

ENV 797 - Time Series Analysis for Energy Data | Spring 2024

Assignment 5 - Due date 02/19/24

David Robinson

Directions

You should open the .rmd file corresponding to this assignment on RStudio. The file is available on our class repository on Github. And to do so you will need to fork our repository and link it to your RStudio.

Once you have the file open on your local machine the first thing you will do is rename the file such that it includes your first and last name (e.g., “LuanaLima_TSA_A05_Sp23.Rmd”). Then change “Student Name” on line 4 with your name.

Then you will start working through the assignment by **creating code and output** that answer each question. Be sure to use this assignment document. Your report should contain the answer to each question and any plots/tables you obtained (when applicable).

When you have completed the assignment, **Knit** the text and code into a single PDF file. Submit this pdf using Sakai.

R packages needed for this assignment: “readxl”, “ggplot2”, “forecast”, “tseries”, and “Kendall”. Install these packages, if you haven’t done yet. Do not forget to load them before running your script, since they are NOT default packages.\

```
#Load/install required package here
```

```
library(readxl)
library(forecast)
```

```
## Warning: package 'forecast' was built under R version 4.3.2
```

```
## Registered S3 method overwritten by 'quantmod':
##   method      from
##   as.zoo.data.frame zoo
```

```
library(tseries)
library(ggplot2)
library(Kendall)
library(lubridate)
```

```
##
## Attaching package: 'lubridate'
```

```
## The following objects are masked from 'package:base':
##
##   date, intersect, setdiff, union
```

```
library(tidyverse) #load this package so yon clean the data frame using pipes
```

```
## -- Attaching core tidyverse packages ----- tidyverse 2.0.0 --
## v dplyr 1.1.3 v stringr 1.5.0
## v forcats 1.0.0 v tibble 3.2.1
## v purrr 1.0.2 v tidyr 1.3.0
## v readr 2.1.4
```

```
## -- Conflicts ----- tidyverse_conflicts() --
## x dplyr::filter() masks stats::filter()
## x dplyr::lag() masks stats::lag()
## i Use the conflicted package (<http://conflicted.r-lib.org/>) to force all conflicts to become errors
```

Decomposing Time Series

Consider the same data you used for A04 from the spreadsheet “Table_10.1_Renewable_Energy_Production_and_Consumption”. The data comes from the US Energy Information and Administration and corresponds to the December 2023 Monthly Energy Review.

```
#Importing data set - using xlsx package
```

```
energy_data <- read.csv(file="./Data/Table_10.1_Renewable_Energy_Production_and_Consumption_by_Source.csv")
```

```
#Now let's extract the column names from row 11 only
```

```
read_col_names <- read_excel("./Data/Table_10.1_Renewable_Energy_Production_and_Consumption_by_Source.xlsx", sheet = 11, col_names = TRUE)
```

```
## New names:
```

```
## * ' -> '...1'
## * ' -> '...2'
## * ' -> '...3'
## * ' -> '...4'
## * ' -> '...5'
## * ' -> '...6'
## * ' -> '...7'
## * ' -> '...8'
## * ' -> '...9'
## * ' -> '...10'
## * ' -> '...11'
## * ' -> '...12'
## * ' -> '...13'
## * ' -> '...14'
```

```
colnames(energy_data) <- read_col_names
head(energy_data)
```

```
##           Month Wood Energy Production Biofuels Production
## 1 1973 January           129.630           Not Available
## 2 1973 February          117.194           Not Available
## 3   1973 March           129.763           Not Available
## 4   1973 April           125.462           Not Available
## 5    1973 May            129.624           Not Available
```

## 6	1973 June	125.435	Not Available
##	Total Biomass Energy Production	Total Renewable Energy Production	
## 1		129.787	219.839
## 2		117.338	197.330
## 3		129.938	218.686
## 4		125.636	209.330
## 5		129.834	215.982
## 6		125.611	208.249
##	Hydroelectric Power Consumption	Geothermal Energy Consumption	
## 1		89.562	0.490
## 2		79.544	0.448
## 3		88.284	0.464
## 4		83.152	0.542
## 5		85.643	0.505
## 6		82.060	0.579
##	Solar Energy Consumption	Wind Energy Consumption	Wood Energy Consumption
## 1	Not Available	Not Available	129.630
## 2	Not Available	Not Available	117.194
## 3	Not Available	Not Available	129.763
## 4	Not Available	Not Available	125.462
## 5	Not Available	Not Available	129.624
## 6	Not Available	Not Available	125.435
##	Waste Energy Consumption	Biofuels Consumption	
## 1	0.157	Not Available	
## 2	0.144	Not Available	
## 3	0.176	Not Available	
## 4	0.174	Not Available	
## 5	0.210	Not Available	
## 6	0.176	Not Available	
##	Total Biomass Energy Consumption	Total Renewable Energy Consumption	
## 1		129.787	219.839
## 2		117.338	197.330
## 3		129.938	218.686
## 4		125.636	209.330
## 5		129.834	215.982
## 6		125.611	208.249

```
nobs=nrow(energy_data)
nvar=ncol(energy_data)

t <- 1:nobs
```

Q1

For this assignment you will work only with the following columns: Solar Energy Consumption and Wind Energy Consumption. Create a data frame structure with these two time series only and the Date column. Drop the rows with *Not Available* and convert the columns to numeric. You can use filtering to eliminate the initial rows or convert to numeric and then use the `drop_na()` function. If you are familiar with pipes for data wrangling, try using it!

```
energy_data$Month <- as.Date(paste(energy_data$Month, "01"),
                             format = "%Y %B %d")
```

```
energy_data_solar_wind <- energy_data[,8:9]
energy_data_dates <- energy_data[,1]
energy_data <- cbind(energy_data_dates,energy_data_solar_wind)

energy_data[, 2:3] <- lapply(energy_data[, 2:3], as.numeric)
```

```
## Warning in lapply(energy_data[, 2:3], as.numeric): NAs introduced by coercion
```

```
## Warning in lapply(energy_data[, 2:3], as.numeric): NAs introduced by coercion
```

```
energy_data <- drop_na(energy_data)

nobs=nrow(energy_data)
t <- 1:nobs

ts_energy_data <- ts(energy_data[t,1:3], frequency=12,start=c(1984,1))
```

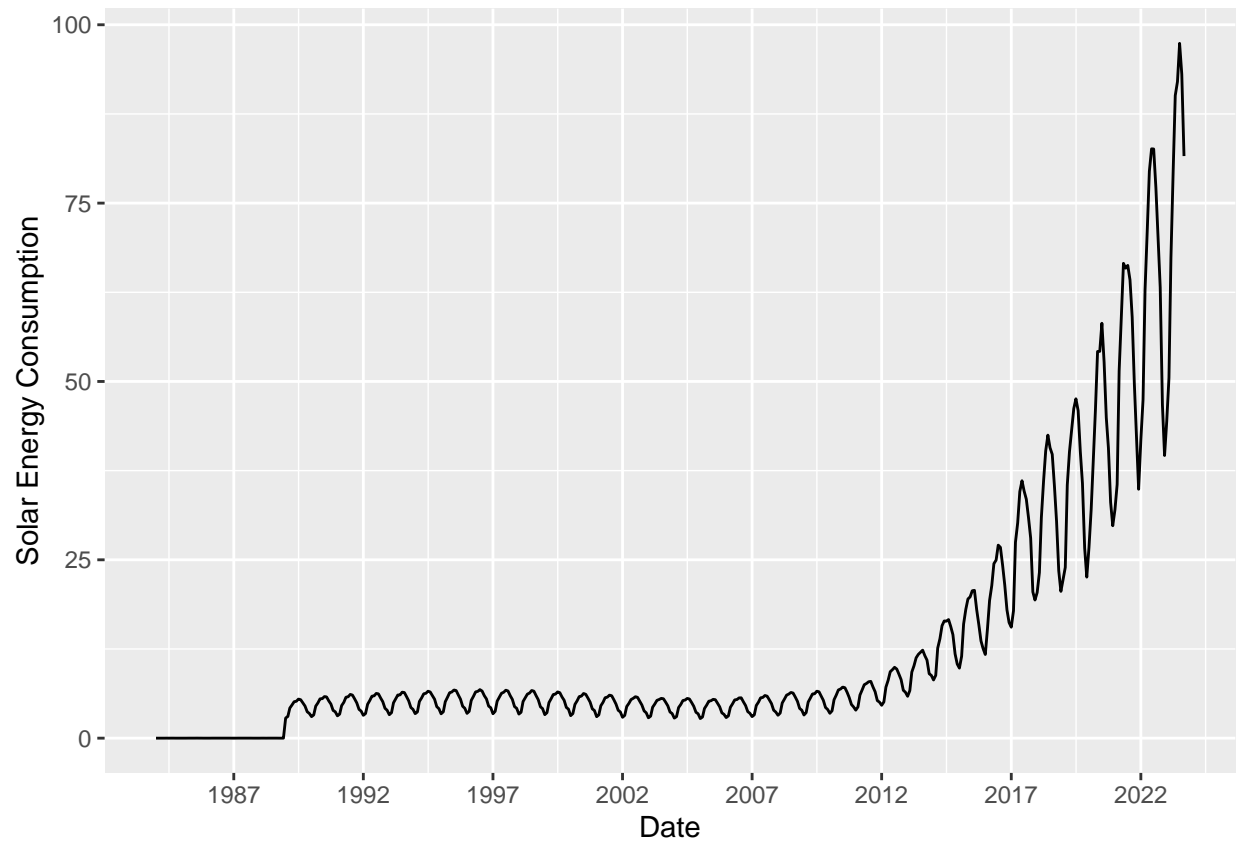
Q2

Plot the Solar and Wind energy consumption over time using ggplot. Plot each series on a separate graph. No need to add legend. Add informative names to the y axis using `ylab()`. Explore the function `scale_x_date()` on ggplot and see if you can change the x axis to improve your plot. Hint: use `scale_x_date(date_breaks = "5 years", date_labels = "%Y")`)

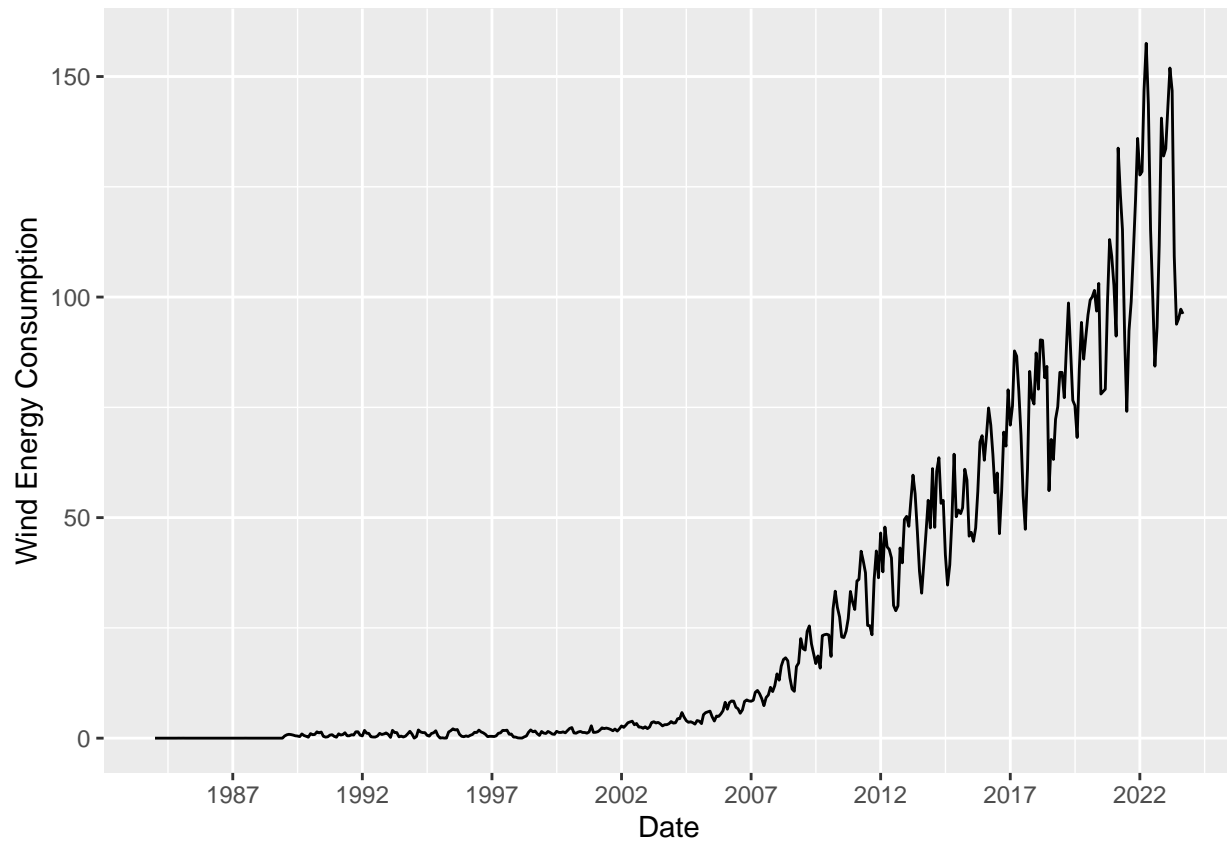
```
plot_solar <- ggplot(energy_data,aes(x = energy_data_dates,
                                     y = energy_data[,2])) +
  geom_line() +
  labs(x = "Date", y = "Solar Energy Consumption") +
  scale_x_date(date_breaks = "5 years", date_labels = "%Y")

plot_wind <- ggplot(energy_data, aes(x = energy_data_dates,
                                     y = energy_data[,3])) +
  geom_line() +
  labs(x = "Date", y = "Wind Energy Consumption") +
  scale_x_date(date_breaks = "5 years", date_labels = "%Y")

plot_solar
```



plot_wind

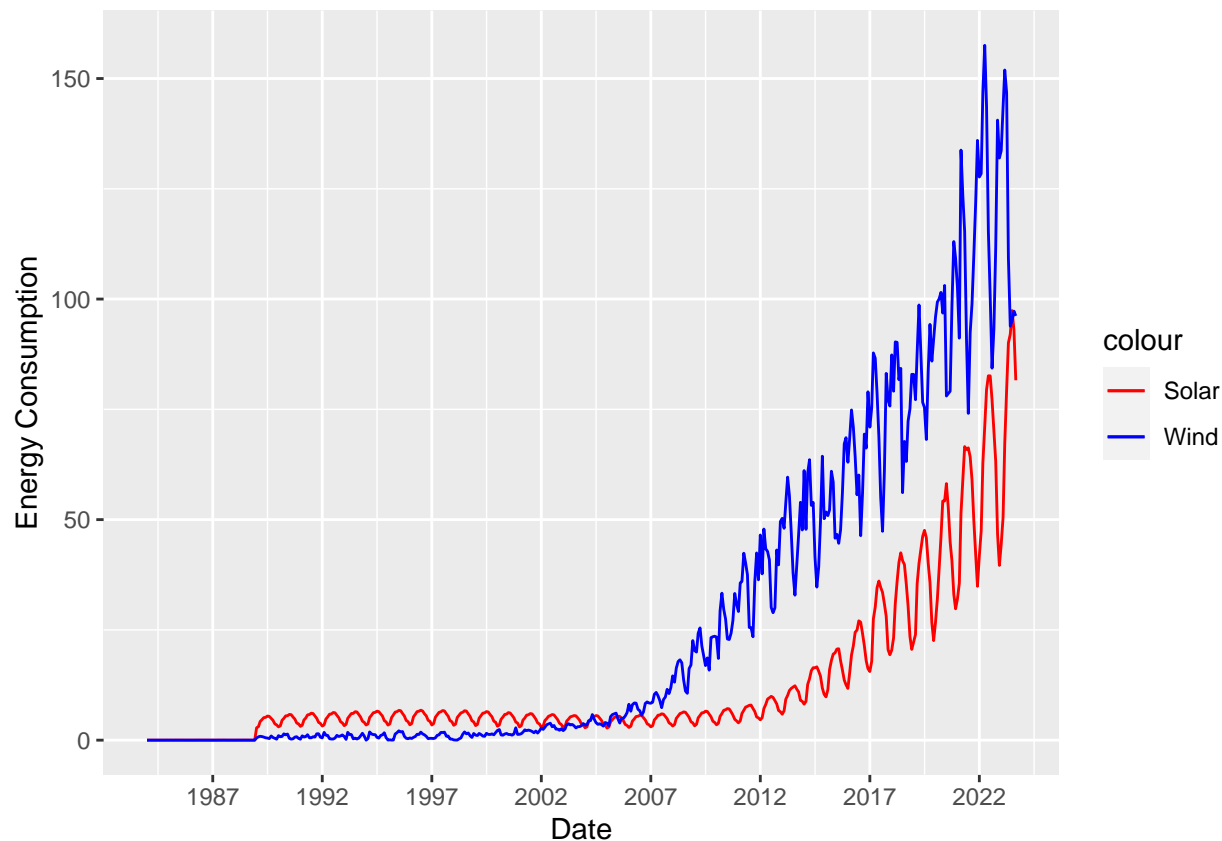


Q3

Now plot both series in the same graph, also using `ggplot()`. Use function `scale_color_manual()` to manually add a legend to `ggplot`. Make the solar energy consumption red and wind energy consumption blue. Add informative name to the y axis using `ylab("Energy Consumption")`. And use function `scale_x_date()` to set x axis breaks every 5 years.

```
combined_plot <- ggplot(energy_data, aes(x = energy_data_dates)) +
  geom_line(aes(y = energy_data[,2], color = "Solar")) +
  geom_line(aes(y = energy_data[,3], color = "Wind")) +
  labs(x = "Date", y = "Energy Consumption") +
  scale_color_manual(values = c("Solar" = "red", "Wind" = "blue")) +
  scale_x_date(date_breaks = "5 years", date_labels = "%Y")

combined_plot
```



Decomposing the time series

The stats package has a function called `decompose()`. This function only take time series object. As the name says the decompose function will decompose your time series into three components: trend, seasonal and random. This is similar to what we did in the previous script, but in a more automated way. The random component is the time series without seasonal and trend component.

Additional info on `decompose()`.

- 1) You have two options: alternative and multiplicative. Multiplicative models exhibit a change in frequency over time.
- 2) The trend is not a straight line because it uses a moving average method to detect trend.
- 3) The seasonal component of the time series is found by subtracting the trend component from the original data then grouping the results by month and averaging them.
- 4) The random component, also referred to as the noise component, is composed of all the leftover signal which is not explained by the combination of the trend and seasonal component.

Q4

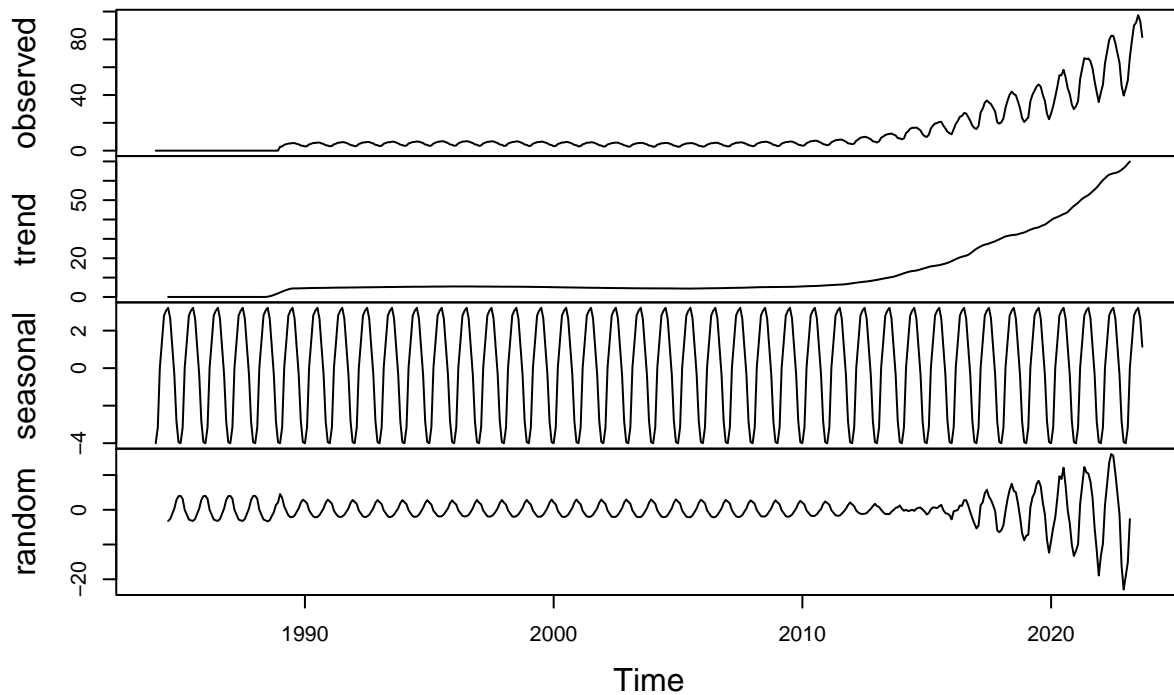
Transform wind and solar series into a time series object and apply the decompose function on them using the additive option, i.e., `decompose(ts_data, type = "additive")`. What can you say about the trend component? What about the random component? Does the random component look random? Or does it appear to still have some seasonality on it?

```
ts_solar <- ts(energy_data[,2], frequency = 12, start = c(1984, 1))
ts_wind <- ts(energy_data[,3], frequency = 12, start = c(1984, 1))

decomposed_ts_solar_add <- decompose(ts_solar,"additive")
decomposed_ts_wind_add <- decompose(ts_wind,"additive")

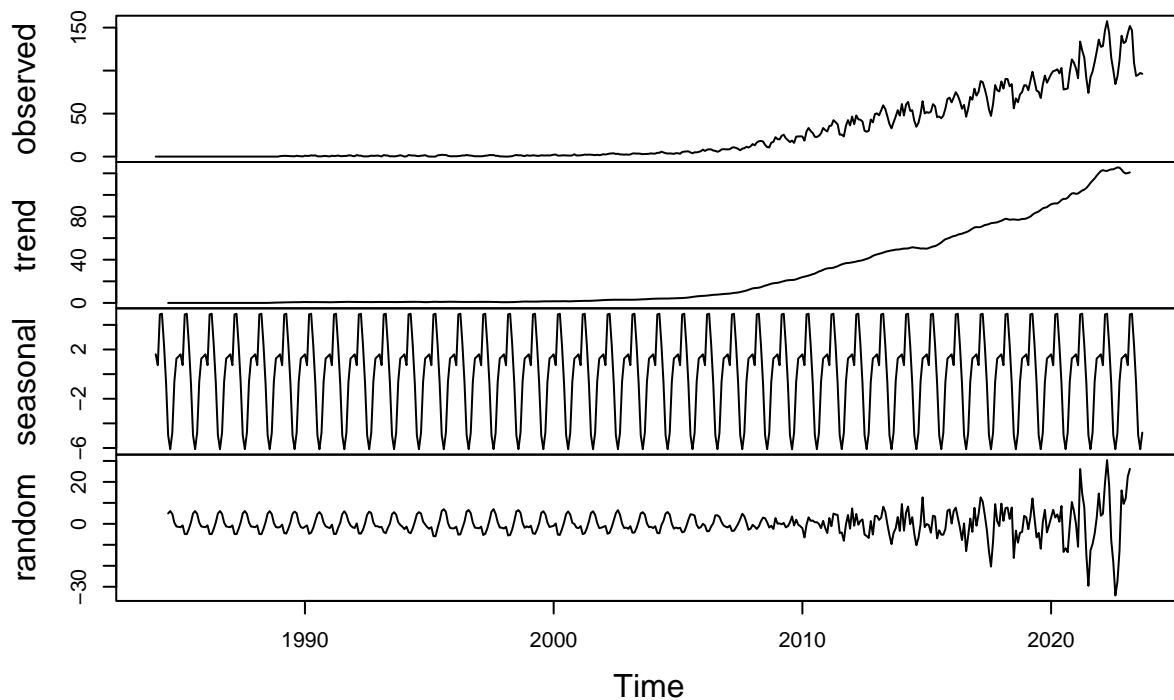
plot(decomposed_ts_solar_add)
```

Decomposition of additive time series



```
plot(decomposed_ts_wind_add)
```


Decomposition of additive time series



*#For solar, the trend component is present and steadily increasing over time.
#The random component doesn't look that random and appears to still have some
#seasonality.*

*#For wind, the trend component is present and steadily increasing over time.
#The random component doesn't look that random and appears to still have some
#seasonality.*

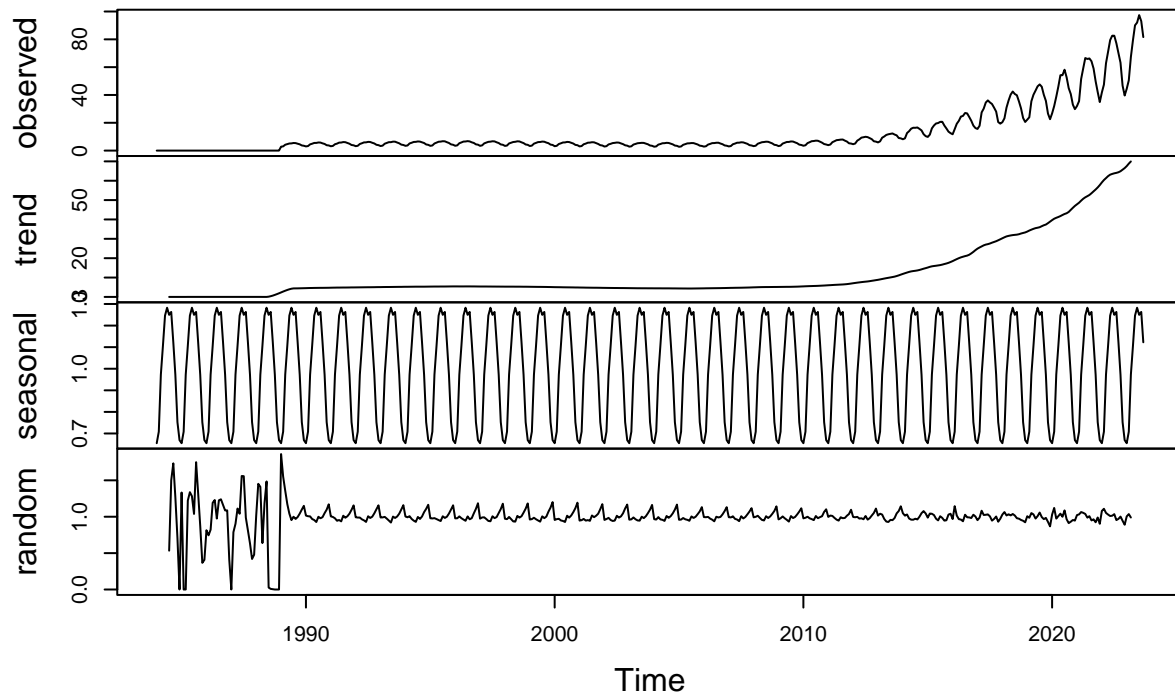
Q5

Use the `decompose` function again but now change the type of the seasonal component from additive to multiplicative. What happened to the random component this time?

```
decomposed_ts_solar_mult <- decompose(ts_solar,"multiplicative")
decomposed_ts_wind_mult <- decompose(ts_wind,"multiplicative")

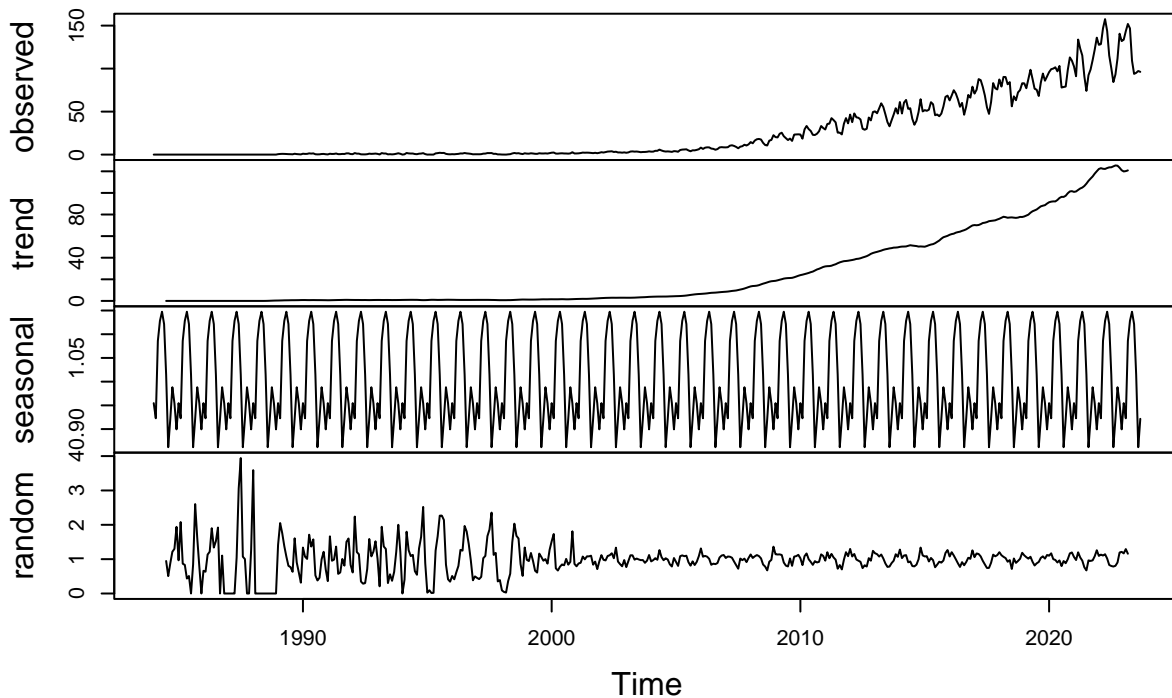
plot(decomposed_ts_solar_mult)
```

Decomposition of multiplicative time series



```
plot(decomposed_ts_wind_mult)
```

Decomposition of multiplicative time series



*#For solar, the random component now appears more random at the outset
#(before 1990) and then turns into a pattern that looks somewhat seasonal
#in nature.*

*#For wind, the random component now appears more random at the outset
#(before 2000) and then turns into a pattern that looks somewhat seasonal
#in nature.*

Q6

When fitting a model to this data, do you think you need all the historical data? Think about the data from 90s and early 20s. Are there any information from those years we might need to forecast the next six months of Solar and/or Wind consumption. Explain your response.

Answer: Solar and wind resource / consumption has increased dramatically since the 90s and early 20s. Therefore, forecasting the next sixth months 20+ years later, it is unlikely that we need any information from those years.

Q7

Create a new time series object where historical data starts on January 2012. Hint: use `filter()` function so that you don't need to point to row numbers, i.e, `filter(yyyy, year(Date) >= 2012)`. Apply the `decompose` function `type=additive` to this new time series. Comment the results. Does the random component look random? Think about our discussion in class about seasonal components that depends on the level of the series.

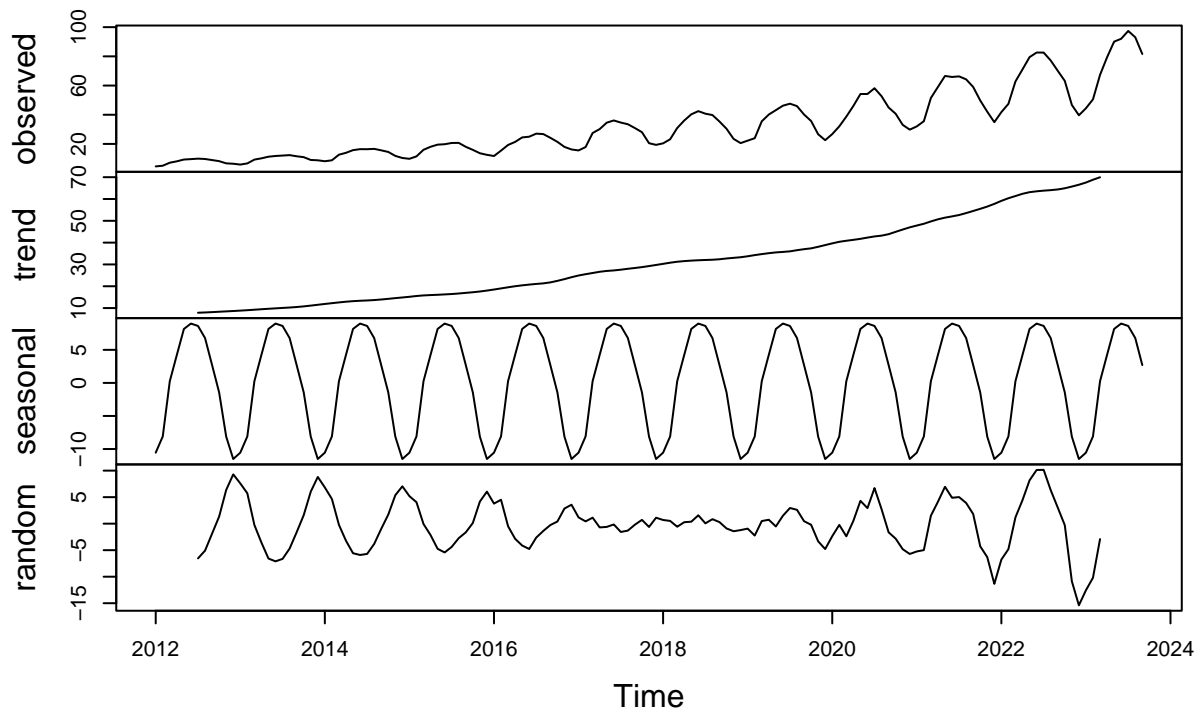
```
energy_data_2012 <- filter(energy_data, year(energy_data[,1]) >= 2012)

ts_solar_2012 <- ts(energy_data_2012[,2], frequency = 12, start = c(2012, 1))
ts_wind_2012 <- ts(energy_data_2012[,3], frequency = 12, start = c(2012, 1))

decomposed_ts_solar_2012_add <- decompose(ts_solar_2012,"additive")
decomposed_ts_wind_2012_add <- decompose(ts_wind_2012,"additive")

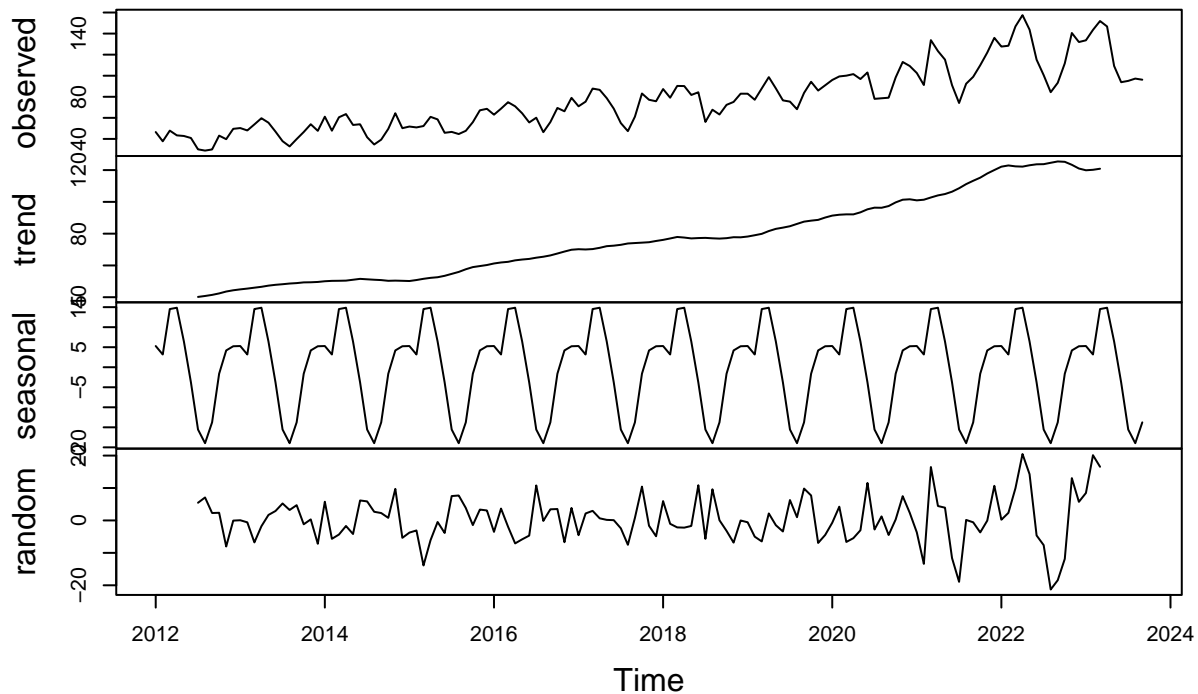
plot(decomposed_ts_solar_2012_add)
```

Decomposition of additive time series



```
plot(decomposed_ts_wind_2012_add)
```

Decomposition of additive time series



Answer: #For solar, the random component does not look that random. There is a clear wave-like pattern between 2012 and 2017, then again from 2019 to 2023. #For wind, the random component looks a bit more random. There is a bit of a wavelike pattern, but it appears more random in nature with less regularities in wave pattern.

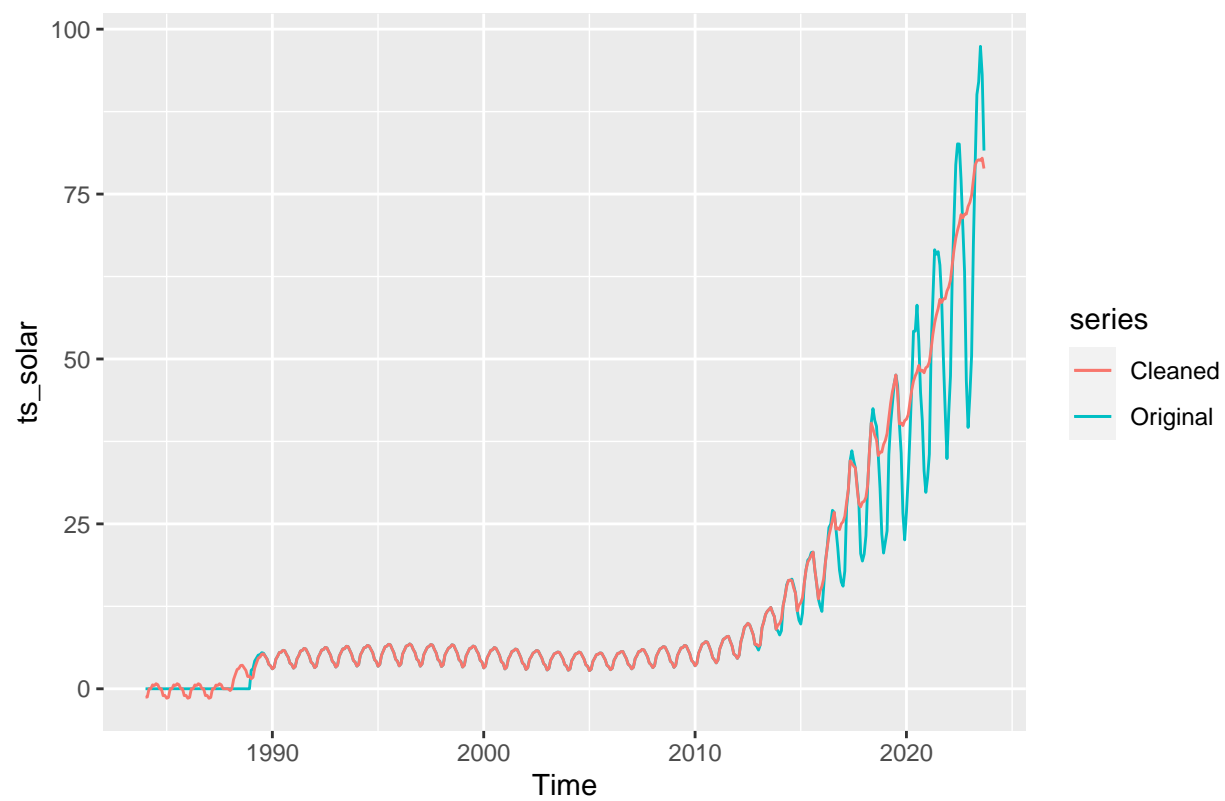
Identify and Remove outliers

Q8

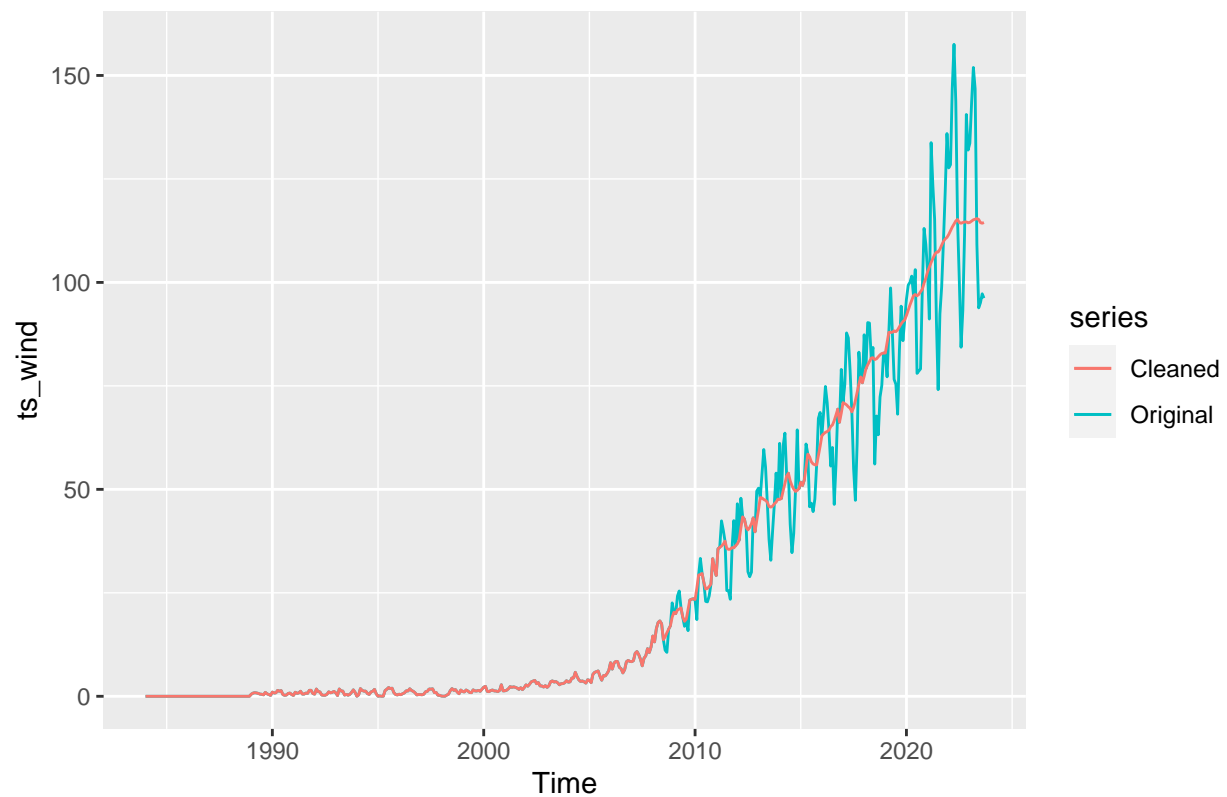
Apply the `tsclean()` to both series from Q7. Did the function removed any outliers from the series? Hint: Use `autoplot()` to check if there is difference between cleaned series and original series.

```
clean_ts_solar <- tsclean(ts_solar)
clean_ts_wind <- tsclean(ts_wind)

autoplot(ts_solar, series = "Original")+
  autolayer(clean_ts_solar, series = "Cleaned")
```



```
autoplot(ts_wind, series = "Original")+  
  autolayer(clean_ts_wind, series = "Cleaned")
```



*#For solar, the function did remove outliers from the series starting in 2013
#and steadily removing more over time.*

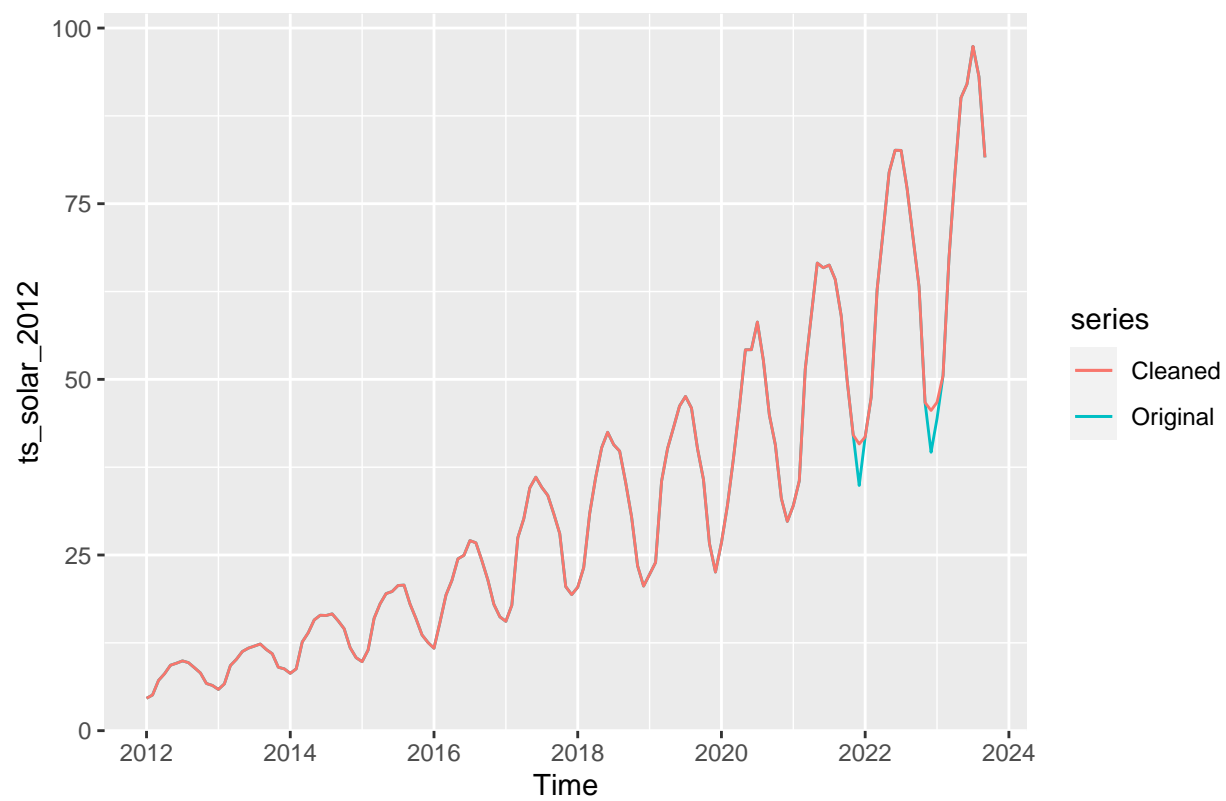
*#For wind, the function did remove outliers starting in 2009 and steadily
#removing more over time.*

Q9

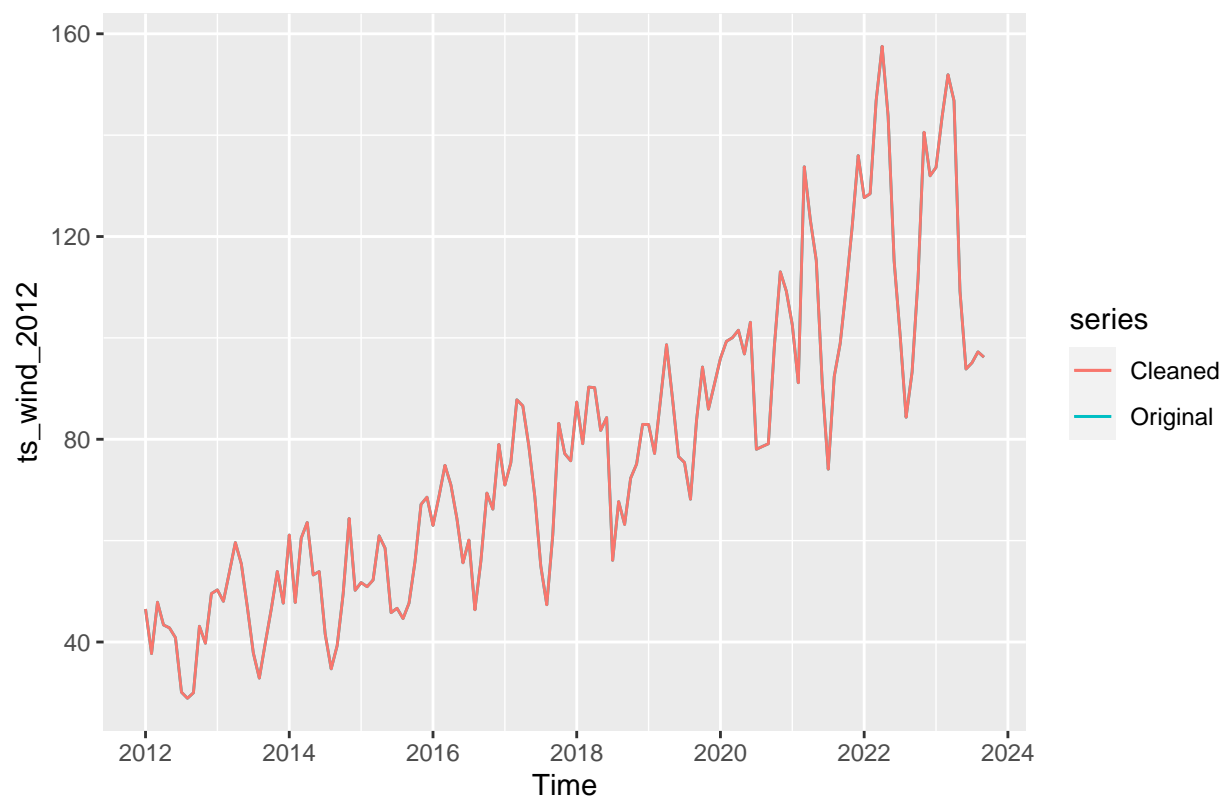
Redo number Q8 but now with the time series you created on Q7, i.e., the series starting in 2014. Using what `autoplot()` again what happened now? Did the function removed any outliers from the series?

```
clean_ts_solar_2012 <- tsclean(ts_solar_2012)
clean_ts_wind_2012 <- tsclean(ts_wind_2012)

autoplot(ts_solar_2012, series = "Original")+
  autolayer(clean_ts_solar_2012, series = "Cleaned")
```



```
autoplot(ts_wind_2012, series = "Original")+  
  autolayer(clean_ts_wind_2012, series = "Cleaned")
```

Answer: For solar, the function did remove some outliers from the series around 2022 and 2023.
For wind, the function did not remove any outliers.