

Stat 443: Time Series and Forecasting

Assignment 3: Time Series Models

Caden Hewlett

March 18, 2024

Question 1: El Niño Forecasting

The file `NIN034.csv` contains the monthly El Niño 3.4 index from 1870 to 2023. The El Niño 3.4 index represents the average equatorial sea surface temperature (in degrees Celsius) from around the international dateline to the coast of South America.

Part a

Perform exploratory data analysis.

Part a.1

Import the data into R and create a time-series object for the El Niño 3.4 index.

```
# import data
setwd('C:/Users/caden/OneDrive/Desktop/STAT_443/STAT443')
df <- read.csv("data/NIN034.csv")

# raw data is in a bad format, need to pivot it without years
dflong <- df %>%
  pivot_longer(cols = -Year, names_to = "Month", values_to = "Value")
# re-translate into dates with month abbreviations
dflong <- dflong %>%
  mutate(Date = make_date(Year, match(Month, month.abb), 1))
# then, format into a time series!
nino_ts <- ts(dflong$Value,
              start = c(1870, 1), frequency = 12)

# time series length = data frame length = 154 years * 12 months per year
stopifnot(all.equal(length(nino_ts), nrow(df)*12, 154*12))
```

Break the time series object into a training and test set. You can use the function `window()` on a `ts` object to split the data. Let the training set be from January 1870 to December 2021, and let the test set start in January 2022 and end in November 2023.

```

# split data into train and test
nino_train = window(nino_ts, start = c(1870, 1), end = c(2021, 12))
nino_test = window(nino_ts, start = c(2022, 1), end = c(2023, 11))
# verify split
stopifnot(all.equal(
  length(nino_test)+length(nino_train), length(na.remove(nino_ts))
))

```

Part a.2

Plot the training data as well as its acf and pacf.

Training Data

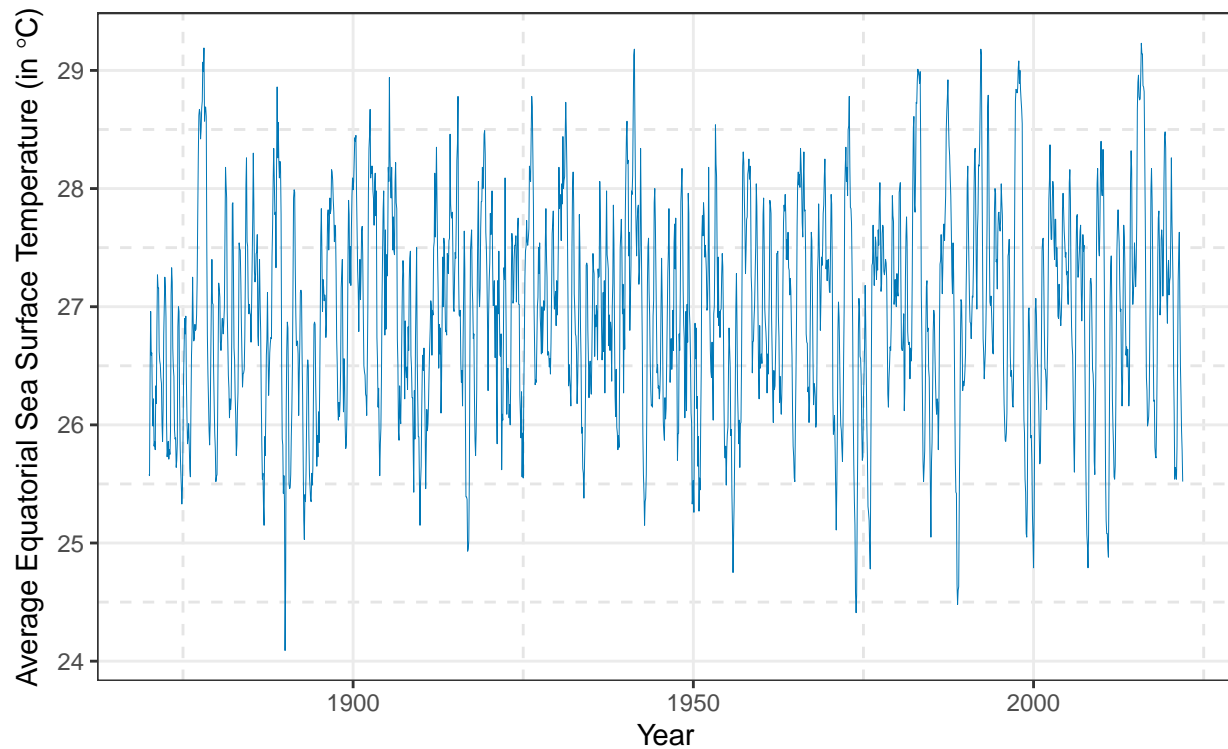
```

p1data = fortify.zoo(nino_train)
p1 <- ggplot(p1data, aes(x = Index, y = nino_train)) +
  geom_line(color = "#0077b6", linewidth = 0.1) +
  labs(
    title = "El Nino 3.4 Index from Jan. 1870 to Dec. 2021",
    subtitle = expression(
      paste("Index Represents Average Equatorial Sea Surface Temperature (in ",
        degree, "C)")),
    y = expression(
      paste("Average Equatorial Sea Surface Temperature (in ", degree, "C)")),
    x = "Year"
  ) + theme_bw() +
  theme(panel.grid.minor = element_line(
    color = "grey90",
    linetype = "dashed",
    linewidth = 0.5
  ))
print(p1)

```

El Nino 3.4 Index from Jan. 1870 to Dec. 2021

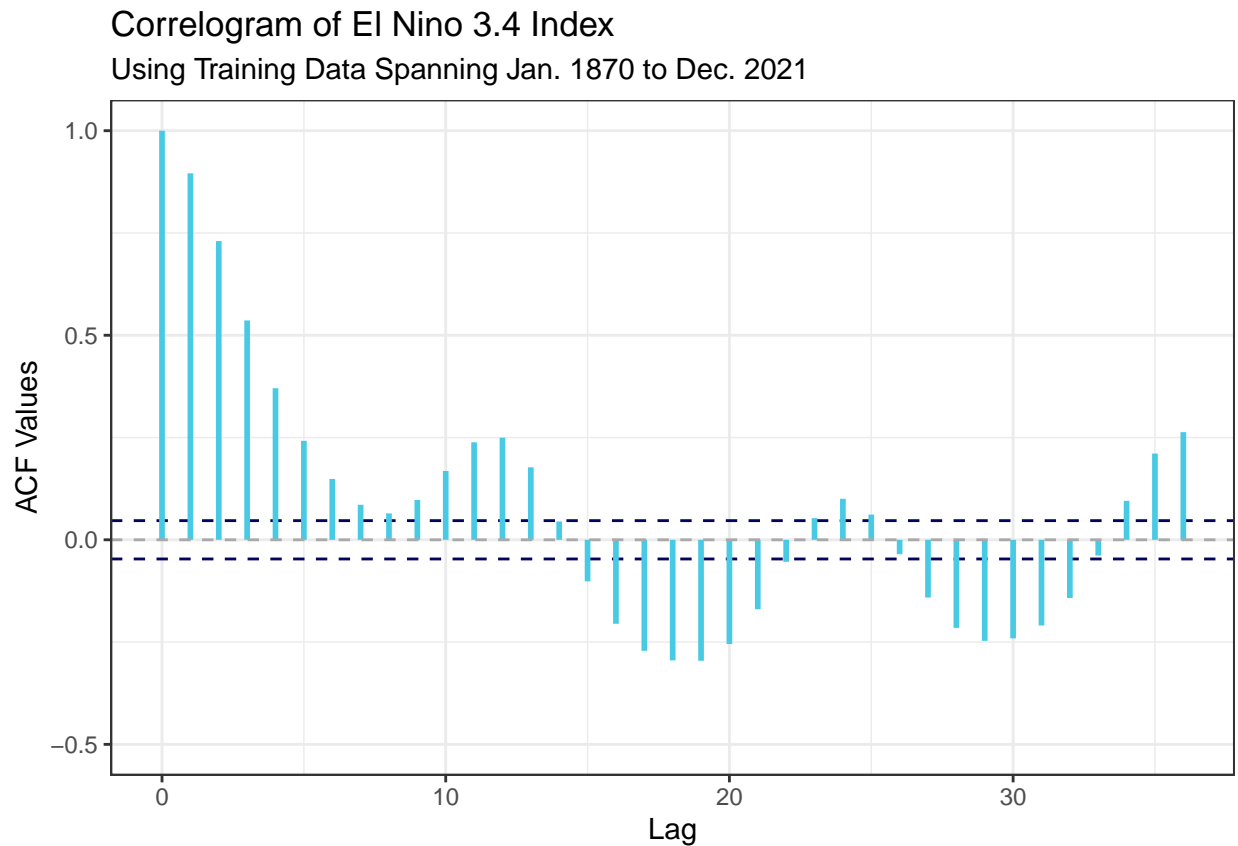
Index Represents Average Equatorial Sea Surface Temperature (in °C)



Autocorrelation Function

```
p2data = data.frame(  
  h = 0:36,  
  rh = acf(nino_train, plot = FALSE, lag.max = 36)$acf  
)  
n = length(nino_train)  
p2 <- ggplot(p2data, aes(x = h, y = rh)) +  
  geom_hline(yintercept = 2/sqrt(n),  
             linetype = "dashed",  
             col = "#03045e") +  
  geom_hline(yintercept = -2/sqrt(n),  
             linetype = "dashed",  
             col = "#03045e") +  
  ylim(-0.5, 1) +  
  geom_segment(aes(xend = h, yend = 0),  
              color = "#48cae4",  
              linewidth = 1) +  
  geom_hline(yintercept = 0,  
             linetype = "dashed",  
             color = "darkgray") +  
  labs(x = "Lag", y = "ACF Values",  
       title = "Correlogram of El Nino 3.4 Index",  
       subtitle = "Using Training Data Spanning Jan. 1870 to Dec. 2021") +  
  theme_bw()
```

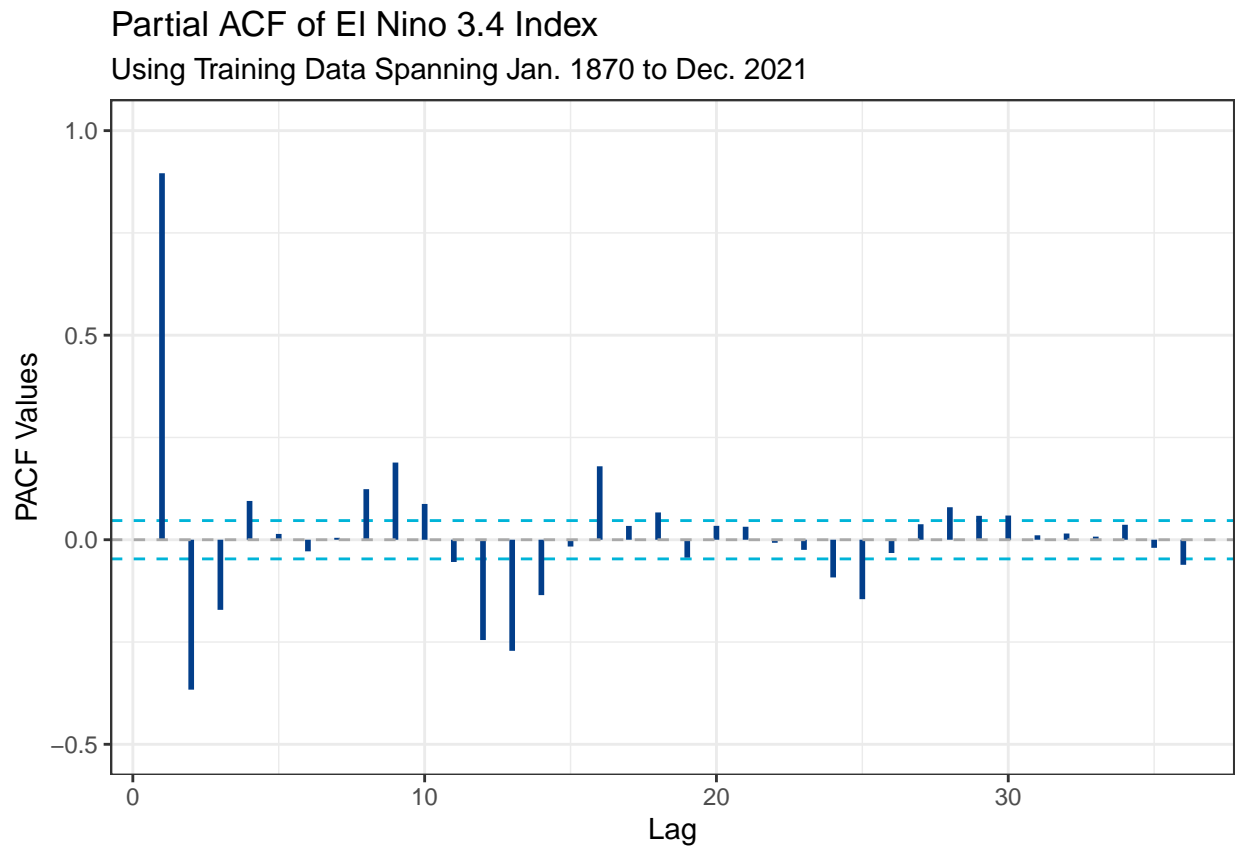
```
print(p2)
```



Partial Autocorrelation Function

```
p3data = data.frame(  
  h = 1:36,  
  rhh = pacf(nino_train, plot = FALSE, lag.max = 36)$acf  
)  
  
p3 <- ggplot(p3data, aes(x = h, y = rhh)) +  
  geom_hline(yintercept = 2/sqrt(n),  
             linetype = "dashed",  
             col = "#00b4d8") +  
  geom_hline(yintercept = -2/sqrt(n),  
             linetype = "dashed",  
             col = "#00b4d8") +  
  ylim(-0.5, 1) +  
  geom_segment(aes(xend = h, yend = 0),  
              color = "#023e8a",  
              linewidth = 1) +  
  geom_hline(yintercept = 0,  
             linetype = "dashed",  
             color = "darkgray") +  
  labs(x = "Lag", y = "PACF Values",  
       title = "Partial ACF of El Nino 3.4 Index",  
       subtitle = "Using Training Data Spanning Jan. 1870 to Dec. 2021") +
```

```
theme_bw()
print(p3)
```



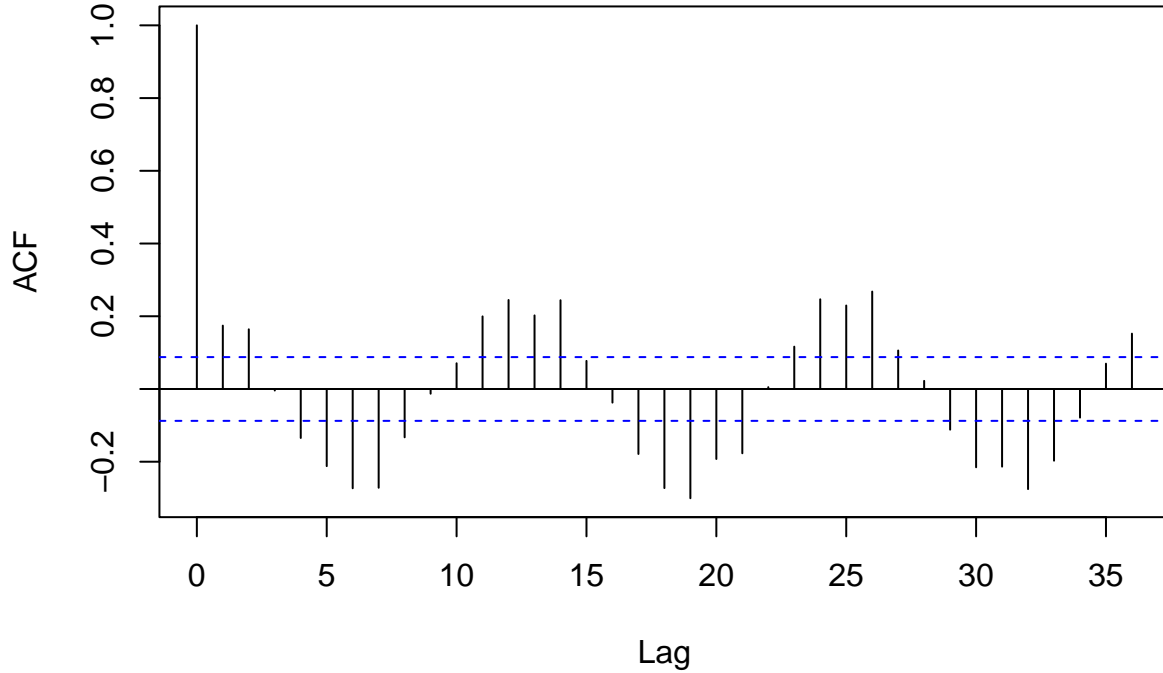
Comments

From the ACF, it is very likely that the original series exhibits a trend. This is due to the fact that there is a clear sinusoidal component of the ACF that is not dampening at a high rate (as we could potentially expect in models with AR components with $\alpha_i < 0$ for $i \in [1, p]$).

In fact, we can replicate the pattern in the sample ACF quite closely by simply using a sine function alongside random noise. Let $\{Z_t\}_{t=1}^{500} \stackrel{\text{iid}}{\sim} N(0, 1)$ and let $X_t = Z_t + \sin(t/2)$, for $t \in \{1, 2, \dots, 500\} \subset \mathbb{N}$. A plot of this artificial additive seasonal model is below:

```
set.seed(443); Z = rnorm(500); t = seq(1:500)
acf(Z + sin(t/2), main = "Artificial ACF", lag.max = 36)
```

Artificial ACF



As you can see, this artificial ACF closely matches our observed sample ACF from earlier. As such, it is very likely that there exists some seasonal term s_t in the time series of sea surface temperature data. We would intuitively anticipate a seasonal component to these data, as it is very likely that sea temperatures vary over the course of the year due to seasons as measurements are coming from the same area (the international dateline to the coast of South America.) Moreover, from a more informal observational perspective, a plot of the original series data doesn't seem to support a non-constant seasonal amplitude - as in there's no changing peak heights over time - which may indicate an additive seasonality rather than multiplicative.

In addition, neither the plot of the data nor the ACF seem to indicate a significant trend in the data. If there were a trend component m_t , we would see some consistent change over time in the overall direction of the series beyond simple seasonality. It doesn't seem like there is anything like that in this series; **however**, a purely visual analysis isn't comprehensive so this doesn't mean that a trend component doesn't exist.

Finally, to determine whether or not the series is stationary, we recall the definition of a weakly stationary stochastic process. Specifically, we will consider the first property of weak stationarity - that the mean is constant. We defined this formally in Assignment 1 previously, and was given similar to the following:

$$\text{Weak Stationarity Property One: } \exists \mu \in \mathbb{R} \text{ s.t. } \forall t \in \mathbb{Z}, \mathbb{E}(X_t) = \mu$$

The presence of *either* a seasonal component or a trend causes a contradiction, as the expected value of the series becomes some function of t (and hence cannot be a constant.)

Thus, for our particular series, we have the following implication, letting s_t be the seasonal component, $c \in \mathbb{R}$ be some real constant and $f(t)$ be a non-constant real function (as in, it changes with t .) Then,

$$\exists s_t \implies \forall t \mathbb{E}(X_t) = f(t) + c \implies \nexists \mu \in \mathbb{R} \text{ s.t. } \forall t \mathbb{E}(X_t) = \mu \therefore X_t \text{ is not stationary. } \square$$

In short, due to the presence of a trend, the time series for the training data is not stationary. It should be noted that not much can be concluded with results from the PACF, as the series is non-stationary.

Part b

Forecast sea surface temperature for 2022 and 2023 using the Box-Jenkins method and the data from 1870-2021.

Part b.1.

Remove any seasonal variation and trend from the training data, if there is any, using the `stl` function in R.

Plot the filtered data set, as well as its acf and pacf.

In experimenting with only de-seasonalizing and only de-trending the data, or both, the best performance (i.e. stationary/near stationary in the ACF) came from removing both season and trend.

```
nino_stl = stl(nino_train, s.window = "periodic")
loess_decomp = data.frame(nino_stl$time.series)
nino_filtered = nino_train - loess_decomp$seasonal # - loess_decomp$trend
acf_filtered = acf(nino_filtered, lag.max = 30, plot = F)$acf
```

```
pacf_filtered = pacf(nino_filtered, lag.max = 30, plot = F)$acf
```

```
# plot(nino_stl)
```

ideas:

```
# fit1 = arima(nino_filtered, order = c(6, 0, 0)) # -278.9
# fit2 = arima(nino_filtered, order = c(17, 0, 0)) # -278.9
# arima(nino_filtered, order = c(0, 0, 12))
# ?arima
```

```
# fit2$loglik
```

Bonus: Stepwise Model Optimization According to AIC

If we wanted to, we can also use the `auto.arima` function from the `forecast` package to pick the best model according to AIC. The function performs stepwise non-approximation for AIC and returns the model corresponding to the minimum. It takes a bit of time to run (since no approximations are done.)

```
# bonus_model = auto.arima(nino_train, ic = "aic", stepwise = TRUE, approximation = FALSE)
```

Part c

Forecast sea surface temperature for 2022 and 2023 using the Holt-Winters method and the data from 1870-2021.

Part c.1.

Use the `HoltWinters` function in R to fit an appropriate model to the training data.

As in our original series (and with the Box-Jenkins method) we did not observe a significant trend, we will construct a Holt-Winters model with $\beta = 0$.

We will use an additive seasonal effect, which inherits the **Frequency** of the time series as the period. In this case, then, $p = 12$.

We then have an additive seasonal effect I_t for t given by the following:

$$I_t = \gamma(x_t - L_t) + (1 - \gamma)I_{t-p}$$

Where L_t is the level component that defines Holt-Winters smoothing/forecasting methods. Since $\beta = 0$, it is constructed without a T component, and is given by:

$$L_t = \alpha(x_t - I_{t-p}) + (1 - \alpha)L_{t-1}$$

For $\alpha, \gamma \in [0, 1]$. Further, for $\ell \in \mathbb{Z}$, the ℓ -step ahead forecast at time t is given by the following.

$$\hat{x}_t(\ell) = L_t + I_{t-p+\ell}$$

We report the coefficients from `HoltWinters` in the table below.

```
hw_no_trend = HoltWinters(nino_train, beta = 0, seasonal = "additive") # no trend
hw_results_coefs = (data.frame(
  Alpha = round(hw_no_trend$alpha, 3),
  Beta = paste(0, ".000", sep = ""),
  Gamma = paste(hw_no_trend$gamma, ".000", sep = "")
))
rownames(hw_results_coefs) = NULL
kable(hw_results_coefs, caption = "Holt-Winters Coefficients")
```

Table 1: Holt-Winters Coefficients

Alpha	Beta	Gamma
0.961	0.000	1.000

Our fitted $\hat{\gamma} = 1$, implying that the model places maximum smoothing weight on the current state of the season (i.e. the $(1 - \gamma)I_{t-p}$ component has zero weight according to the fitting process.) Further, our fitted $\hat{\alpha} \approx 0.96$ - this means that the model is placing a lot of significance on the current level (and period) and much less on levels in the past when fitting.

Part c.2.

Use this model to predict sea surface temperature from Jan 2022 through Nov 2023.

```
hw_predictions = data.frame(
  predict(hw_no_trend, n.ahead = 23, prediction.interval = TRUE, level = 0.95))

# prepare the hw forecast data
p4data = data.frame(
  Time = as.Date(time(nino_test)),
  Observed = as.numeric(nino_test),
  Forecast = as.numeric(hw_predictions$fit),
  Lower = as.numeric(hw_predictions$lwr),
  Upper = as.numeric(hw_predictions$upr)
)
```



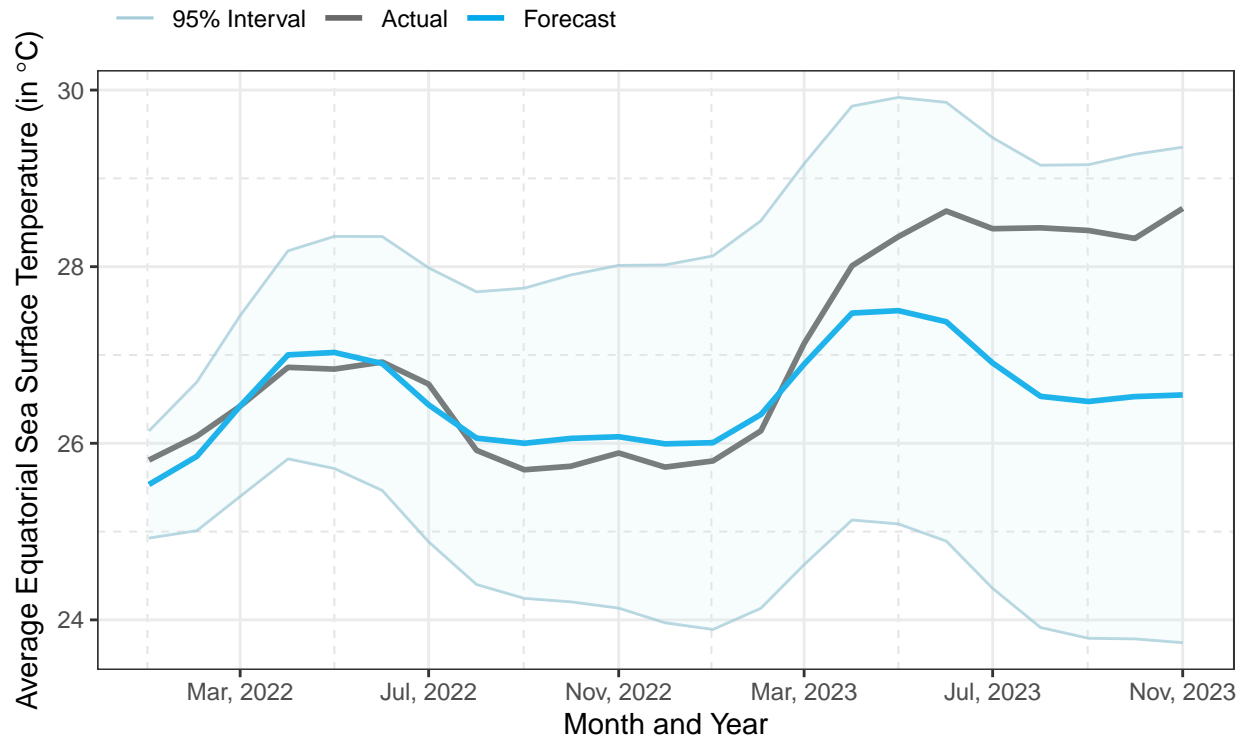
```

# make fancy as heck plot
p4 = ggplot(data = p4data, aes(x = Time)) +
  # lines for obsv, pred and interval bounds
  geom_line(aes(y = Observed, color = "Actual"), lwd = 1) +
  geom_line(aes(y = Forecast, color = "Forecast"), lwd = 1) +
  geom_line(aes(y = Lower, color = "95% Interval"), alpha = 0.75) +
  geom_line(aes(y = Upper, color = "95% Interval"), alpha = 0.75) +
  # light blue fill area between prediction intervals
  geom_ribbon(aes(ymin = Lower, ymax = Upper), fill = "#caf0f8", alpha = 0.15) +
  # legend and colour assignment
  scale_color_manual(name = "",
    values = c(
      "Actual" = "#696969",
      "Forecast" = "#02a9ea",
      "95% Interval" = "#9fc8d6"
    )) +
  # customize x-axis for nice dates
  scale_x_date(date_breaks = "4 month", date_labels = "%b, %Y") +
  # add titles with units
  labs(
    title = "Forecast and 95% Prediction Interval of Test Data",
    subtitle = "El Nino 3.4 Index from Jan. 2022 to Nov. 2023",
    x = "Month and Year",
    y = expression(
      paste("Average Equatorial Sea Surface Temperature (in ", degree, "C)")
    )
  ) +
  theme_bw() +
  theme(panel.grid.minor = element_line(
    color = "grey90",
    linewidth = 0.35,
    linetype = "dashed"
  ), legend.position = "top", legend.justification = "left",
    legend.margin = margin(0,0,0,0))
print(p4)

```

Forecast and 95% Prediction Interval of Test Data

El Nino 3.4 Index from Jan. 2022 to Nov. 2023



Part c.3.

Calculate the mean squared prediction error.

The mean squared prediction error for our $N_\ell = 23$ total ℓ step ahead forecast is given by:

$$\text{MSE}_{\text{pred}} = \frac{1}{N} \sum_{\text{Holdout}} \left(\text{Truth} - \text{Forecast} \right)^2 = \frac{1}{N_\ell} \sum_{\ell=1}^{N_\ell} \left(x_{t+\ell} - \hat{x}_t(\ell) \right)^2$$

```
N = length(nino_test)
# verify equality in set cardinalities
stopifnot(all.equal(N, length(hw_predictions$fit)))
# calculate ms prediction error
ms_pre_hw = mean( (hw_predictions$fit - nino_test)^2 )
ms_pre_hw
```

```
## [1] 0.8965602
```

Part c.4.

How does it compare to the Box-Jenkins models above?

```
# TODO: Comparison code goes here...
```

Question 2 : Hours Worked Forecasting

In this question we will predict the time series of monthly average values of the usual hours worked across all industries in Canada for the period from January 1987 until December 2023, which was explored in Assignment 1, using the file `usual_hours_worked_ca.csv`.

We'll use the Box-Jenkins method and Holt-Winters method.

Part 1: Data Preparation

Read in the data and create a time-series object for the mean monthly working hours:

```
hDF = read.csv("usual_hours_worked_ca.csv")
hours_series = ts(hDF$Hours, start = c(1987, 1), frequency = 12)
```

Separate the data into train and test. The training dataset should include all observations up to and including December 2020; The test dataset should include all observations from January 2021 to December 2023

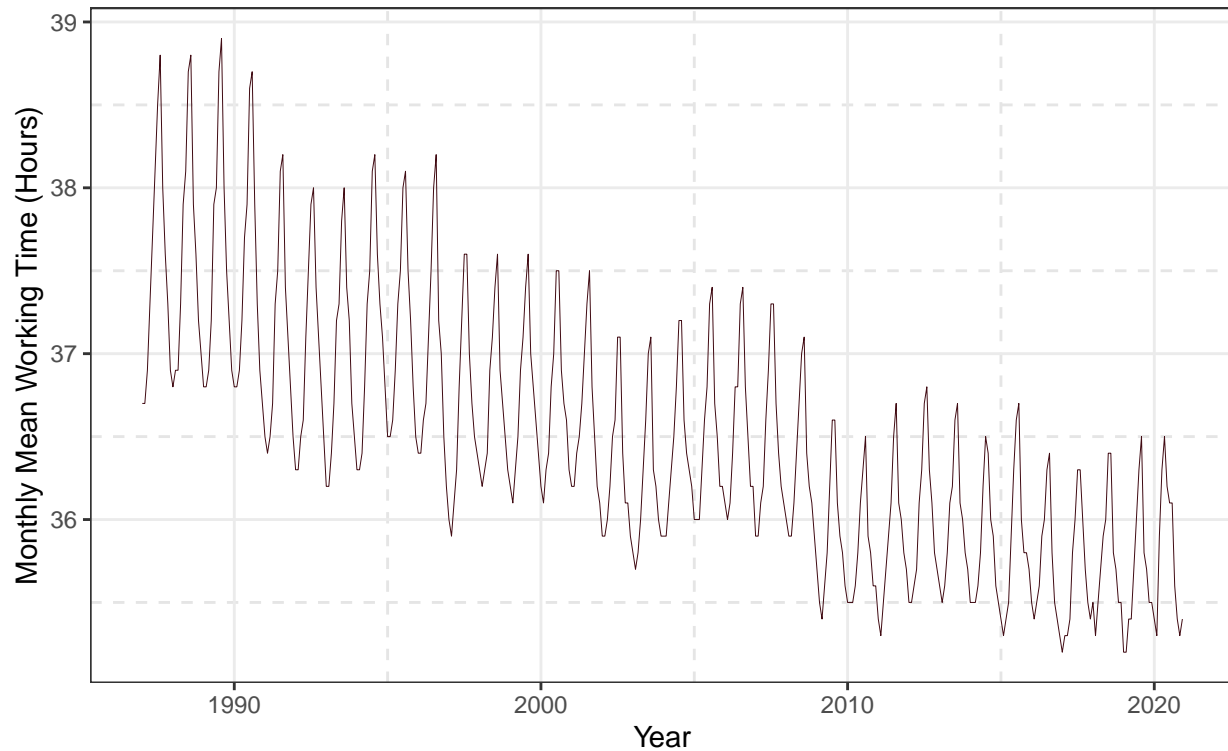
```
hours_train = window(hours_series, end = c(2020, 12))
hours_test = window(hours_series, start = c(2021, 1), end = c(2023, 12))
# verify split is valid
stopifnot(all.equal(length(hours_test)+length(hours_train),
                     length(na.remove(hours_series))))
```

Plot the training data.

```
p5data = fortify.zoo(hours_train)
n = length(hours_train)
p5 <- ggplot(p5data, aes(x = Index, y = hours_train)) +
  geom_line(color = "#410B13", linewidth = 0.1) +
  labs(
    title = "Monthly Average of Usual Hours Worked in Canada",
    subtitle = "Across all Industries from January 1987 to December 2020",
    y = "Monthly Mean Working Time (Hours)",
    x = "Year"
  ) + theme_bw() +
  theme(panel.grid.minor = element_line(
    color = "grey90",
    linetype = "dashed",
    linewidth = 0.5
  ))
print(p5)
```

Monthly Average of Usual Hours Worked in Canada

Across all Industries from January 1987 to December 2020



Part 2: Box-Jenkins Method

In this part, we select and fit a $SARIMA(p, d, q) \times (P, D, Q)_s$ model and make forecasts using the fitted model.

Part A: Differencing

Difference the training set time series at lag 1.

```
hours_train_diff = diff(hours_train, lag = 1, differences = 1)
```

Plot the new time series and its correlogram.

New Series

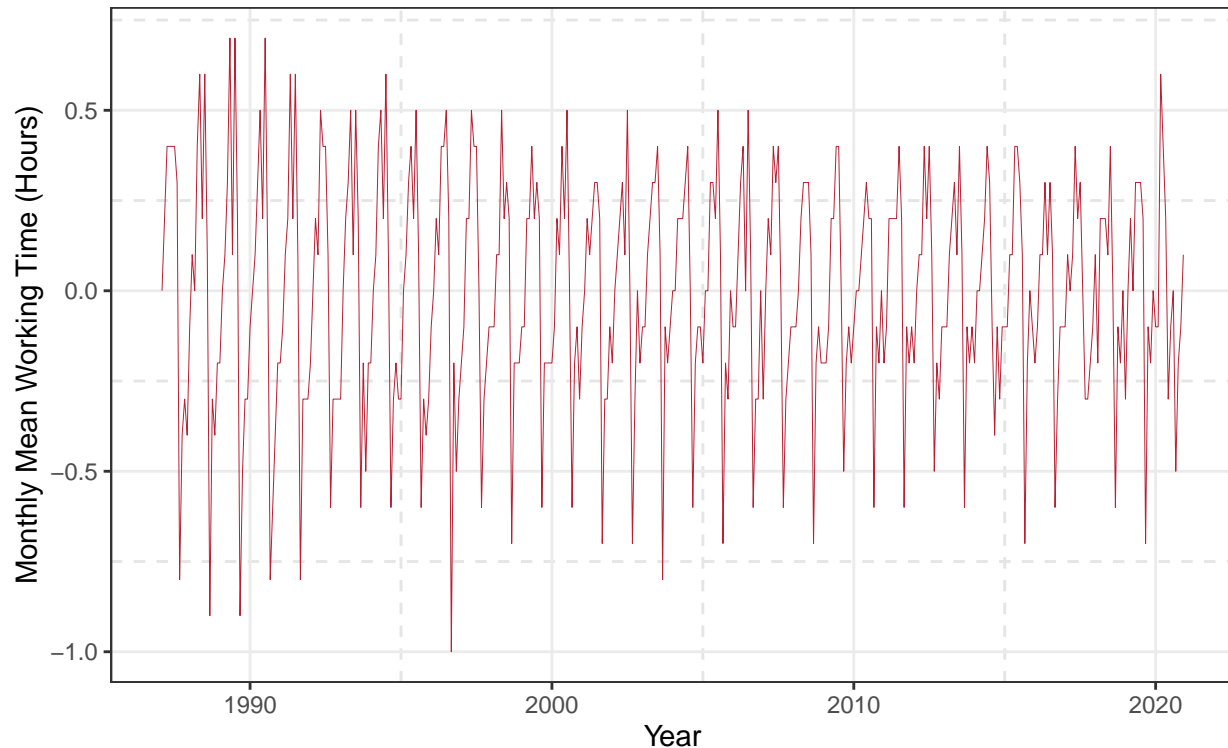
```
p6data = fortify.zoo(hours_train_diff)
p6 <- ggplot(p6data, aes(x = Index, y = hours_train_diff)) +
  geom_line(color = "#BA1F33", linewidth = 0.1) +
  labs(
    title = "Differenced Series of Monthly Average Hours Worked in Canada",
    subtitle = "Across all Industries from January 1987 to December 2020",
    y = "Monthly Mean Working Time (Hours)",
    x = "Year"
  ) + theme_bw() +
```

```

theme(panel.grid.minor = element_line(
  color = "grey90",
  linetype = "dashed",
  linewidth = 0.5
))
print(p6)

```

Differenced Series of Monthly Average Hours Worked in Canada
Across all Industries from January 1987 to December 2020



Correlogram

```

p7data = data.frame(
  h = 0:36,
  rh = acf(hours_train_diff, plot = FALSE, lag.max = 36)$acf
)
p7 <- ggplot(p7data, aes(x = h, y = rh)) +
  geom_hline(yintercept = 2/sqrt(n),
    linetype = "dashed",
    col = "#421820") +
  geom_hline(yintercept = -2/sqrt(n),
    linetype = "dashed",
    col = "#421820") +
  ylim(-0.6, 1) +
  geom_segment(aes(xend = h, yend = 0),
    color = "#CD5D67",
    linewidth = 1) +
  geom_hline(yintercept = 0,
    linetype = "dashed",

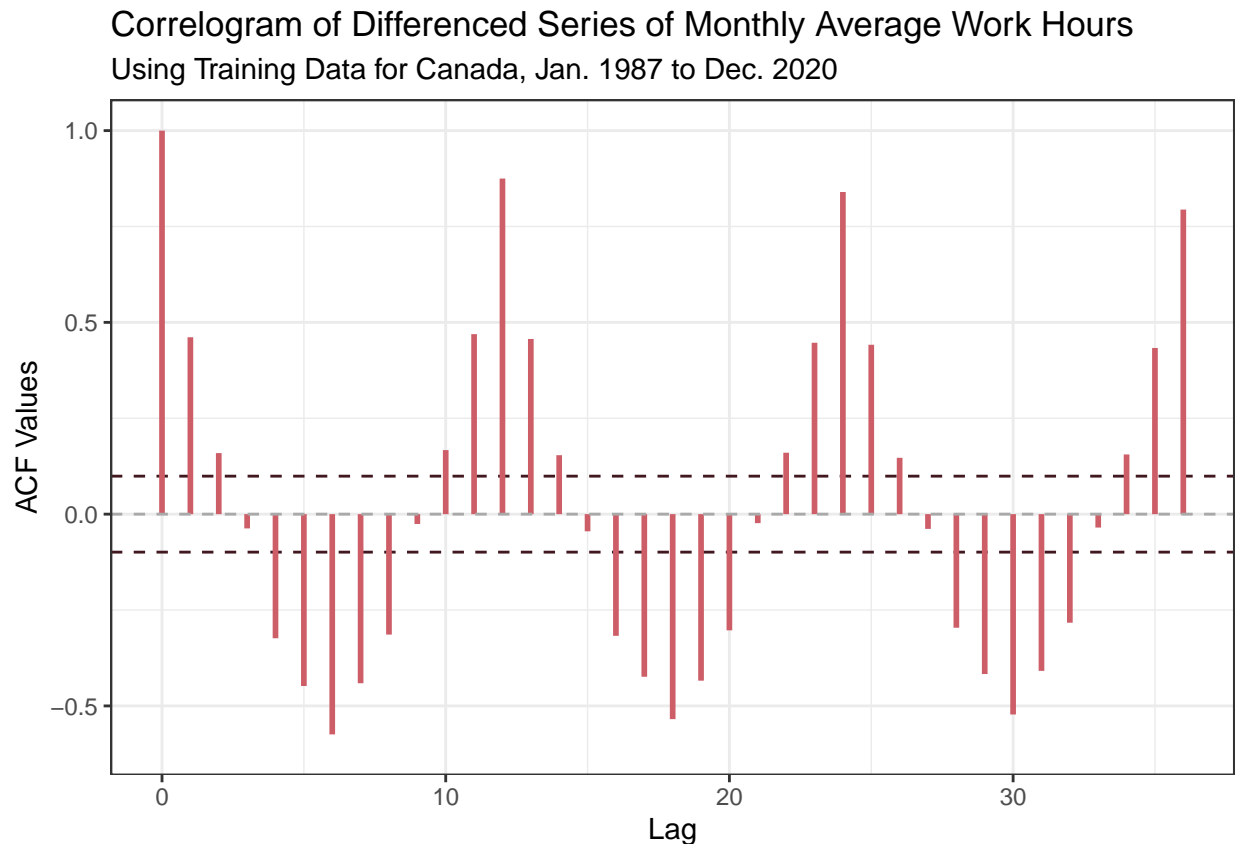
```

```

        color = "darkgray") +
labs(x = "Lag", y = "ACF Values",
     title = "Correlogram of Differenced Series of Monthly Average Work Hours",
     subtitle = "Using Training Data for Canada, Jan. 1987 to Dec. 2020") +

theme_bw()
print(p7)

```



Comment on what you observe.

While the single-iteration ($d = 1$) of sequential differencing at lag 1 appears to have helped to mitigate the trend in the original training data, as seen in the plot of the differenced training time series, there still appears to be a notable seasonal component in the data. We know that this is likely the case due to the sinusoidal component in the correlogram, which is not really decreasing as a function of lag as we would like to see. This periodicity is indicative of a lingering seasonal component that we still need to difference out or remove via other methods.

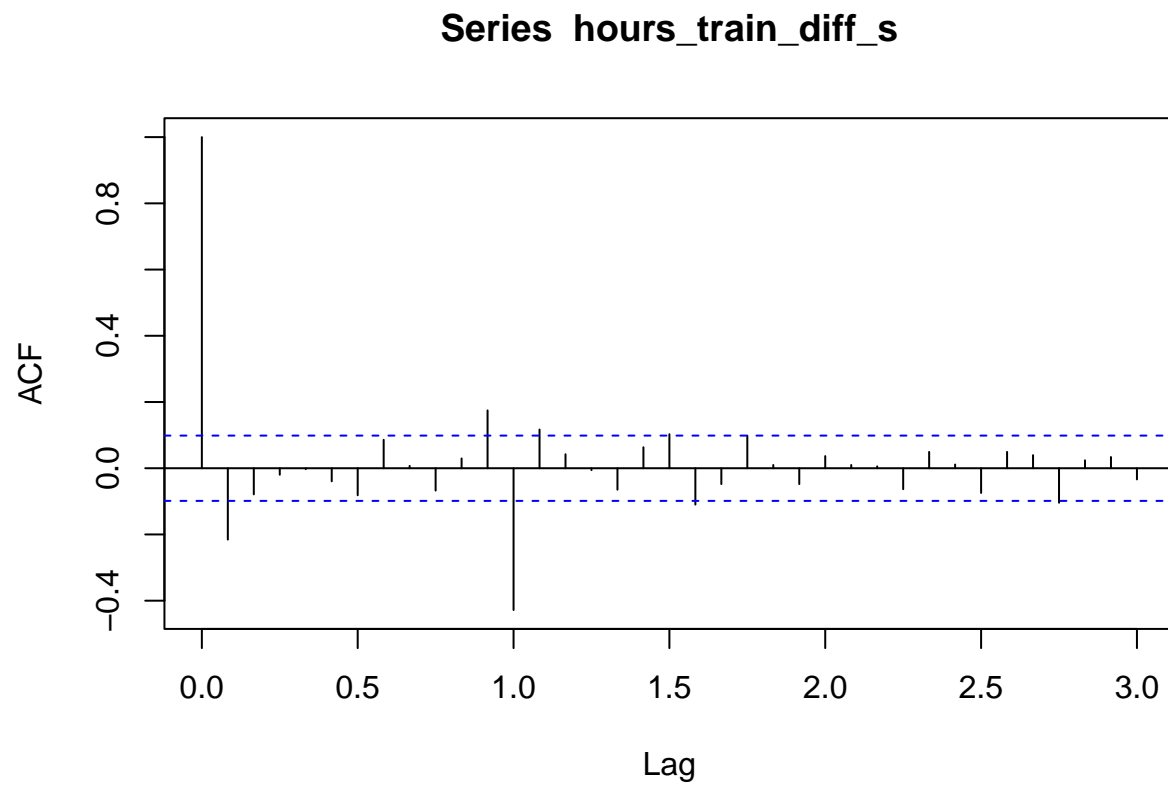
Part B: Seasonal Differencing

Apply seasonal differencing to remove seasonal variation. Since we have monthly data, $s = 12$ seems sensible.

```
hours_train_diff_s = diff(hours_train_diff, lag = 12, differences = 1)
```

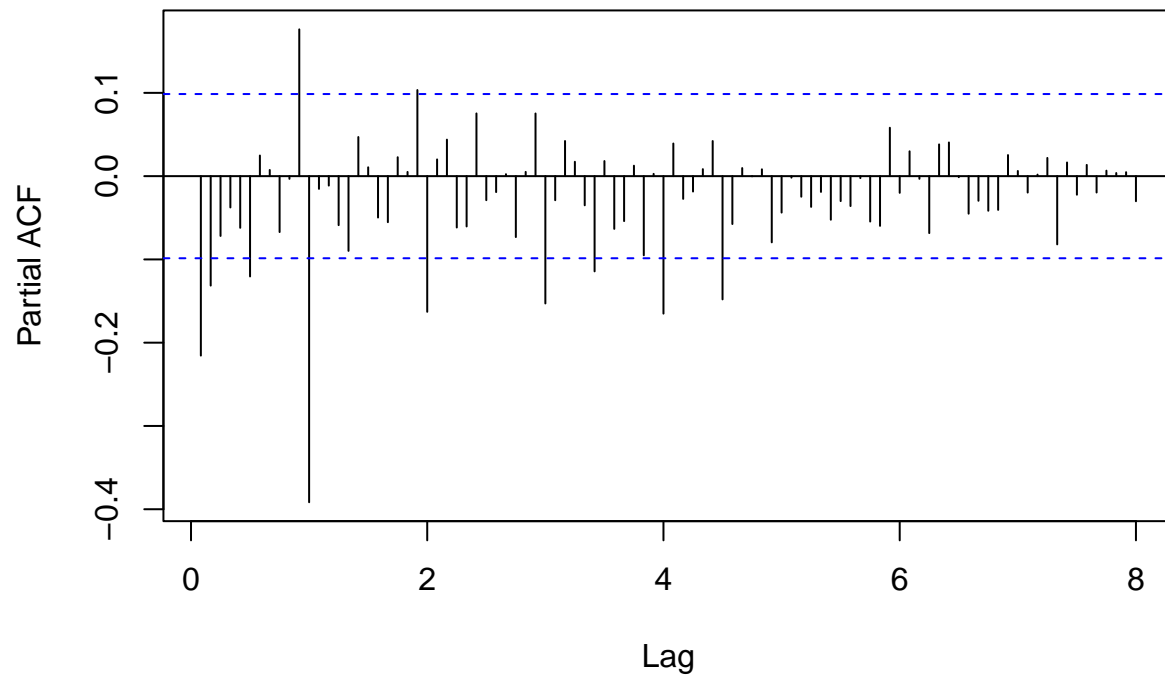
Plot the resulting differenced time series along with its sample acf and pacf.

```
acf(hours_train_diff_s, lag.max = 36) # MA at 12, Q = 1, MA at 1, q = 1
```



```
pacf(hours_train_diff_s, lag.max = 96) # PACF doesn't really show a
```

Series hours_train_diff_s



```
test = pacf(hours_train_diff_s, lag.max = 96, plot = F)$acf # P = 12 perhaps
# which( abs(test) > 2/sqrt(n) )
```

Comment on what you observe

```
auto.arima(hours_train)
```

```
## Series: hours_train
## ARIMA(1,0,2)(1,1,2)[12] with drift
##
## Coefficients:
##          ar1          ma1          ma2          sar1          sma1          sma2          drift
##          0.8692    -0.0525    -0.0083    -0.4152    -0.3223    -0.2947    -0.0047
## s.e.    0.0368     0.0670     0.0647     0.9249     0.9196     0.6688     0.0009
##
## sigma^2 = 0.01132:  log likelihood = 324.62
## AIC=-633.24   AICc=-632.87   BIC=-601.39
```

Part C: Difference Parameters.

Based on the results of Part II (a) and (b), specify the values of d , D , and s .

Part D: SARIMA Parameters

Based on the plots in Part II (b), suggest possible values of p , P , q , and Q , justifying your choices.

Part E: Iterative AIC

Now use the Akaike's Information Criterion (AIC) to select the model based on the training dataset in Part I. Fix the values of p and P as your suggestions in Part II (d),

Consider $q \in [0, 5]$ and $Q \in [0, 5]$.

Select the values of q and Q according to the AIC values.

Fit the model you choose and print the values of the estimated parameters along with the AIC value for the model.

First, we declare our fixed p , P , s , d and D values, then iterate through the different 25 models that can be developed and report the AIC for each.

```
p = 0; P = 1; s = 12; d = 1; D = 1
# declare empty matrix
results = matrix(0, nrow = 5, ncol = 5)
# with legible names
rownames(results) = c("q=1", "q=2", "q=3", "q=4", "q=5")
colnames(results) = c("Q=1", "Q=2", "Q=3", "Q=4", "Q=5")
# iterate through all candidates
for (q in 1:5){
  for (Q in 1:5){
    results[q, Q] = arima(hours_train,
                          order = c(p, d, q),
                          seasonal = c(P, D, Q))$aic
  }
}
```

The AIC results are summarized in the table below.

```
kable(round(results, 3), caption = "AIC Values for Varying (q, Q)")
```

Table 2: AIC Values for Varying (q, Q)

	Q=1	Q=2	Q=3	Q=4	Q=5
q=1	-614.916	-613.101	-611.330	-615.098	-617.794
q=2	-616.286	-614.499	-612.818	-616.883	-619.617
q=3	-621.400	-619.611	-617.792	-622.606	-625.450
q=4	-623.142	-621.333	-619.571	-624.456	-627.612
q=5	-621.897	-620.038	-618.350	-623.523	-626.358

```
P
```

```
## [1] 1
```

According to this table, the best model corresponds to:

```
which(results == min(results), arr.ind = TRUE)
```

```
##      row col
## q=4    4   5
```

Which corresponds to $q = 4$ and $Q = 5$. It has an AIC of -627.61

The full model is given below:

```
best_fit = arima(hours_train, order = c(0, 1, 4), seasonal = c(1, 1, 5))
best_fit
```

```
##
## Call:
## arima(x = hours_train, order = c(0, 1, 4), seasonal = c(1, 1, 5))
##
## Coefficients:
##          ma1          ma2          ma3          ma4          sar1          sma1          sma2          sma3
##      -0.1418   -0.1250   -0.1414   -0.1056   -0.6401   -0.1123   -0.4599   -0.0664
## s.e.    0.0525    0.0505    0.0581    0.0517    0.1678    0.1689    0.1346    0.0643
##          sma4          sma5
##      -0.0142    0.1977
## s.e.    0.0596    0.0559
##
## sigma^2 estimated as 0.01097:  log likelihood = 324.81,  aic = -627.61
```

```
best_fit$coef
```

```
##          ma1          ma2          ma3          ma4          sar1          sma1
## -0.14177013 -0.12500312 -0.14137016 -0.10562719 -0.64012373 -0.11229262
##          sma2          sma3          sma4          sma5
## -0.45986072 -0.06644207 -0.01419793  0.19767961
```