Time Series Analysis of Rainfall in the Rio Grande Valley, Texas, 1941-present

John Knight

Final Project – STAT 6380

07/09/2024

Introduction

The Earth is getting hotter,¹ but how does this affect rainfall? An increase in temperature could make an area more desert-like,² or it could increase severe weather events such as thunderstorms and hurricanes.³ In this study, a time series of precipitation levels over the past century in the Rio Grande Valley, Texas, is analyzed.

The first objective of this study is to assess whether there is a discernable trend in annual precipitation levels across time. The second objective is to show seasonal patterns in precipitation in the Rio Grande Valley. The third objective is to investigate autocorrelation in precipitation totals between consecutive months and years, and the final objective is to use spectral analysis to reveal the underlying frequencies of rainfall patterns.

Methods

Daily precipitation amounts, in inches, were downloaded from the free archive of the National Centers for Environmental Information⁴ (NCEI) for four weather stations: McAllen, McAllen Airport, Weslaco, and Harlingen (see Table 1). The start date of 1941-06-01 was chosen because this was the earliest available data for McAllen (both Harlingen and Weslaco data go back further). Missing values were imputed by taking the adjusted mean of the other stations on each date. Note that McAllen Airport data was only used for imputation purposes, and was omitted from further analysis to avoid overweighting the McAllen area relative to the other two stations. The daily mean precipitation was then calculated for each date, and a monthly mean was calculated as the average of the daily means for each month (to avoid bias toward months with more days).

Station	Start Date	Total Rows	NA Rows	Coverage (%)	Mean Prcp.
Harlingen	1941-06-01	30331	1809	94.0	0.0728
McAllen	1941-06-01	30331	2459	91.9	0.0571
McAllen Airport	1961-01-14	23164	597	97.4	0.0594
Weslaco	1941-06-01	30331	3484	88.5	0.0667

Table 1. Coverage and mean daily precipitation for the four weather stations
--

All statistical analysis was performed using R software. The time series was decomposed into its trend, seasonal, and remainder components using STL decomposition. A range of trend window parameters were tested to find an appropriate smoothing balance. Boxplots were created to visualize the seasonal distribution by month. The remainder component was then assessed for autocorrelation by inspection of the correlogram. The strengths of the trend and seasonality were measured using the following formulae:

$$F_t = \max\left(0, 1 - \frac{VAR(Z_t)}{VAR(M_t + Z_t)}\right) \text{ with range } [0,1] \text{ (no trend to very strong trend).}$$

$$F_s = \max\left(0, 1 - \frac{VAR(Z_t)}{VAR(S_t + Z_t)}\right) \text{ with range } [0,1] \text{ (no seasonality to very strong seasonality).}$$

The autocorrelation function (ACF) and partial ACF of the series were inspected, and an ARIMA model was fitted to the data with the auto.arima function using parameters approximation=FALSE and stepwise=FALSE, thereby instructing a more thorough search for appropriate models. The residuals of the resulting model underwent a visual inspection as well as

a Ljung-Box test for autocorrelation of residuals. This model was then used to create a forecast of precipitation for the next 5 years (60 months).

Finally, spectral analysis was performed on the time series. A smoothed periodogram was created using the Daniell kernel, with various sized windows tested to find an appropriate level of smoothing. Confidence intervals were calculated for frequencies of prominent peaks in the resulting periodogram to add context to their significance.







The time series is summarized in Figure 1, which shows the raw series, plus the trend, seasonal, and remainder components. A trend window of 201 was used to achieve the smooth trend seen above – smaller values resulted in more fluctuation. The trend can be seen to peak in the 1970s, decreasing thereafter before rising again in the past 10 years. The strength of the trend was found to be very weak, 0.012, and the strength of the seasonality was 0.215. In the correlogram of the remainders, there was significant autocorrelation at lag 1, while in the correlogram of the raw series, there was clear evidence of seasonal autocorrelation at lags 12 and 24 (see Figures 4 and 5 in Appendix).

The seasonal boxplots in Figure 2 reveal that September is the month with the highest precipitation on average. May and June are also months with high precipitation. The lowest rainfall occurs from November through to April.

The auto.arima function found the optimal model to be ARIMA(1,0,0)(2,0,0)[12]. This aligns with earlier inspection of the correlograms which showed autocorrelation at lags 1, 12, and 24. The coefficients of the ARIMA model were AR1 = 0.127, SAR1 = 0.223, SAR2 = 0.141, with a

mean of 0.064. The Ljung-Box test for autocorrelation returned a p-value of 0.267, indicating that there is no significant autocorrelation in the residuals. The Akaike Information Criterion (AIC) was -2656.61. This model was used to forecast the next 5 years with a relatively flat trend and prediction intervals shown in Figure 6 in the Appendix.



Figure 2. Boxplot displaying monthly distribution of precipitation.

The periodogram for the spectral analysis can be seen in Figure 3. A Daniell kernel with window size (13, 13) and tapering of 0.1 was found to be suitable for a smooth and clear spectrum. The most prominent peak, as expected, was at frequency 1 (annual) but there was also a prominent peak at frequency 3. Smaller peaks were observed at frequencies 2, 4, and 5. Confidence intervals were calculated for frequencies 1, 2, 3, 4, and 5 and they can be seen along with the estimated spectral densities in Table 2.

Frequency	Spectral Density	95% CI
1	0.000396	(0.000211, 0.000639)
2	0.000441	(0.000235, 0.000711)
3	0.000362	(0.000193, 0.000584)
4	0.000249	(0.000132, 0.000401)
5	0.000253	(0.000135, 0.000408)

Table 2. Spectral density estimates and 95% confident intervals for peak values in periodogram.



Discussion

The study found no significant evidence of a long-term change in the trend of precipitation in the Rio Grande Valley. However, the autocorrelation in the ARIMA model suggested the presence of cycles in rainfall patterns. These may be related to various weather phenomena such as El Niño and La Niña.⁵ The highest seasonal rainfall in September evidently coincides with the peak of the Atlantic hurricane season.⁶ The prominent peak at frequency 1 in the spectral analysis is not surprising, since the Earth's climate is known to experience annual seasons. However, further peaks at frequencies 2, 3, 4, and 5 suggest more complex underlying systems that interact throughout the year. This helps to explain why the seasonal totals from January to December to not resemble a simple curve, as the equivalent measures for temperature would do.

The strength of this study is that it leverages raw data, direct from source, to analyze the climate in a specific region of the United States whose climate is significantly hotter than other areas. This method can uncover more accurate insights than simply by looking at national or global climate summaries. The use of ARIMA provides a robust model to detect autocorrelation among observations in a time series, helping to explain monthly and yearly changes in weather patterns.

A possible limitation is the unknown reliability of the data. An equivalent dataset for temperature was explored, and it contained many outliers which were obvious errors (for example, a temperature of 0° F or a temperature vastly different from a neighboring city). However, these errors are harder to detect in precipitation, because a value of 0 is always valid, and it is plausible for two cities a short distance apart to experience drastically different rainfall on a given day. Some possible avenues for further exploration of this topic would be to look at the number of major rain events, rather than the total precipitation, or a multivariate analysis comparing temperature and precipitation.

References

- 1. NASA Earth Observatory. (2023, January 13). 2022 tied for fifth warmest year on record. NASA. Retrieved from <u>https://earthobservatory.nasa.gov/images/150828/2022-tied-for-fifth-warmest-year-on-record</u>
- Doyle, A. (2016, October 27). Global warming to turn parts of Europe into desert by end of the century, report warns. The Independent. Retrieved from https://www.independent.co.uk/climate-change/news/global-warming-desert-europeclimate-change-a7375276.html
- 3. NASA Earth Observatory. (2023, January 13). *How climate change may be impacting storms over Earth's tropical oceans*. NASA. Retrieved from https://science.nasa.gov/earth-science/oceanography/living-ocean/hurricanes-and-climate-change
- 4. National Centers for Environmental Information (NCEI). (n.d.). Data access. NOAA. Retrieved July 2, 2024, from https://www.ncei.noaa.gov/access/search/
- 5. National Oceanic and Atmospheric Administration (NOAA). (2024, June 16). What are El Niño and La Niña? NOAA. Retrieved from https://oceanservice.noaa.gov/facts/ninonina.html
- 6. National Oceanic and Atmospheric Administration (NOAA). (2016, August 22). The peak of the hurricane season – why now? NOAA. Retrieved from https://www.noaa.gov/stories/peak-of-hurricane-season-why-now

<u>Appendix</u>



Figure 4. Correlogram showing ACF or remainder component after trend and seasonality have been removed.





Figure 6. Forecast of precipitation for next 5 years (60 months).

R Code

John Knight - STAT 6380

Final Project - RGV Climate

options(scipen = 999) # To avoid scientific notation in plots

```
# Read the CSV file
weather_data <- read.csv("weather_data.csv")</pre>
```

Remove the STATION column
weather data <- weather data %>% select(-STATION)

Convert DATE column to Date type weather data\$DATE <- mdy(weather data\$DATE)

Filter out rows where DATE is before 1941-06-01
weather_data <- weather_data %>% filter(DATE >= as.Date("1941-06-01"))

Pivot the data so that the stations are columns

pivoted_data <- weather_data %>%

pivot_wider(names_from = STATION_NAME, values_from = c(PRCP, TMAX)) %>%
rename_with(~gsub(" ", "_", .))

Data cleaning and imputation

1. Detect outliers in TMAX and set them to NA.

Detect outliers in TMAX. For each station, get diff with mean of other three stations. pivoted_data\$McAllen_Airport_Diff <- pivoted_data\$TMAX_McAllen_Airport rowMeans(pivoted_data[, c("TMAX_Weslaco", "TMAX_Harlingen", "TMAX_McAllen")], na.rm = TRUE) pivoted_data\$McAllen_Diff <- pivoted_data\$TMAX_McAllen - rowMeans(pivoted_data[, c("TMAX_Weslaco", "TMAX_Harlingen", "TMAX_McAllen_Airport")], na.rm = TRUE) pivoted_data\$Weslaco_Diff <- pivoted_data\$TMAX_Weslaco - rowMeans(pivoted_data[, c("TMAX_McAllen", "TMAX_Harlingen", "TMAX_McAllen_Airport")], na.rm = TRUE) pivoted_data\$Herlingen_Diff <- pivoted_data\$TMAX_McAllen_Airport")], na.rm = TRUE) pivoted_data\$Harlingen_Diff <- pivoted_data\$TMAX_Harlingen - rowMeans(pivoted_data[, c("TMAX_Weslaco", "TMAX_McAllen", "TMAX_McAllen_Airport")], na.rm = TRUE) pivoted_data\$Herlingen_Diff <- pivoted_data\$TMAX_Harlingen - rowMeans(pivoted_data[, c("TMAX_Weslaco", "TMAX_McAllen", "TMAX_McAllen_Airport")], na.rm = TRUE)

If yes, see how many diffs are +/# If 2v2 or 1v1, ignore
If 3v1 or 2v1, turn the 1 into NA
n <- dim(pivoted_data)[1]
for (i in 1:n) {
 diffs <- pivoted_data[i, c("Weslaco_Diff", "Harlingen_Diff", "McAllen_Airport_Diff",
 "McAllen_Diff")]</pre>

Check if any absolute value >= 20
if (any(abs(diffs) >= 20, na.rm = TRUE)) {
 # Count negative and nonnegative values
 neg_count <- sum(diffs < 0, na.rm = TRUE)</pre>

```
pos count \leq sum(diffs \geq 0, na.rm = TRUE)
  if (neg count != pos count) {
   if (pos count > neg count) {
    # More positive than negative, turn the negative value's corresponding TMAX to NA
    if (!is.na(diffs["Weslaco_Diff"]) && diffs["Weslaco Diff"] < 0) pivoted data[i,
"TMAX Weslaco"] <- NA
    if (!is.na(diffs["Harlingen Diff"]) && diffs["Harlingen Diff"] < 0) pivoted data[i,
"TMAX Harlingen"] <- NA
    if (!is.na(diffs["McAllen Airport Diff"]) && diffs["McAllen Airport Diff"] < 0)
pivoted data[i, "TMAX McAllen Airport"] <- NA
    if (!is.na(diffs["McAllen Diff"]) && diffs["McAllen Diff"] < 0) pivoted data[i,
"TMAX McAllen"] <- NA
   } else {
    # More negative than positive, turn the positive value's corresponding TMAX to NA
    if (!is.na(diffs["Weslaco Diff"]) && diffs["Weslaco Diff"] >= 0) pivoted data[i,
"TMAX Weslaco"] <- NA
    if (!is.na(diffs["Harlingen Diff"]) && diffs["Harlingen Diff"] >= 0) pivoted data[i,
"TMAX Harlingen"] <- NA
    if (!is.na(diffs["McAllen Airport Diff"]) && diffs["McAllen Airport Diff"] >= 0)
pivoted data[i, "TMAX McAllen Airport"] <- NA
    if (!is.na(diffs["McAllen_Diff"]) && diffs["McAllen_Diff"] >= 0) pivoted_data[i,
"TMAX McAllen"] <- NA
   }
  }
 }
}
```

2. Get the coverage % for each of the 4 stations for PRCP and TMAX.

```
# Function to get column statistics for a single column
get_column_stats <- function(data, column_name) {
    # Get total rows in the original data
    total rows <- nrow(data)</pre>
```

```
# Find the earliest non-NA date
non_na_rows <- data[!is.na(data[[column_name]]), ]
start_date <- min(non_na_rows$DATE, na.rm = TRUE)</pre>
```

```
# Filter the data from the start date onwards
filtered_data <- data[data$DATE >= start_date, ]
filtered_total_rows <- nrow(filtered_data)
na_count <- sum(is.na(filtered_data[[column_name]]))
non na percentage <- (1 - na count / filtered total rows) * 100</pre>
```

data.frame(

```
Column = column name,
  Start Date = start date,
  Total Rows = filtered total rows,
  NA Count = na count,
  Non NA Percentage = round(non na percentage, 2)
 )
}
# Specify columns to analyze
columns to analyze <- c("TMAX Harlingen", "TMAX McAllen", "TMAX McAllen Airport",
"TMAX Weslaco",
             "PRCP Harlingen", "PRCP McAllen", "PRCP McAllen Airport",
"PRCP Weslaco")
coverage <- data.frame()</pre>
# Loop through columns and get stats
for (col in columns to analyze) {
 col stats <- get column stats(pivoted data, col)
 coverage <- rbind(coverage, col stats)
}
print(coverage)
# 3. Find rows with no PRCP or no TMAX for any station.
library(dplyr)
pivoted data <- pivoted data %>%
 mutate(
  TMAX all na = ifelse(
   rowSums(is.na(select(., TMAX Harlingen, TMAX McAllen, TMAX McAllen Airport,
TMAX_Weslaco))) == 4,
   1,
   0
  ),
  PRCP all na = ifelse(
   rowSums(is.na(select(., PRCP Harlingen, PRCP McAllen, PRCP McAllen Airport,
PRCP Weslaco))) == 4,
   1,
   0
  )
 )
```

4. Get the mean PRCP and TMAX for each station, only using rows where all 4 have values.

```
mean prcp <- pivoted data %>%
filter(!is.na(PRCP_Harlingen) &
      !is.na(PRCP McAllen) &
      !is.na(PRCP McAllen Airport) &
     !is.na(PRCP Weslaco)) %>%
 summarise(
  PRCP Harlingen = mean(PRCP Harlingen, na.rm = TRUE),
  PRCP McAllen = mean(PRCP McAllen, na.rm = TRUE),
  PRCP McAllen Airport = mean(PRCP McAllen Airport, na.rm = TRUE),
  PRCP Weslaco = mean(PRCP Weslaco, na.rm = TRUE)
 )
mean prcp df <- as.data.frame(mean prcp)
mean prcp df \leq- data.frame(lapply(mean prcp df, function(x) sprintf("%.4f", x)))
print(mean prcp df)
mean temps <- pivoted data %>%
 filter(!is.na(TMAX Harlingen) &
      !is.na(TMAX McAllen) &
     !is.na(TMAX McAllen Airport) &
      !is.na(TMAX Weslaco)) %>%
 summarise(
  TMAX Harlingen = mean(TMAX Harlingen, na.rm = TRUE),
  TMAX McAllen = mean(TMAX McAllen, na.rm = TRUE),
  TMAX McAllen Airport = mean(TMAX McAllen Airport, na.rm = TRUE),
  TMAX Weslaco = mean(TMAX Weslaco, na.rm = TRUE)
 )
```

mean_temps_df <- as.data.frame(mean_temps)
mean_temps_df <- data.frame(lapply(mean_temps_df, function(x) sprintf("%.3f", x)))</pre>

print(mean_temps_df)

5. Impute missing values of PRCP and TMAX based on available values, with an adjustment based on means.

Calculate the overall mean for precipitation and temperature
overall_mean_prcp <- mean(as.numeric(mean_prcp_df), na.rm = TRUE)
overall_mean_temps <- mean(as.numeric(mean_temps_df), na.rm = TRUE)</pre>

Calculate the ratio adjustments for PRCP
prcp_adjustments <- as.numeric(mean_prcp_df) / overall_mean_prcp
temp_adjustments <- as.numeric(mean_temps_df) - overall_mean_temp</pre>

```
# Rename adjustments for easier reference
names(prcp adjustments) <- names(mean prcp df)
names(temp adjustments) <- names(mean temps df)
# Function to adjust PRCP values by ratio
adjust prcp values <- function(value, adjustment) {
 if (is.na(value)) {
  return(NA)
 } else {
  return(value / adjustment)
 }
}
# Function to adjust TMAX values by difference
adjust tmax values <- function(value, adjustment) {
 if (is.na(value)) {
  return(NA)
 } else {
  return(value - adjustment)
 }
}
# Apply the adjustments to each row
pivoted data <- pivoted data %>%
 rowwise() %>%
 mutate(
  Adjusted PRCP McAllen Airport = adjust prcp values(PRCP McAllen Airport,
prcp adjustments["PRCP McAllen Airport"]),
  Adjusted PRCP Harlingen = adjust prcp values(PRCP Harlingen,
prcp adjustments["PRCP Harlingen"]),
  Adjusted PRCP McAllen = adjust prcp values(PRCP McAllen,
prcp adjustments["PRCP McAllen"]),
  Adjusted PRCP Weslaco = adjust prcp values(PRCP Weslaco,
prcp adjustments["PRCP Weslaco"]),
  Adjusted TMAX McAllen Airport = adjust tmax values(TMAX McAllen Airport,
temp adjustments["TMAX McAllen Airport"]),
  Adjusted TMAX Harlingen = adjust tmax values(TMAX Harlingen,
temp adjustments["TMAX Harlingen"]),
  Adjusted TMAX McAllen = adjust tmax values(TMAX McAllen,
temp adjustments["TMAX McAllen"]),
  Adjusted TMAX Weslaco = adjust tmax values(TMAX Weslaco,
temp adjustments["TMAX Weslaco"])
 )%>%
 ungroup()
```

Calculate the mean of the non-NA values for the adjusted PRCP and TMAX columns mean adjusted prcp <- pivoted data %>% select(Adjusted PRCP Harlingen, Adjusted PRCP McAllen, Adjusted PRCP McAllen Airport, Adjusted PRCP Weslaco) %>% summarise(across(everything(), ~mean(., na.rm = TRUE))) mean adjusted tmax <- pivoted data %>% select(Adjusted TMAX Harlingen, Adjusted TMAX McAllen, Adjusted TMAX McAllen Airport, Adjusted TMAX Weslaco) %>% summarise(across(everything(), ~mean(., na.rm = TRUE))) # Create the imputed columns based on row means of non-NA adjusted values pivoted data <- pivoted data %>% rowwise() %>% mutate(Imputed PRCP Harlingen = ifelse(is.na(Adjusted PRCP Harlingen), ifelse(all(is.na(c(Adjusted PRCP Harlingen, Adjusted PRCP McAllen, Adjusted PRCP McAllen Airport, Adjusted PRCP Weslaco))), NA, mean(c(Adjusted PRCP McAllen, Adjusted PRCP McAllen Airport, Adjusted PRCP Weslaco), na.rm = TRUE)), Adjusted PRCP Harlingen), Imputed PRCP McAllen = ifelse(is.na(Adjusted PRCP McAllen), ifelse(all(is.na(c(Adjusted PRCP Harlingen, Adjusted PRCP McAllen, Adjusted PRCP McAllen Airport, Adjusted PRCP Weslaco))), NA, mean(c(Adjusted PRCP Harlingen, Adjusted PRCP McAllen Airport, Adjusted PRCP Weslaco), na.rm = TRUE)), Adjusted PRCP McAllen), Imputed PRCP McAllen Airport = ifelse(is.na(Adjusted PRCP McAllen Airport), ifelse(all(is.na(c(Adjusted PRCP Harlingen, Adjusted PRCP McAllen, Adjusted PRCP McAllen Airport, Adjusted PRCP Weslaco))), NA, mean(c(Adjusted PRCP Harlingen, Adjusted PRCP McAllen, Adjusted PRCP Weslaco), na.rm = TRUE)), Adjusted PRCP McAllen Airport), Imputed PRCP Weslaco = ifelse(is.na(Adjusted PRCP Weslaco), ifelse(all(is.na(c(Adjusted PRCP Harlingen, Adjusted PRCP McAllen, Adjusted PRCP McAllen Airport, Adjusted PRCP Weslaco))), NA, mean(c(Adjusted PRCP Harlingen, Adjusted PRCP McAllen, Adjusted PRCP McAllen Airport), na.rm = TRUE)), Adjusted_PRCP_Weslaco), Imputed TMAX Harlingen = ifelse(is.na(Adjusted TMAX Harlingen), ifelse(all(is.na(c(Adjusted TMAX Harlingen, Adjusted TMAX McAllen, Adjusted TMAX McAllen Airport, Adjusted TMAX Weslaco))), NA, mean(c(Adjusted TMAX McAllen, Adjusted TMAX McAllen Airport, Adjusted TMAX Weslaco), na.rm = TRUE)),

```
Adjusted TMAX Harlingen),
  Imputed TMAX McAllen = ifelse(is.na(Adjusted TMAX McAllen),
                 ifelse(all(is.na(c(Adjusted TMAX Harlingen, Adjusted TMAX McAllen,
Adjusted TMAX McAllen Airport, Adjusted TMAX Weslaco))), NA,
                     mean(c(Adjusted TMAX Harlingen,
Adjusted TMAX McAllen Airport, Adjusted TMAX Weslaco), na.rm = TRUE)),
                  Adjusted TMAX McAllen),
  Imputed TMAX McAllen Airport = ifelse(is.na(Adjusted TMAX McAllen Airport),
                      ifelse(all(is.na(c(Adjusted TMAX Harlingen,
Adjusted TMAX McAllen, Adjusted TMAX McAllen Airport, Adjusted TMAX Weslaco))),
NA,
                         mean(c(Adjusted TMAX Harlingen, Adjusted TMAX McAllen,
Adjusted TMAX Weslaco), na.rm = TRUE)),
                      Adjusted TMAX McAllen Airport),
  Imputed TMAX Weslaco = ifelse(is.na(Adjusted TMAX Weslaco),
                 ifelse(all(is.na(c(Adjusted TMAX Harlingen, Adjusted TMAX McAllen,
Adjusted TMAX McAllen Airport, Adjusted TMAX Weslaco))), NA,
                     mean(c(Adjusted TMAX Harlingen, Adjusted TMAX McAllen,
Adjusted TMAX McAllen Airport), na.rm = TRUE)),
                 Adjusted TMAX Weslaco)
 )%>%
 ungroup()
```

TIME SERIES ANALYSIS

First, get the mean for each day, then the mean for each month and make time series. # Ignore rows where all values are NA

```
# Step 1: Calculate row-wise means excluding McAllen Airport
pivoted_data <- pivoted_data %>%
rowwise() %>%
mutate(
    Mean_PRCP = ifelse(
    all(is.na(c(Imputed_PRCP_Harlingen, Imputed_PRCP_McAllen,
Imputed_PRCP_Weslaco))),
    NA,
    mean(c(Imputed_PRCP_Harlingen, Imputed_PRCP_McAllen, Imputed_PRCP_Weslaco),
na.rm = TRUE)
    ),
    Mean_TMAX = ifelse(
    all(is.na(c(Imputed_TMAX_Harlingen, Imputed_TMAX_McAllen,
Imputed_TMAX_Weslaco))),
    NA,
```

```
mean(c(Imputed TMAX Harlingen, Imputed TMAX McAllen,
Imputed TMAX Weslaco), na.rm = TRUE)
  )
 )%>%
 ungroup()
# Step 2: Create a time series for monthly mean PRCP and TMAX
# Convert DATE to yearmon format for aggregation by month
pivoted data$YearMonth <- as.yearmon(pivoted data$DATE)
# Aggregate to calculate the mean PRCP and TMAX for each month
monthly data <- pivoted data %>%
 group by(YearMonth) %>%
 summarise(
  Monthly Mean PRCP = mean(Mean PRCP, na.rm = TRUE),
  Monthly Mean TMAX = mean(Mean TMAX, na.rm = TRUE)
 ) %>%
 ungroup()
# Create time series objects
monthly prcp ts <- ts(monthly data$Monthly Mean PRCP, start =
c(year(min(pivoted data$DATE)), month(min(pivoted data$DATE))), frequency = 12)
monthly tmax ts <- ts(monthly data$Monthly Mean TMAX, start =
c(year(min(pivoted data$DATE)), month(min(pivoted data$DATE))), frequency = 12)
# Plot the time series
autoplot(monthly prcp ts) +
 labs(caption="Caption", x="Year", y="Mean daily precipitation (inches)")
# Create boxplots to view seasons
monthly prcp df <- data.frame(
 Date = as.Date(time(monthly prcp ts)),
 Precipitation = as.numeric(monthly prcp ts)
)
monthly prcp df$Month <- factor(format(monthly prcp df$Date, "%b"), levels = month.abb)
ggplot(monthly prcp df, aes(x = Month, y = Precipitation)) +
 geom boxplot(fill='lightblue') +
 labs(x = "Month",
    y = "Average daily precipitation (inches)",
    caption="Figure 2. Boxplot displaying monthly distribution of precipitation.")
# Decompose the monthly PRCP time series using STL
stl prcp <- stl(monthly prcp ts, s.window = "periodic", t.window=201) # Parameters need to be
odd
autoplot(stl prcp) +
```

labs(caption="Figure 1. STL decomposition of monthly mean precipitation for Rio Grande Valley, 1941-2024.")

Assess autocorrelation among the random component after removing trend and seasonality. remainder_component <- stl_prcp\$time.series[, "remainder"] # Extract the remainder component ggAcf(remainder_component) +

labs(title="",

caption="Figure 4. Correlogram showing ACF or remainder component after trend and seasonality have been removed.")

Measure strength of trend and seasonality (final slide of lecture 2)

trend_component <- stl_prcp\$time.series[, "trend"] # Extract components
seasonal_component <- stl_prcp\$time.series[, "seasonal"]
remainder_component <- stl_prcp\$time.series[, "remainder"]</pre>

var_remainder <- var(remainder_component, na.rm = TRUE) # Calculate variances var_trend_plus_remainder <- var(trend_component + remainder_component, na.rm = TRUE) var_seasonal_plus_remainder <- var(seasonal_component + remainder_component, na.rm = TRUE)

Ft <- max(0, 1 - (var_remainder / var_trend_plus_remainder)) # Calculate strength of trend

Fs <- max(0, 1 - (var_remainder / var_seasonal_plus_remainder)) # Calculate strength of seasonality

cat("Strength of Trend (Ft):", Ft, "\n") cat("Strength of Seasonality (Fs):", Fs, "\n")

ACF and partial ACF
ggAcf(monthly_prcp_ts) +
labs(title="")
ggPacf(monthly_prcp_ts) +
labs(title = "", caption = "Figure 5. Partial autocorrelation function of monthly precipitation.")

ARIMA

ndiffs and nsdiffs to determine the number of differences
ndiffs(monthly_prcp_ts) # 0
nsdiffs(monthly_prcp_ts) # 0

my_arima <- auto.arima(monthly_prcp_ts, approximation = FALSE, stepwise = FALSE) # Parameters cause slower, more thorough search.

Check residuals of ARIMA model by plotting residuals and Portmanteau test (Ljung-Box) residuals <- residuals(my_arima)

```
ggAcf(residuals) +
 labs(title = "ACF of Residuals")
ggPacf(residuals) +
 labs(title = "PACF of Residuals")
autoplot(residuals) +
 labs(title = "Residuals of ARIMA Model", x = "Time", y = "Residuals")
Box.test(residuals, lag = 20, type = "Ljung-Box")
# Forecast next 5 years
forecast horizon <- 60
forecast <- forecast(my arima, h = forecast horizon)
# Filter the time series to include only data from 2010 onwards for plotting
start year <- 2010
filtered ts <- window(monthly prcp ts, start = c(start year, 1))
# Plot the filtered time series and the forecast
autoplot(filtered ts, series = "Observed", color='darkgrey') +
 autolayer(forecast, series = "Forecast") +
 labs(title = "",
    x = "Year"
    y = "Mean Daily Precipitation (inches)",
    caption = "Figure 6. Forecast of precipitation for next 5 years (60 months).")
## Spectral Analysis
spec <- spec.pgram(monthly prcp ts, plot = FALSE) # Periodogram
plot(spec, main = "Periodogram of Monthly Precipitation", xlab = "Frequency", ylab = "Spectral
Density")
# Smoothed periodogram with Daniell kernel
spec_smooth <- spec.pgram(monthly prcp ts, spans = c(13, 13), taper = 0.1, log = "no")
# Extract the frequency and spectral density
freq <- spec smooth$freq
spec density <- spec smooth$spec</pre>
# Create a data frame for ggplot2
spec data <- tibble(Frequency = freq, Spectral Density = spec density)
ggplot(spec data, aes(x = Frequency, y = Spectral Density)) +
 geom line() +
 labs(title = "",
    x = "Frequency",
    y = "Spectral Density",
```

caption = "Figure 3. Periodogram of monthly precipitation time series smoothed using Daniell kernel.")

```
# Get confidence intervals at frequencies 1, 2, 3, 4, 5.
freqs_of_interest <- c(1, 2, 3, 4, 5) / 12 # Convert to frequency in cycles per month
indices <- sapply(freqs_of_interest, function(f) which.min(abs(freq - f)))
spec_values <- spec_density[indices]
dof <- 2 * 13 # 2 times the smoothing span, spans = c(13, 13)
ci_multiplier <- qchisq(c(0.025, 0.975), df = dof) / dof
ci_lower <- spec_values * ci_multiplier[1]
ci_upper <- spec_values * ci_multiplier[2]
result <- tibble(Frequency = freqs_of_interest * 12, # Convert back to original scale
        Spectral_Density = spec_values,
        CI_Lower = ci_lower,
        CI_Upper = ci_upper)
print(result)
```