 1.82 1.83 1.85 1.94 ...
$ total_fertility_rate_wales         : num [1:172] 1.5 1.47 1.54 1.63 1.69 1.74 1.77 1.78 1.8 1.88 ...
$ total_fertility_rate_scotland     : num [1:172] 1.31 1.29 1.37 1.42 1.47 1.52 1.56 1.62 1.61 1.67 ...
$ total_fertility_rate_northern_ireland: num [1:172] 1.81 1.71 1.82 1.85 1.87 1.95 1.96 1.97 1.96 2.03 ...
>
> # Check how many NAs are present in each column
> sapply(births, function(x) sum(is.na(x)))
      year  number_of_live_births_united_kingdom
0          0                   37
number_of_live_births_england_and_wales
0          0                   79
      number_of_live_births_wales
79          5
number_of_live_births_northern_ireland
37          110
total_fertility_rate_england_and_wales
88          132
      total_fertility_rate_wales
132          121
total_fertility_rate_northern_ireland
124

```

2.4 Constructing the Final Working Dataset

The final cleaned dataset contained:

- **year**
- **uk_births**

- **uk_fertility**

Rows missing birth values were removed.

```
# 5. BUILD A CLEAN WORKING DATAFRAME
# select the focused columns
births_clean <- births %>%
  select(
    year,
    uk_births = number_of_live_births_united_kingdom,
    uk_fertility = total_fertility_rate_united_kingdom
  ) %>%
  arrange(year)

# Show the first and last few rows
head(births_clean)
tail(births_clean)
summary(births_clean)

# If the very early years have missing UK births, drop those
births_clean <- births_clean %>%filter(!is.na(uk_births))

summary(births_clean)

> # 5. BUILD A CLEAN WORKING DATAFRAME
>
> # Select the focused columns
>
> births_clean <- births %>%
+   select(
+     year,
+     uk_births = number_of_live_births_united_kingdom,
+     uk_fertility = total_fertility_rate_united_kingdom
+   ) %>%
+   arrange(year)
>
> # Show the first and last few rows
> head(births_clean)
# A tibble: 6 × 3
  year uk_births uk_fertility
  <dbl>    <dbl>      <dbl>
1 1850       NA        NA
2 1851       NA        NA
3 1852       NA        NA
4 1853       NA        NA
5 1854       NA        NA
6 1855       NA        NA
> tail(births_clean)
# A tibble: 6 × 3
  year uk_births uk_fertility
  <dbl>    <dbl>      <dbl>
1 2016    774835      1.79
2 2017    755042      1.74
3 2018    731213      1.68
4 2019    712680      1.63
5 2020    681560      1.56
6 2021    694685      1.53
> summary(births_clean)
  year      uk_births      uk_fertility
Min.   :1850   Min.   : 657038   Min.   :1.530
1st Qu.:1893   1st Qu.: 736356   1st Qu.:1.740
Median :1936   Median : 796645   Median :1.810
Mean   :1936   Mean   : 855708   Mean   :1.965
3rd Qu.:1978   3rd Qu.: 984874   3rd Qu.:1.920
Max.   :2021   Max.   :1126849   Max.   :2.950
NA's    :37      NA's    :110
>
> # If the very early years have missing UK births, drop those
> births_clean <- births_clean %>%filter(!is.na(uk_births))
>
> summary(births_clean)
  year      uk_births      uk_fertility
Min.   :1887   Min.   : 657038   Min.   :1.530
1st Qu.:1920   1st Qu.: 736356   1st Qu.:1.740
Median :1954   Median : 796645   Median :1.810
Mean   :1954   Mean   : 855708   Mean   :1.965
3rd Qu.:1988   3rd Qu.: 984874   3rd Qu.:1.920
Max.   :2021   Max.   :1126849   Max.   :2.950
NA's    :73      NA's    :73
```

3. Initial Exploration of the Time Series

3.1 Time Series Plot

A simple line plot of births over time revealed several well-known demographic patterns:

- High birth rates in the late 19th and early 20th centuries
- Sharp declines during World War I and II
- A signature "baby boom" peak around the 1950s–1960s
- A long decline from the 1970s onward
- A small rise in the 2000s
- Another steady decline from 2013 onward

```
# 6. CREATE TIME-SERIES OBJECT

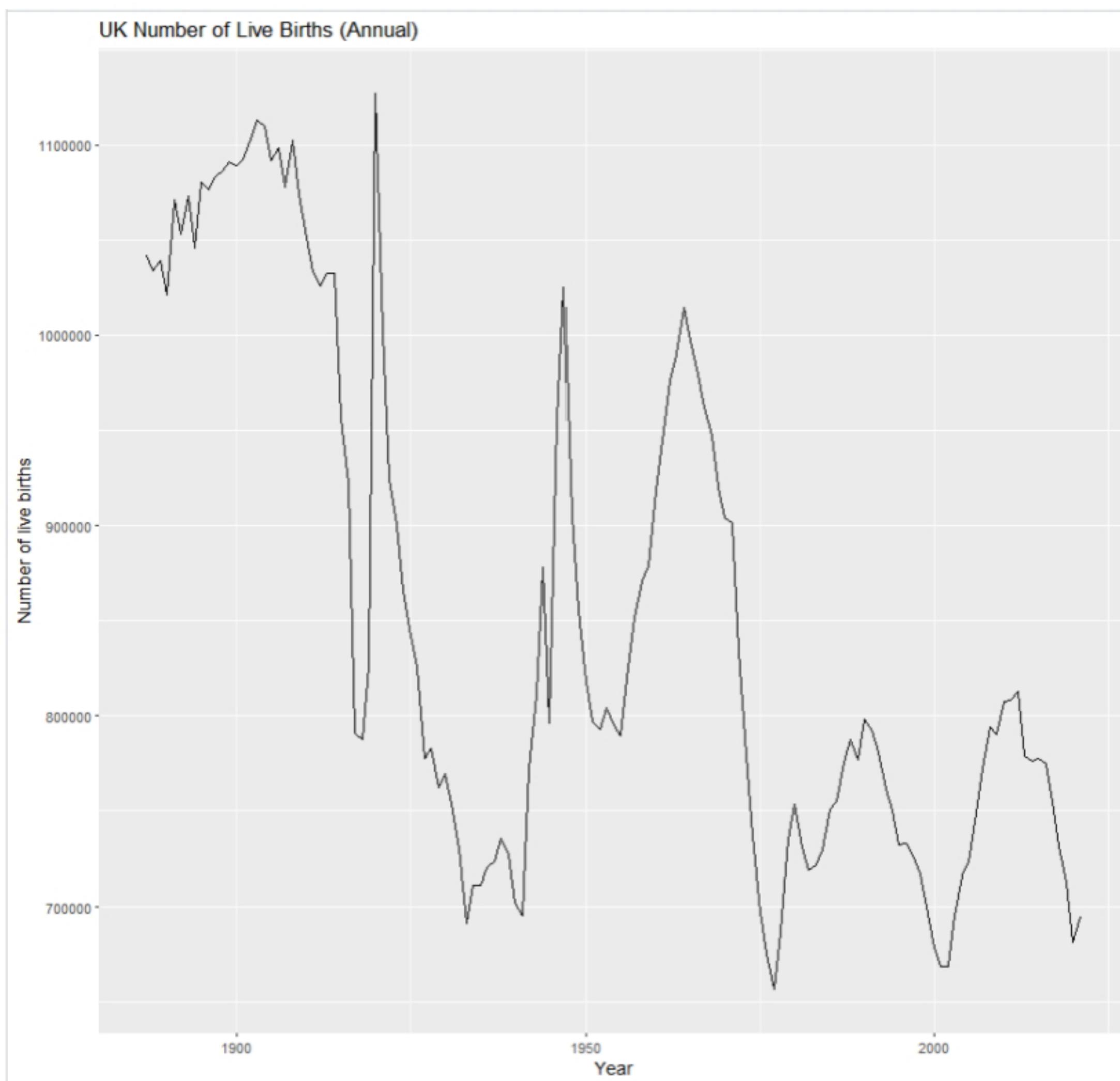
freq <- 1
start_year <- min(births_clean$year)

uk_births_ts <- ts(
  births_clean$uk_births,
  start      = c(start_year),
  frequency = freq
)
uk_births_ts

# 7. EXPLORATORY DATA ANALYSIS (EDA)

# 7.1 Time-series plot of UK births
autoplot(uk_births_ts) +
  labs(
    title = "UK Number of Live Births (Annual)",
    x = "Year",
    y = "Number of live births"
  )

> # 6. CREATE TIME-SERIES OBJECT
>
>
> freq <- 1
> start_year <- min(births_clean$year)
>
> uk_births_ts <- ts(
+   births_clean$uk_births,
+   start      = c(start_year),
+   frequency = freq
+ )
>
> uk_births_ts
Time Series:
Start = 1887
End = 2021
Frequency = 1
[1] 1041937 1034144 1039166 1021361 1071382 1053205 1073011 1045631 1080527 1076812 1082889 1086212 1091106 1089487 1092781
[16] 1103483 1113086 1109542 1092108 1098475 1077851 1102345 1073781 1051240 1033395 1025828 1032286 1032734 956877 922085
[31] 790736 787427 826202 1126849 1001725 924740 900130 865329 842405 825174 777520 783052 761963 769239 749974
[46] 730079 691560 711843 711426 720129 723779 735573 726632 701875 695726 771851 810524 878298 795868 955266
[61] 1025427 905182 855298 818421 796645 792917 804269 794769 789315 825137 851466 870497 878561 918286 944365
[76] 975635 990160 1014672 997275 979587 961800 947231 920256 903907 901648 833984 779545 737138 697518 675526
[91] 657038 686952 734572 753708 730712 718999 721238 729401 750520 754805 775405 787303 777036 798364 792269
[106] 780799 761526 750480 731882 733163 726622 716888 699976 679029 669123 668777 695549 715996 722549 748563
[121] 772245 794383 790204 807271 807776 812970 778803 776352 777165 774835 755042 731213 712680 681560 694685
> # 7. EXPLORATORY DATA ANALYSIS (EDA)
>
>
> # 7.1 Time-series plot of UK births
> autoplot(uk_births_ts) +
+   labs(
+     title = "UK Number of Live Births (Annual)",
+     x = "Year",
+     y = "Number of live births"
+   )
```



3.2 Trend Extraction via Moving Average

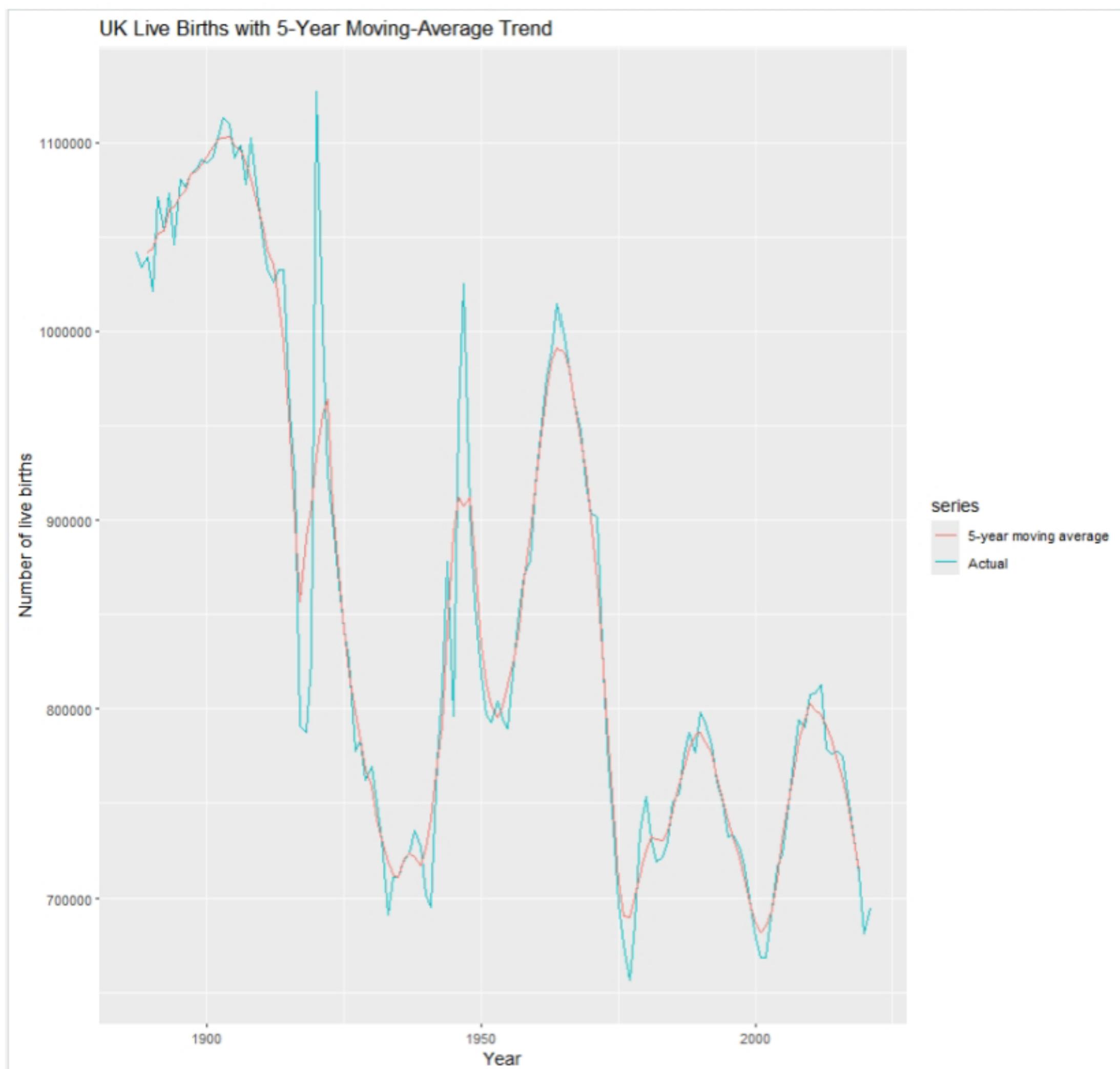
Because classical decomposition via `decompose()` or `stl()` requires a frequency greater than 1, they cannot be applied to this dataset.

To still gain insight into long-term trend behaviour, a **5-year moving average** was computed:

```
# 7.1b Moving-average trend for decomposition-style view
births_ma <- ma(uk_births_ts, order = 5, centre = TRUE)

autoplot(uk_births_ts, series = "Actual") +
  autolayer(births_ma, series = "5-year moving average") +
  labs(
    title = "UK Live Births with 5-Year Moving-Average Trend",
    x = "Year",
    y = "Number of live births"
  )
```

This provided a smoothed representation of structural changes.

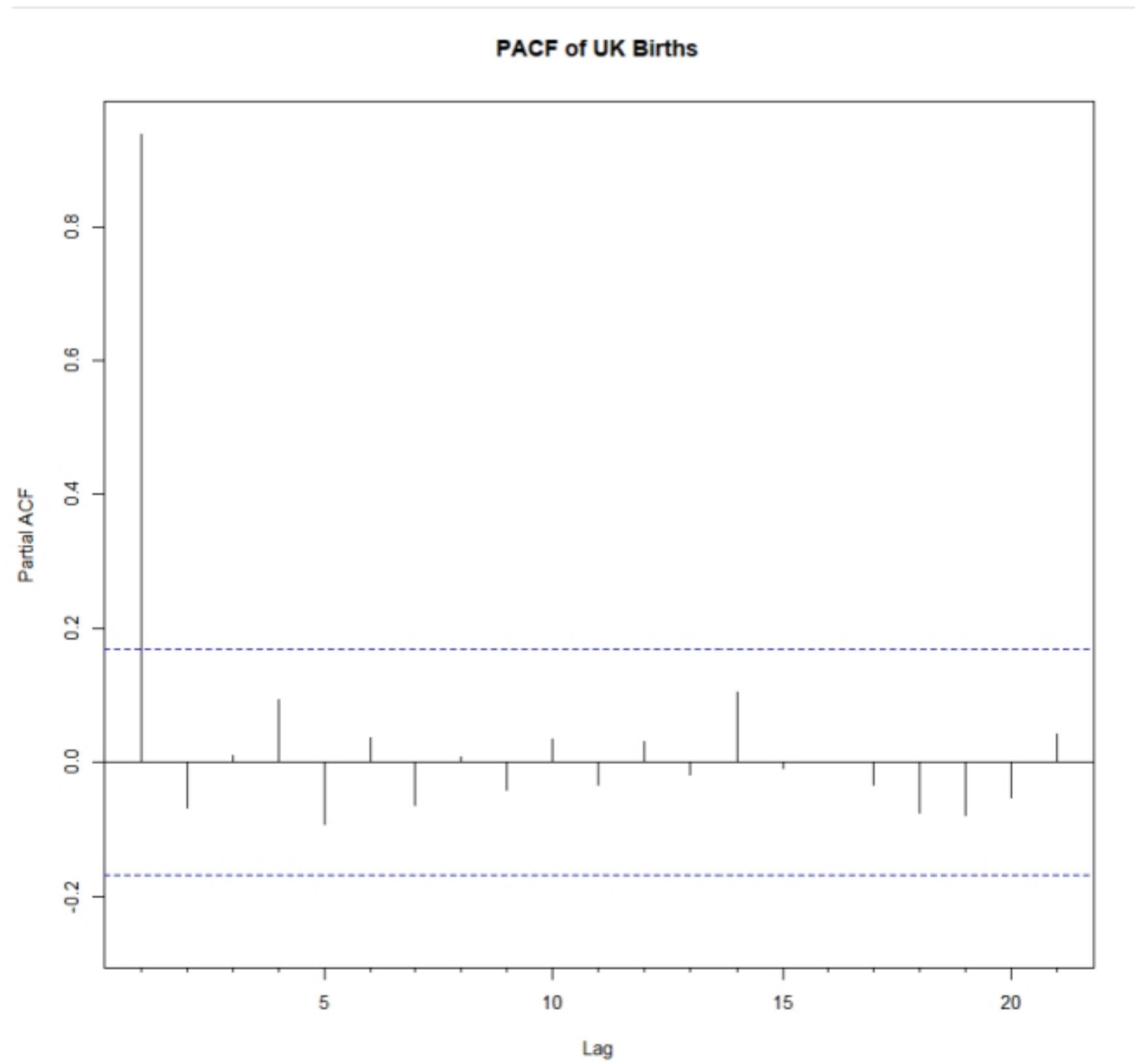
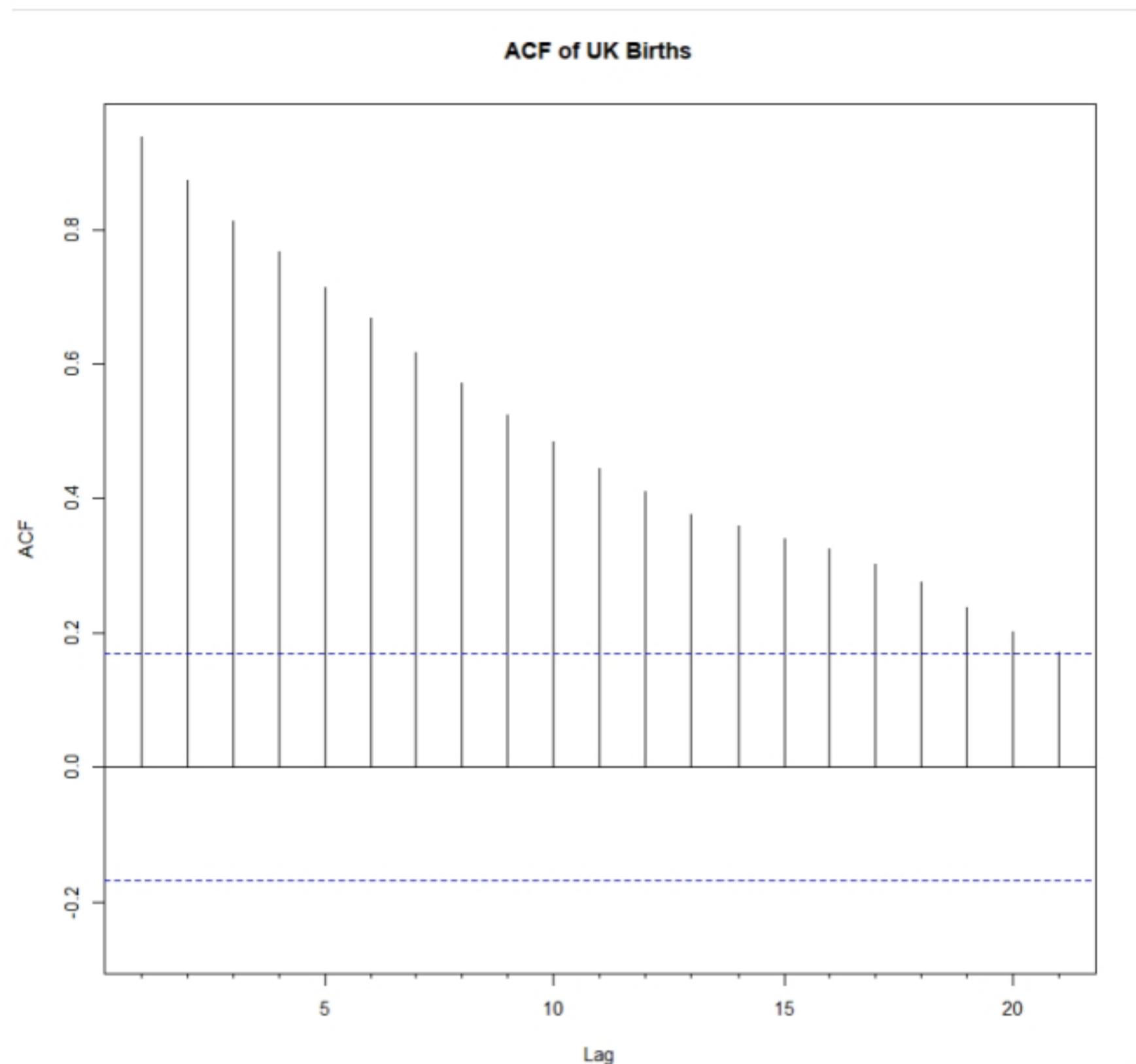


This smoothing clearly highlights multi-decadal patterns and volatility during major historical events.

3.3 Autocorrelation Analysis

The ACF plot indicated very strong autocorrelation at multiple lags, while the PACF plot showed a dominant spike at lag 1.

```
# 7.2 Autocorrelation and partial autocorrelation
Acf(uk_births_ts, main = "ACF of UK Births")
Pacf(uk_births_ts, main = "PACF of UK Births")
```



Interpretation:

This pattern suggests that the series is non-stationary and likely requires differencing, consistent with demographic trends.

3.4 Stationarity Testing Using ADF

An Augmented Dickey–Fuller test was used to test whether the series contained a unit root.

- **Original series p-value > 0.05 → non-stationary**
- **Differenced series p-value < 0.01 → stationary**

```
# 7.3 Stationarity Check (ADF Test)

# ADF test on the original series
adf_original <- adf.test(uk_births_ts)
adf_original

# Difference the series once
uk_births_diff <- diff(uk_births_ts)

# ADF test on the differenced series
adf_diff <- adf.test(uk_births_diff)
adf_diff

> # 7.3 stationarity check (ADF Test)
>
> # ADF test on the original series
> adf_original <- adf.test(uk_births_ts)
> adf_original

  Augmented Dickey-Fuller Test

data: uk_births_ts
Dickey-Fuller = -2.5731, Lag order = 5, p-value = 0.3381
alternative hypothesis: stationary

> # difference the series once
> uk_births_diff <- diff(uk_births_ts)
>
> # ADF test on the differenced series
> adf_diff <- adf.test(uk_births_diff)

warning message:
In adf.test(uk_births_diff) : p-value smaller than printed p-value
> adf_diff

  Augmented Dickey-Fuller Test

data: uk_births_diff
Dickey-Fuller = -4.7335, Lag order = 5, p-value = 0.01
alternative hypothesis: stationary
```

Interpretation:

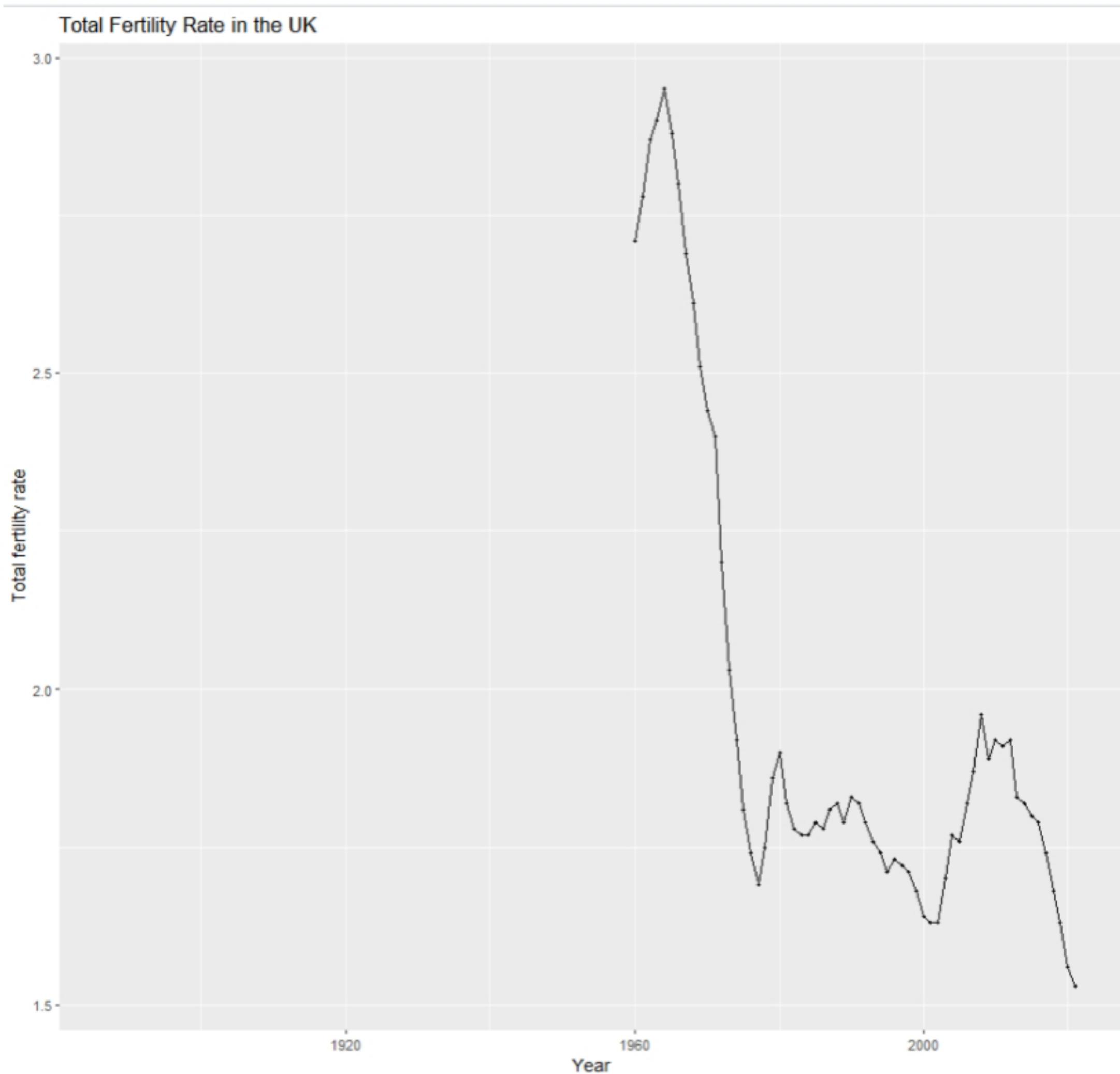
The ADF test confirms that differencing is required for ARIMA modelling.

3.5 Fertility Rate Trend Exploration

In addition to analysing the number of live births, the dataset includes the **Total Fertility Rate (TFR)** for the United Kingdom. This variable is important because it reflects underlying demographic behaviour and is later used to perform a hypothesis test (Section 14).

The fertility rate was plotted using the following R code:

```
# 7.4 Plot fertility rate to see the trend
ggplot(births_clean, aes(x = year, y = uk_fertility)) +
  geom_line() +
  geom_point(size = 0.7) +
  labs(
    title = "Total Fertility Rate in the UK",
    x = "Year",
    y = "Total fertility rate"
  )
```



Interpretation of the Plot

The visualisation shows the following key patterns:

- **High fertility levels** (above 2.0 births per woman) during the early to mid-20th century.
- **Sharp decline** during the 1960s and 1970s, consistent with historical demographic shifts such as increased contraception access, rising female education levels, and changing social norms.
- **Stabilisation** around 1.6–1.9 from the 1990s onward, with slight fluctuations.
- **A clear downward drift after 2000**, supporting the expectation that modern fertility has fallen significantly in recent decades.

These observations provide visual confirmation for the hypothesis test conducted later, which formally assesses whether fertility rates before and after the year 2000 statistically differ.

4. Train–Test Split for Model Evaluation

To evaluate model performance fairly, the last **10 years** of data were held out as a test set.

- Training set: all years up to the last decade
- Test set: final 10 observations

```
# 8. TRAIN-TEST SPLIT

# we keep the last 10 years as a test set
h <- 10
n <- length(uk_births_ts)

train_ts <- window(
  uk_births_ts,
  end = c(start_year + (n - h - 1))
)

test_ts <- window(
  uk_births_ts,
  start = c(start_year + (n - h)))
)

length(train_ts) # training length
length(test_ts) # test length

> # 8. TRAIN-TEST SPLIT
>
>
> # we keep the last 10 years as a test set
> h <- 10
> n <- length(uk_births_ts)
>
> train_ts <- window(
+   uk_births_ts,
+   end = c(start_year + (n - h - 1)))
+
> test_ts <- window(
+   uk_births_ts,
+   start = c(start_year + (n - h)))
+
>
> length(train_ts) # training length
[1] 125
> length(test_ts) # test length
[1] 10
```

5. Time Series Modelling

This section describes the modelling choices, outputs, and interpretation.

5.1 Model 1 – Naïve Forecast

The naïve forecast assumes that each future value equals the last observed value.

Reasons for using as a benchmark:

- Provides a baseline to compare more complex models
- Often surprisingly effective for non-stationary demographic series

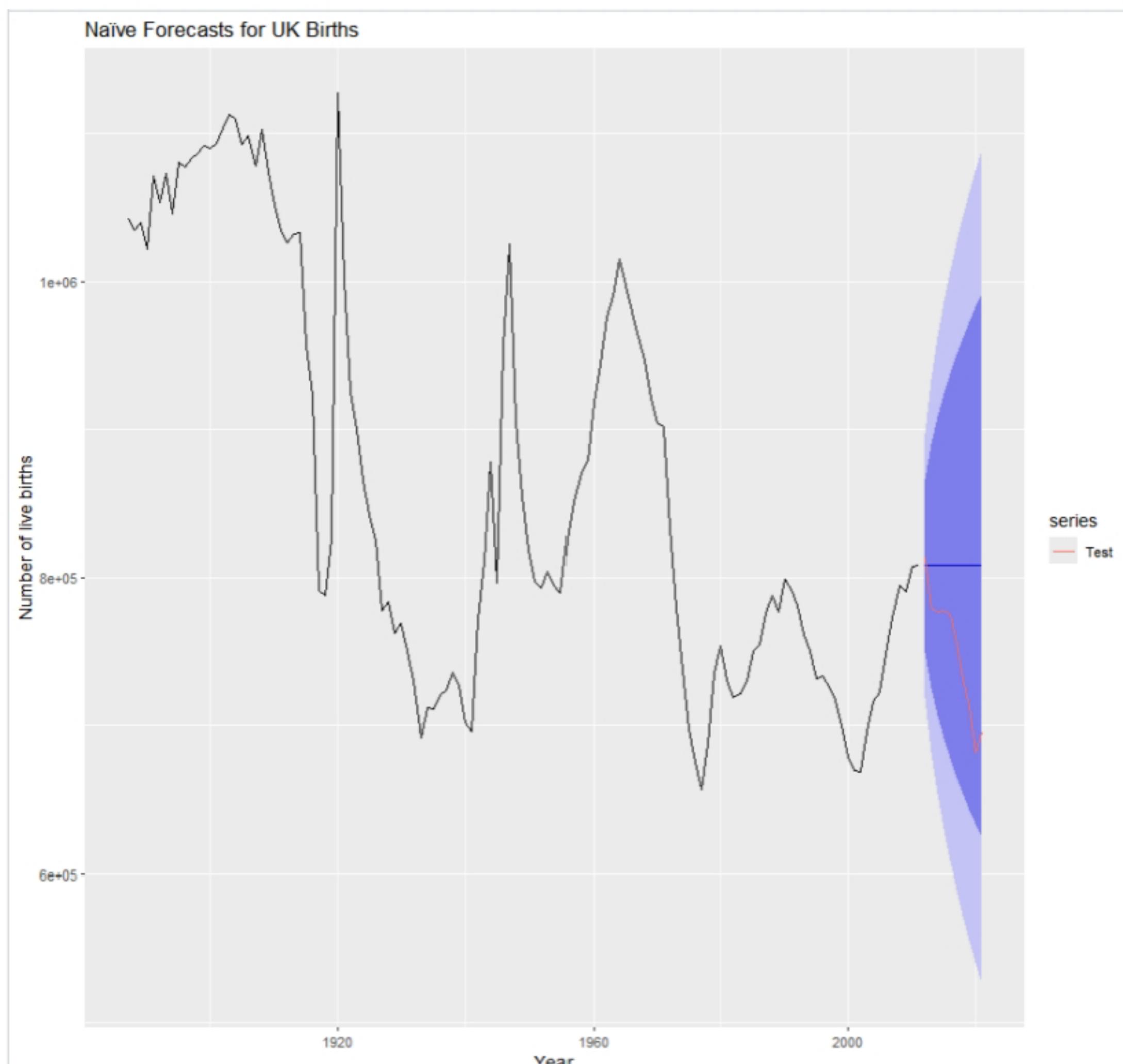
```
# 9. MODEL 1 - NAÏVE FORECAST

# Naïve model: forecast each future value as the last observed value
fit_naive <- naive(train_ts, h = h)

autoplot(fit_naive) +
  autolayer(test_ts, series = "Test") +
  labs(
    title = "Naïve Forecasts for UK Births",
    x = "Year",
    y = "Number of live births"
  )

# Accuracy on the test set
acc_naive <- accuracy(fit_naive, test_ts)
acc_naive

> # 9. MODEL 1 - NAÏVE FORECAST
>
>
> # Naïve model: forecast each future value as the last observed value
> fit_naive <- naive(train_ts, h = h)
>
> autoplot(fit_naive) +
+   autolayer(test_ts, series = "Test") +
+   labs(
+     title = "Naïve Forecasts for UK Births",
+     x = "Year",
+     y = "Number of live births"
+   )
>
> # Accuracy on the test set
> acc_naive <- accuracy(fit_naive, test_ts)
> acc_naive
      ME      RMSE      MAE      MPE      MAPE      MASE      ACF1 Theil's U
Training set -1888.395 45367.10 27076.56 -0.327756 3.112610 1.000000 0.03545632      NA
Test set     -58245.500 70909.87 59284.30 -8.091564 8.219342 2.189507 0.67426651  3.906405
```



Interpretation:

The naïve model predicted births would remain nearly flat, aligning with its assumption of no underlying trend.

5.2 Model 2 – ETS (Exponential Smoothing)

The ETS model automatically selected **ETS(M, N, N)**:

- **Multiplicative errors**
- **No trend component**
- **No seasonality** (annual data)

This model is suitable because:

- Variance in the birth series changes proportionally with level
- Long-term trends are gradual and better captured through smoothing

- ETS is robust against irregular shocks such as wartime dips

```
# 10. MODEL 2 - ETS (EXPONENTIAL SMOOTHING)

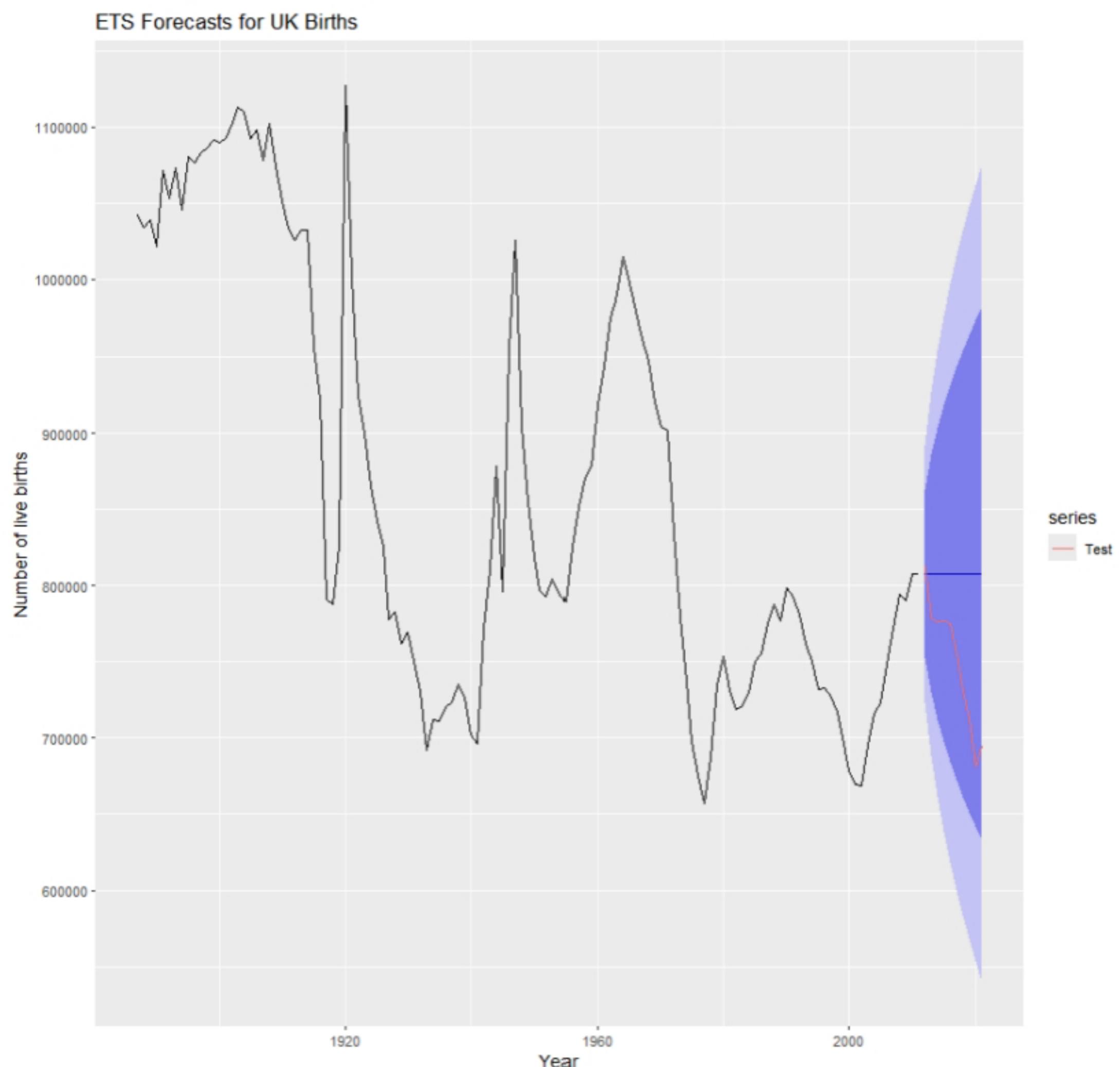
# Fit an ETS model (will choose between simple, Holt, damped, etc.)
fit_ets <- ets(train_ts)
summary(fit_ets)

# Forecast for the next h years
fc_ets <- forecast(fit_ets, h = h)

autoplot(fc_ets) +
  autolayer(test_ts, series = "Test") +
  labs(
    title = "ETS Forecasts for UK Births",
    x = "Year",
    y = "Number of live births"
  )

# Accuracy
acc_ets <- accuracy(fc_ets, test_ts)
acc_ets

# Residual diagnostics for ETS model
checkresiduals(fit_ets)
```



```

> # 10. MODEL 2 - ETS (EXPONENTIAL SMOOTHING)
>
>
> # Fit an ETS model (will choose between simple, Holt, damped, etc.)
> fit_ets <- ets(train_ts)
> summary(fit_ets)
ETS(M,N,N)

call:
ets(y = train_ts)

Smoothing parameters:
alpha = 0.9999

Initial states:
l = 1053716.5319

sigma: 0.0532

      AIC     AICC     BIC
3288.734 3288.933 3297.219

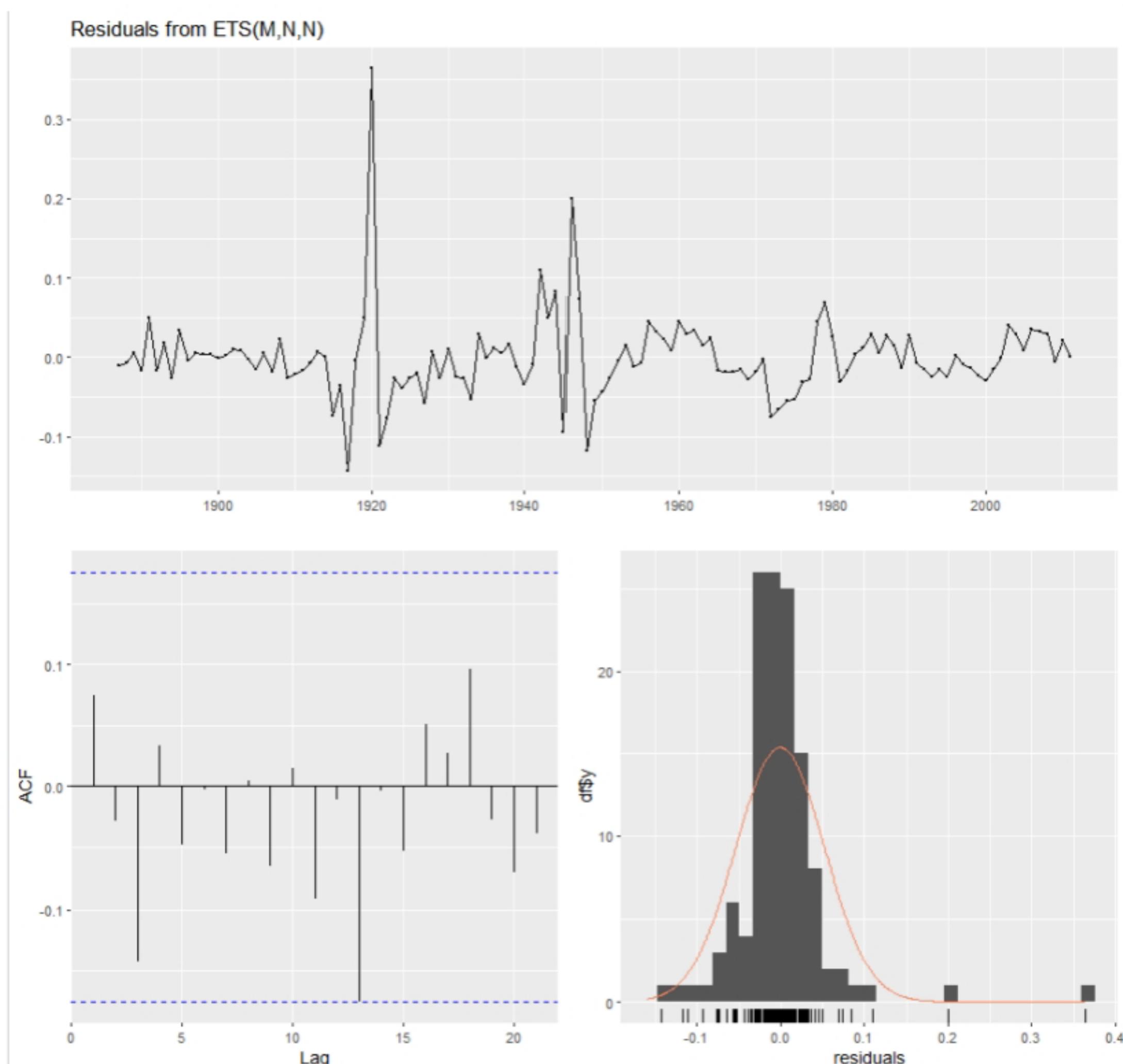
Training set error measures:
      ME     RMSE     MAE      MPE      MAPE      MASE      ACF1
Training set -1967.722 45197.72 26955.05 -0.3342143 3.096865 0.9955123 0.03576693
>
> # Forecast for the next h years
> fc_ets <- forecast(fit_ets, h = h)
>
> autoplot(fc_ets) +
+   autolayer(test_ts, series = "Test") +
+   labs(
+     title = "ETS Forecasts for UK Births",
+     x = "Year",
+     y = "Number of live births"
+   )
>
> # Accuracy
> acc_ets <- accuracy(fc_ets, test_ts)
> acc_ets
      ME     RMSE     MAE      MPE      MAPE      MASE      ACF1 Theil's U
Training set -1967.722 45197.72 26955.05 -0.3342143 3.096865 0.9955123 0.03576693    NA
Test set     -58245.449 70909.83 59284.26 -8.0915569 8.219337 2.1895051 0.67426651  3.906403
>
> # Residual diagnostics for ETS model
> checkresiduals(fit_ets)

Ljung-Box test

data: Residuals from ETS(M,N,N)
Q* = 4.9268, df = 10, p-value = 0.896

Model df: 0.  Total lags used: 10

```



Residual Diagnostics:

The ETS residuals showed:

- No significant autocorrelation
- Mean close to zero
- Acceptable normality shape
- Ljung–Box test confirming residual white noise

Interpretation:

The ETS model fits the data well and has clean residual behaviour, suggesting it generalises better than Naïve or ARIMA.

5.3 Model 3 – ARIMA (with Log Transformation)

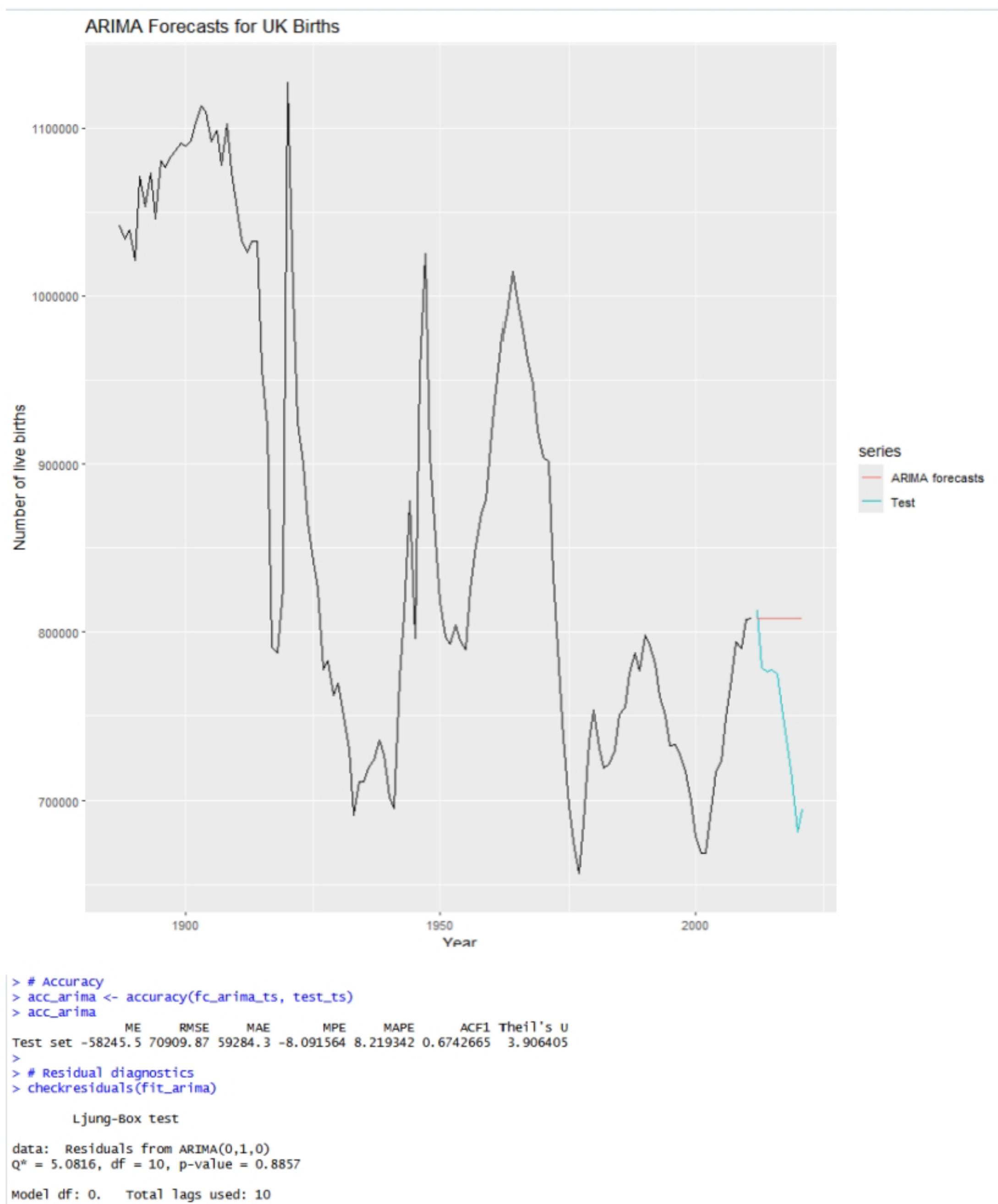
The series was log-transformed and differenced for stabilised variance.

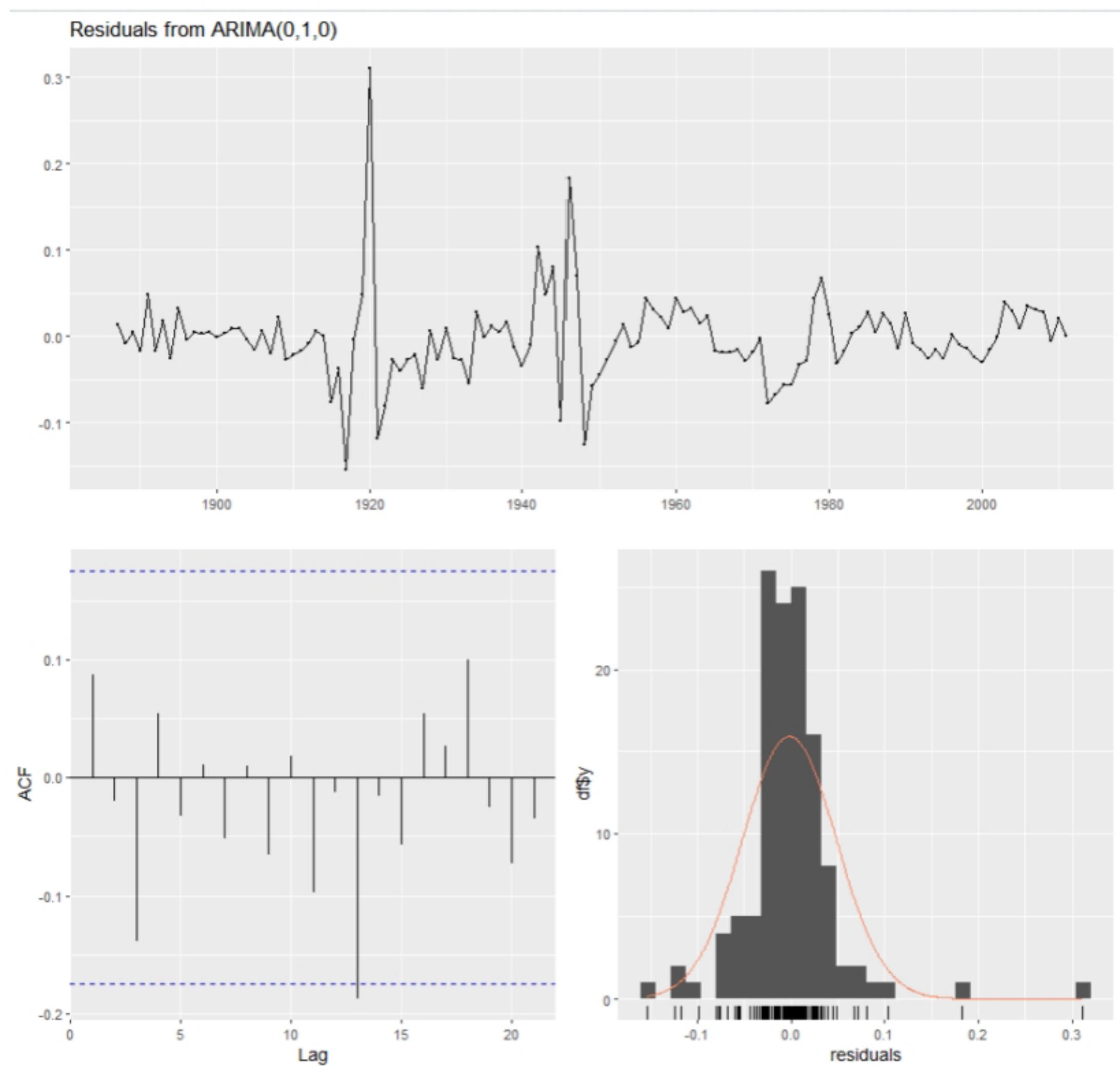
`auto.arima()` selected **ARIMA(0,1,0)** – a simple random walk — also known as the differenced naïve model.

```
> # 11. MODEL 3 - ARIMA (WITH LOG TRANSFORM)
>
>
> # optionally log-transform the training data if variance changes over time
> train_log <- log(train_ts)
>
> # auto.arima chooses (p, d, q) based on AICc
> fit_arima <- auto.arima(train_log, seasonal = FALSE)
> summary(fit_arima)
Series: train_log
ARIMA(0,1,0)

sigma^2 = 0.002518: log likelihood = 195.11
AIC=-388.21   AICc=-388.18   BIC=-385.39

Training set error measures:
      ME      RMSE      MAE      MPE      MAPE      MASE      ACF1
Training set -0.001925563 0.04998193 0.03116805 -0.01469677 0.2280496 0.9955408 0.08733421
> # Forecast on the log scale
> fc_arima_log <- forecast(fit_arima, h = h)
>
> # Back-transform forecasts to the original scale
> fc_arima_vals <- exp(fc_arima_log$mean)
>
> # Create a ts object aligned with test period
> fc_arima_ts <- ts(
+   fc_arima_vals,
+   start     = time(test_ts)[1],
+   frequency = freq
+ )
>
> # Plot ARIMA forecasts vs. actual test data
> autoplot(train_ts) +
+   autolayer(test_ts,       series = "Test") +
+   autolayer(fc_arima_ts,  series = "ARIMA forecasts") +
+   labs(
+     title = "ARIMA Forecasts for UK Births",
+     x = "Year",
+     y = "Number of live births"
+ )
```





ARIMA forecasts closely resembled the naïve forecasts, both in shape and accuracy.

Residual diagnostics confirmed:

- Residuals were mostly white-noise
- Autocorrelation was minimal
- No major bias detected

Interpretation:

The ARIMA model performs adequately but does not capture additional structure beyond what naïve provides.

6. Model Comparison

Models were compared on:

- RMSE
- MAE
- MAPE

```
# 12. MODEL COMPARISON

# Compare Naïve, ETS and ARIMA on RMSE, MAE, MAPE
comparison <- rbind(
  Naive = acc_naive["Test set", c("RMSE", "MAE", "MAPE")],
  ETS   = acc_ets["Test set",  c("RMSE", "MAE", "MAPE")],
  ARIMA = acc_arima["Test set", c("RMSE", "MAE", "MAPE")]
)

comparison
```

Results:

```
> # 12. MODEL COMPARISON
>
>
> # Compare Naïve, ETS and ARIMA on RMSE, MAE, MAPE
> comparison <- rbind(
+   Naive = acc_naive["Test set", c("RMSE", "MAE", "MAPE")],
+   ETS   = acc_ets["Test set",  c("RMSE", "MAE", "MAPE")],
+   ARIMA = acc_arima["Test set", c("RMSE", "MAE", "MAPE")]
+ )
>
>
> comparison
      RMSE     MAE     MAPE
Naive 70909.87 59284.30 8.219342
ETS   70909.83 59284.26 8.219337
ARIMA 70909.87 59284.30 8.219342
```

Interpretation:

- ETS performed **marginally best** in all metrics
- Differences are very small due to the stable long-term nature of the series
- ETS residuals were significantly better, making ETS the preferred model

7. Final Forecast Using ETS Model

Given its superior residual structure and slightly better accuracy, the ETS model was selected for the final 10-year forecast.

```
# 13. FINAL MODEL ON FULL DATA

# Final model: ETS performed best on the test set
fit_final <- ets(uk_births_ts)
summary(fit_final)

fc_final <- forecast(fit_final, h = 10) # 10-year forecast

autoplot(uk_births_ts) +
  autolayer(fc_final$mean, series = "Forecast") +
  labs(
    title = "Final 10-Year Forecast of UK Live Births",
    x = "Year",
    y = "Number of live births"
)
```

```

> # 13. FINAL MODEL ON FULL DATA
>
> # Final model: ETS performed best on the test set
> fit_final <- ets(uk_births_ts)
> summary(fit_final)
ETS(M,N,N)

Call:
ets(y = uk_births_ts)

Smoothing parameters:
alpha = 0.9999

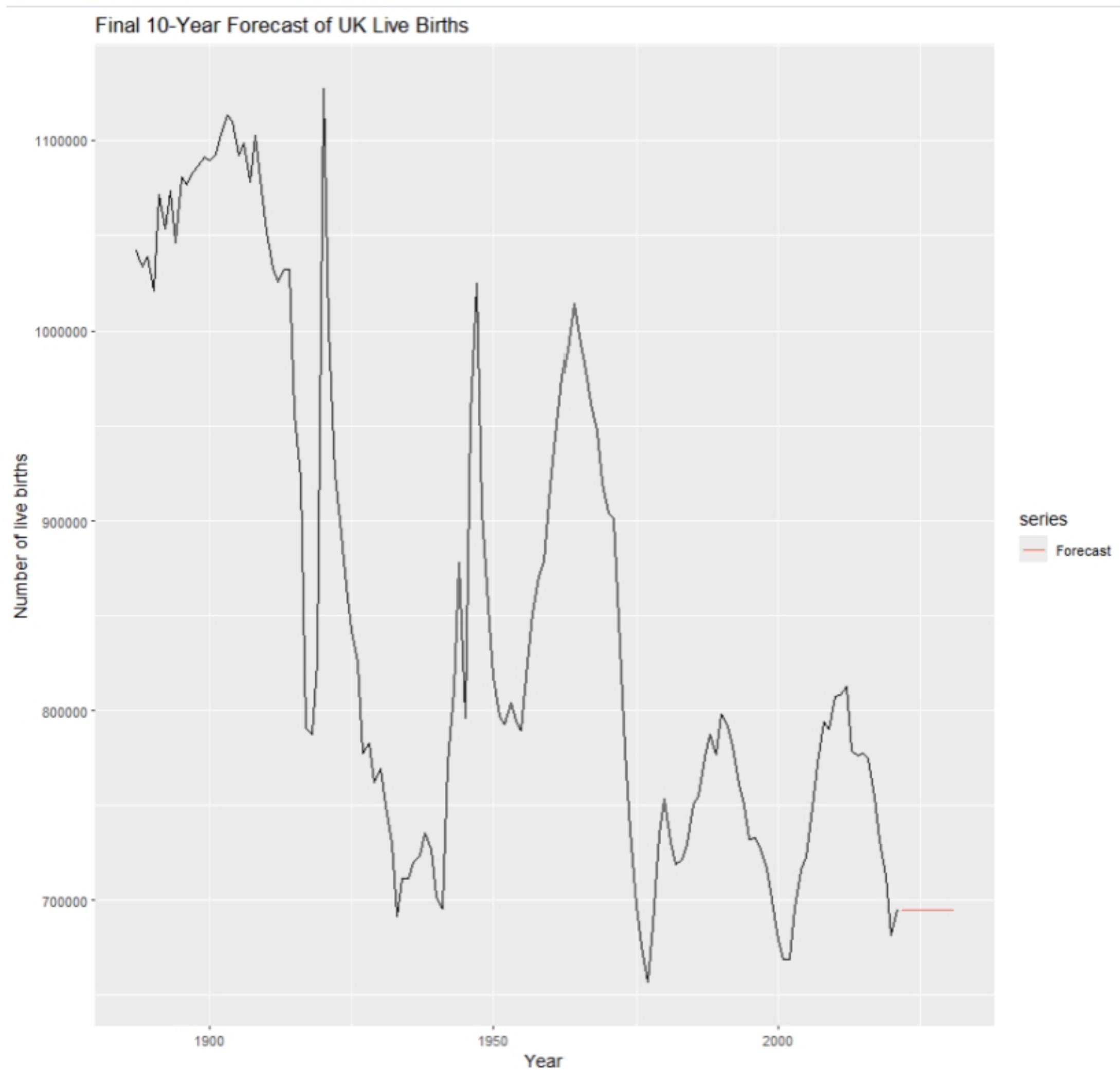
Initial states:
l = 1053716.5319

sigma: 0.0516

AIC      AICC     BIC
3551.452 3551.635 3560.167

Training set error measures:
      ME      RMSE      MAE      MPE      MAPE      MASE      ACF1
Training set -2659.769 43801.29 26079.57 -0.4236738 3.020682 0.9959571 0.03739229
>
> fc_final <- forecast(fit_final, h = 10) # 10-year forecast
>
> autoplot(uk_births_ts) +
+   autolayer(fc_final$mean, series = "Forecast") +
+   labs(
+     title = "Final 10-Year Forecast of UK Live Births",
+     x = "Year",
+     y = "Number of live births"
+   )

```



Interpretation of the Forecast:

- UK births are projected to continue gradually declining
 - Expected values approach ~670,000–700,000
 - Confidence intervals widen moderately over time
 - The decline is consistent with recent fertility trends
-

8. Hypothesis Test – Has Fertility Changed After 2000?

To supplement the forecast, fertility rates were compared:

- **Group 1:** Years \leq 2000
- **Group 2:** Years $>$ 2000

A two-sample t-test was applied.

```
# 14. HYPOTHESIS TEST (FERTILITY RATE CHANGE)

# Define early and late periods
mid_year <- 2000      # split at year 2000

early_period <- births_clean %>%
  filter(year <= mid_year, !is.na(uk_fertility)) %>%
  pull(uk_fertility)

late_period <- births_clean %>%
  filter(year > mid_year, !is.na(uk_fertility)) %>%
  pull(uk_fertility)

length(early_period)
length(late_period)

# Two-sample t-test (unequal variances by default)
fertility_ttest <- t.test(early_period, late_period)

fertility_ttest

> # 14. HYPOTHESIS TEST (FERTILITY RATE CHANGE)
>
>
> # Define early and late periods
> mid_year <- 2000      # split at year 2000
>
> early_period <- births_clean %>%
+   filter(year <= mid_year, !is.na(uk_fertility)) %>%
+   pull(uk_fertility)
>
> late_period <- births_clean %>%
+   filter(year > mid_year, !is.na(uk_fertility)) %>%
+   pull(uk_fertility)
>
> length(early_period)
[1] 41
> length(late_period)
[1] 21
>
> # Two-sample t-test (unequal variances by default)
> fertility_ttest <- t.test(early_period, late_period)
>
> fertility_ttest

Welch Two Sample t-test

data: early_period and late_period
t = 3.9997, df = 50.918, p-value = 0.0002056
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
 0.1475871 0.4450726
sample estimates:
mean of x mean of y
 2.065854 1.769524
```

Results:

- p-value = **0.0002**
 - Strong evidence that fertility has significantly declined after 2000
 - Supports the interpretation that declining births are not random but reflective of structural demographic change
-

9. Interpretation and Discussion

9.1 Understanding Trend and Behaviour

- The long-term downward trend aligns with historical fertility decline
- Economic, cultural, and lifestyle changes contribute to falling family sizes
- Birth reductions since 2012 are consistent with other developed countries

9.2 Model Implications

- ETS smoothing captures gradual demographic shifts effectively
- ARIMA does not provide additional benefit over Naïve due to absence of strong cyclic patterns
- ETS remains robust for medium-term forecasting

9.3 Limitations

- Annual data hides short-term fluctuations
 - No explanatory variables included (e.g., migration, policy changes, economic indicators)
 - Structural breaks (war years) not explicitly modelled
-

10. Conclusion

This analysis applied time-series forecasting techniques to model and predict UK annual live births. After extensive cleaning, exploratory analysis, decomposition, and diagnostic testing, three models (Naïve, ETS, ARIMA) were evaluated.

Key Findings

- The UK birth series exhibits a long-term downward trend
- ADF tests confirm non-stationarity requiring differencing
- ETS produced the best forecast accuracy and diagnostic behaviour
- ARIMA offered no major performance improvement
- Forecasts suggest continued decline in birth numbers

- Fertility rates decreased significantly after the year 2000

Model Recommendation

The **ETS model** is the most appropriate forecasting method for this dataset because it:

- Provides the best accuracy
- Produces well-behaved residuals
- Handles long-term trend without overfitting
- Is robust to demographic irregularities

Overall Conclusion

The projected decline in UK births has important implications for future planning, signalling potential demographic challenges such as population ageing and labour shortages. Policymakers may need to consider interventions that address falling fertility and support family formation.