# StatComp Project 2: Scottish weather

Stewart Ross (s2078228)

# Introduction

The aim of this report is to explore data collected by the Global Historical Climatology Network for eight weather stations in Scotland, spanning a time period from 1960 to 2018. The data includes information about each weather station such as latitude, longitude and elevation, as well as the precipitation and maximum and minimum temperature values for select years during this period. Unfortunately, some of the collected data is incomplete such as missing precipitation values, which is explored later on in the report.
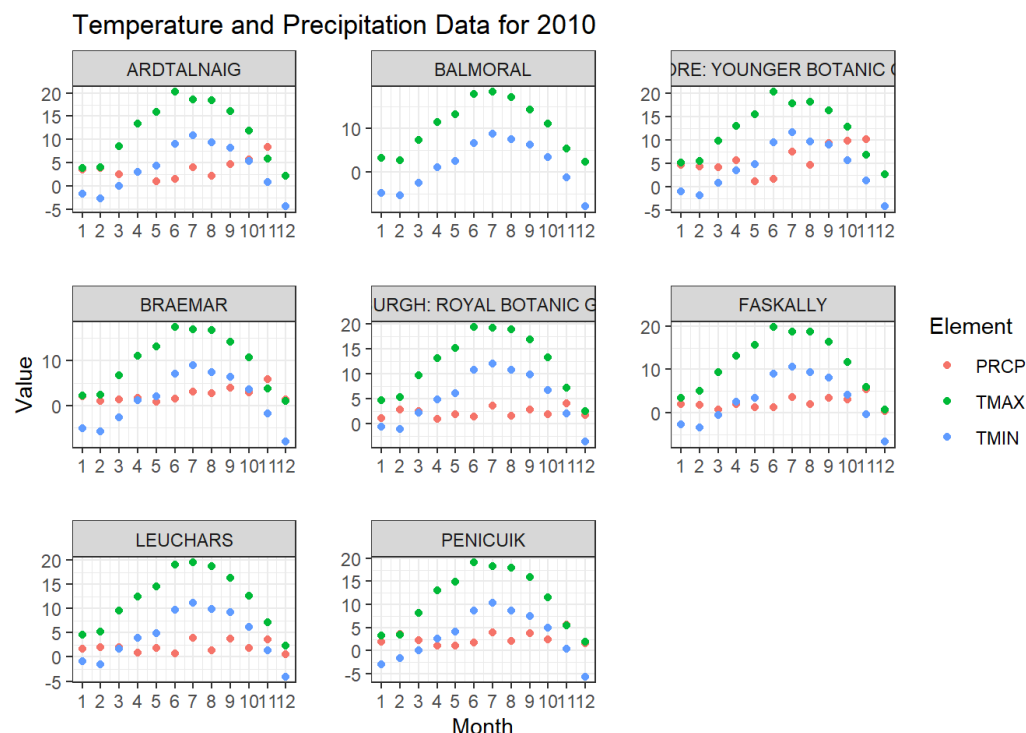
# Seasonal variability

## General analysis

In order to view the station data and the recorded measurements in one data frame, the following code was used to create the new variable `ghcnds`, such as in Lab 07. In addition to this, the column `Season` was added to show if a particular measurement was taken in summer or winter, with summer being the months between April and September and winter the remainder of the year. We also add a seperate column `summer` to display a TRUE statement if a measurement was taken in summer and a FALSE statement if it was taken in winter.

```
ghcnds <- left_join(ghcnd_values, ghcnd_stations, by = "ID") %>%
  mutate(Season = ifelse(Month %in% c(1, 2, 3, 10, 11, 12),"Winter","Summer"))
ghcnds <- ghcnds %>%
  mutate(Summer = Month %in% c(4, 5, 6, 7, 8, 9))
```

The following figure shows the temperature and precipitation data in 2010, with each subplot showing the results for a different station. The data is averaged by month for readability.



Temperature and Precipitation Data for 2010

From the plots above, it can be seen that in 2010 each weather station saw higher maximum and minimum temperatures in the summer months, peaking around July. However, we see that there is generally very little difference in precipitation values between summer and winter, with the notable exception of Benmore: Younger Botanic Gardens and Ardtalnaig, where there is significantly more rain in winter. The average maximum and minimum temperature and precipitation values across all stations and time periods are given in the following table:

| Element | Season | Average |
| --- | --- | --- |
| PRCP | Summer | 2.482327 |
| PRCP | Winter | 3.495553 |
| TMAX | Summer | 15.768091 |
| TMAX | Winter | 7.696776 |

| Element | Season | Average |
|---|---|---|
| TMIN | Summer | 7.234912 |
| TMIN | Winter | 1.418342 |

It can be seen that across all stations and time periods there is on average 40% more rain in winter than summer, and that the maximum and minimum temperatures in summer are around 100% and 400% greater respectively.

## Monte Carlo Permutation Test Approximated P-Values

A Monte Carlo Permutation test is a way of testing a null hypothesis by means of random permutations. We first assume the null hypothesis to be true, in this case that rainfall distribution is the same in winter as in summer, and are looking for evidence to reject this in favour of the alternative hypothesis, that the distributions have different expected values. We set a number of random permutations to be tested, in this case 1000, and permute the values of precipitation in the original data frame. For each permutation, the value of the permuted test statistic $T_{perm}$ = |permuted winter average − permuted summer average| is compared against the observed test statistic $T_{obs}$ = |winter average − summer average| from the original data frame. If $T_{perm}$ is greater than $T_{obs}$, a one is recorded, and after the 1000 permutations, the mean of these Boolean values is taken which is the estimated p-value $\hat{p}$. For a given significance level $\alpha$, if $\hat{p}$ is less than this value we reject the null hypothesis in favour of the alternative hypothesis and, if it is greater, we fail to reject the null hypothesis.

We can see from the permutation test that the number of successes, i.e number of $T_{perm} > T_{obs}$, follows a binomial distribution $X \sim \mathrm{Bin}(N, p)$, for some unknown $p$ and the number of permutations $N$. In the permutation test analysis, we can approximate $p$ by $\hat{p} = \frac{x}{N}$, and hence by the properties of the binomial distribution approximate the Monte Carlo standard deviation as $sd(\hat{p}) = \sqrt{\frac{\hat{p}(1-\hat{p})}{N}}$.

The function `permutation_test()` is used to approximate $\hat{p}$ and $sd(\hat{p})$ for a given station, and this can be used in a loop as shown below to approximate $\hat{p}$ and $sd(\hat{p})$ for each of the eight stations. The results were saved as .rsd files due to the high computational cost.

```
p_values <- numeric(nrow(ghcnd_stations))
SD_values <- numeric(nrow(ghcnd_stations))
for (i in 1:nrow(ghcnd_stations)){
  results <- permute_test(ghcnd_stations$ID[i], 1000)
  p_values[i] <- results$p_value
  SD_values[i] <- results$SD_value}
p_values <-saveRDS(p_values, file = "data/p_values.rds")
SD_values <-saveRDS(SD_values, file = "data/sd_values.rds")
```

The resulting data for each station is given in the following table:

| ID | Name | P_values | SD_values |
|---|---|---|---|
| UKE00105874 | BRAEMAR | 0.000 | 0.0000000 |
| UKE00105875 | BALMORAL | 0.000 | 0.0000000 |
| UKE00105884 | ARDTALNAIG | 0.000 | 0.0000000 |
| UKE00105885 | FASKALLY | 0.000 | 0.0000000 |
| UKE00105886 | LEUCHARS | 0.027 | 0.0051255 |
| UKE00105887 | PENICUIK | 0.000 | 0.0000000 |
| UKE00105888 | EDINBURGH: ROYAL BOTANIC GARDE | 0.677 | 0.0147875 |
| UKE00105930 | BENMORE: YOUNGER BOTANIC GARDE | 0.000 | 0.0000000 |

We can see that for six out of the eight stations $\hat{p}$ is zero, which means there is very strong evidence to reject the null hypothesis. If we assume $\alpha = 0.05$, although the evidence isn't as strong, we can see that as 0.03<0.05, the null hypothesis is also rejected for the Leuchars station. For the Edinburgh: Royal Botanic Gardens station, we see that as 0.68>0.05, we fail to reject the null hypothesis.

The true differences between the average winter and summer precipitation values for each weather station are given below:

```
ghcnds <- ghcnds %>%
  mutate(Summer = Month %in% c(4, 5, 6, 7, 8, 9))

winter_averages <- ghcnds %>%
  filter(Element == "PRCP", Season == "Winter") %>%
  group_by(ID, Name) %>%
  summarise(Winter_Average = mean(Value), .groups = "drop")

summer_averages <- ghcnds %>%
  filter(Element == "PRCP", Season == "Summer") %>%
  group_by(ID, Name) %>%
  summarise(Summer_Average = mean(Value), .groups = "drop")

station_averages <- left_join(winter_averages, summer_averages, by = c("ID", "Name")) %>%
  mutate(Average_Difference = abs(Winter_Average - Summer_Average))

knitr::kable(station_averages)
```

| ID | Name | Winter_Average | Summer_Average | Average_Difference |
|---|---|---|---|---|
| UKE00105874 | BRAEMAR | 2.878589 | 2.089790 | 0.7887997 |
| UKE00105875 | BALMORAL | 2.601891 | 1.994370 | 0.6075213 |
| UKE00105884 | ARDTALNAIG | 4.679448 | 2.668042 | 2.0114061 |
| UKE00105885 | FASKALLY | 2.899627 | 2.130559 | 0.7690686 |
| UKE00105886 | LEUCHARS | 1.931070 | 1.809706 | 0.1213640 |
| UKE00105887 | PENICUIK | 2.751264 | 2.326199 | 0.4250654 |
| UKE00105888 | EDINBURGH: ROYAL BOTANIC GARDE | 1.869338 | 1.844564 | 0.0247740 |
| UKE00105930 | BENMORE: YOUNGER BOTANIC GARDE | 8.353198 | 4.995383 | 3.3578143 |

It can be seen that the differences in average rainfall for the Leuchars and Edinburgh: Royal Botanic Gardens stations are significantly smaller than for the others, leading us to believe that the permutation test results are reasonable.

## Monte Carlo Permutation Test Approximated Confidence Intervals

A 95% confidence interval is the range in which the true population parameter, in this case $p$, lies with 95% probability. In order to approximate a confidence for each $p$, we can use the fact that the number of randomized test statistics being more extreme than the observed test statistic is normally distributed, to approximate the confidence interval as, $CI_p = \hat{p} \pm z_{0.975}\sqrt{\frac{\hat{p}(1-\hat{p})}{N}} = \frac{x}{N} \pm z_{0.975}\sqrt{\frac{x(N-x)}{N^3}}$, where $x$ is the observed number of such randomized test statistics and $N$ the number of permutations. In the case that p=0, we use the fact that a test is rejected only if $P_X(X = 0|p_0) = (1 - p_0)^N < 0.025$ to approximate $CI_p$ as $CIp = (0, 1 - 0.025^{1/N})$. The function `p_value_CI` takes 2 lists of the approximated p values and the sd values, as well as number of permutations, to return a data frame displaying the lower and upper bounds of the approximated confidence interval for each p:

```
knitr::kable(p_value_CI(p_values,SD_values,1000))
```

| ID | Name | P_values | SD_values | CI_lower | CI_upper |
|---|---|---|---|---|---|
| UKE00105874 | BRAEMAR | 0.000 | 0.0000000 | 0.0000000 | 0.0036821 |
| UKE00105875 | BALMORAL | 0.000 | 0.0000000 | 0.0000000 | 0.0036821 |
| UKE00105884 | ARDTALNAIG | 0.000 | 0.0000000 | 0.0000000 | 0.0036821 |
| UKE00105885 | FASKALLY | 0.000 | 0.0000000 | 0.0000000 | 0.0036821 |
| UKE00105886 | LEUCHARS | 0.027 | 0.0051255 | 0.0169542 | 0.0370458 |
| UKE00105887 | PENICUIK | 0.000 | 0.0000000 | 0.0000000 | 0.0036821 |
| UKE00105888 | EDINBURGH: ROYAL BOTANIC GARDE | 0.677 | 0.0147875 | 0.6480170 | 0.7059830 |
| UKE00105930 | BENMORE: YOUNGER BOTANIC GARDE | 0.000 | 0.0000000 | 0.0000000 | 0.0036821 |

We can see that for all stations where $\hat{p}$ is less that 0.05, the upper limit of the approximated confidence interval is less that 0.05, and that for the Edinburgh: Royal Botanic gardens station where $\hat{p}$ is greater than 0.05, the lower limit is above 0.05. This is hence strong evidence that the conclusions drawn from our analysis of the approximated p values is correct, and there is strong evidence to suggest that the only station where there isn't more rain in winter than summer is the Edinburgh: Royal Botanic gardens station.

# Spatial weather prediction

## Spatial Prediction Introduction

In this section we analyse a model to estimate average rainfall per month. As the average monthly precipitation values are heavily skewed, we take the square root of each value, and the DecYear mean for each month . The code for this is shown below:

```
ghcnds_sr <- ghcnds %>%
  filter(Element == "PRCP") %>%
  group_by(ID, Name, Year, Month, Element, Longitude, Latitude, Elevation) %>%
  summarise(Monthly_Average = mean(Value), .groups = "drop") %>%
  mutate(Value_sqrt_avg = sqrt(Monthly_Average))

DecYearAvg <- ghcnds %>%
  filter(Element == "PRCP") %>%
  group_by(Year, Month) %>%
  summarise(DecYearAvg = mean(DecYear), .groups = "drop")

ghcnds_sr <- ghcnds_sr %>%
  left_join(DecYearAvg, by = c("Year", "Month"))
```

We define five models, $M_0$ to $M_5$, with

$$M_0 = \text{Value\_sqrt\_avg} \sim \text{Intercept} + \text{Longitude} + \text{Latitude} + \text{Elevation} + \text{DecYearAvg}$$

and,

$$M_K = \text{Value\_sqrt\_avg} \sim \text{Intercept} + \text{Longitude} + \text{Latitude} + \text{Elevation} + \text{DecYearAvg} + \sum_{k=1}^{K}[\gamma_{c,k}\cos(2\pi kt) + \gamma_{s,k}\sin(2\pi kt)].$$

For each increase of $K$, we see the addition of two covariant trig terms, with frequency $k = 1, 2, \ldots$. As these terms are added to the model, the model should capture increasingly complex seasonal variability. The time variable $t$ is taken to be the average DecYear per month so that the lowest frequency $k = 1$ corresponds to a cosine function with a period of one year.

## Overview of Models and Results of Model Estimation

The function `Model_Generator()` takes a value $K$ which dictates the quantity and frequency of the covariate terms and outputs the given model formula. To examine the coefficients for each of the five models, a loop was created to store the fitted models and their associated coefficients in two lists. The coefficients for each of the models are printed below:

```
model_coefficients <- list()
models_fitted <- list()
for (k in 0:4){
  model_formula_k <- Model_Generator(k)
  model_k <- lm(model_formula_k, data = ghcnds_sr)
  model_coefficients[[k + 1]] <- coef(model_k)
  models_fitted[[k + 1]] <- model_k
}
for (i in 1:length(model_coefficients)) {
  cat("Model M", i - 1, "Coefficients:\n")
  print(model_coefficients[[i]])
}
```

```
## Model M 0 Coefficients:
##   (Intercept)     Longitude      Latitude     Elevation     DecYearAvg
##  2.4767951556 -0.5464154480 -0.1359611349  0.0004021791  0.0024110254
## Model M 1 Coefficients:
##                      (Intercept)                     Longitude
##                     2.4903204712                  -0.5466548373
##                         Latitude                     Elevation
##                    -0.1349762158                   0.0004042257
##                       DecYearAvg I(cos(2 * pi * 1 * DecYearAvg))
##                     0.0023759156                   0.1828164680
## I(sin(2 * pi * 1 * DecYearAvg))
##                    -0.1320694139
## Model M 2 Coefficients:
##                      (Intercept)                     Longitude
##                     2.4769205492                  -0.5466484515
##                         Latitude                     Elevation
##                    -0.1349331071                   0.0004042670
##                       DecYearAvg I(cos(2 * pi * 1 * DecYearAvg))
##                     0.0023814884                   0.1829293709
## I(sin(2 * pi * 1 * DecYearAvg)) I(cos(2 * pi * 2 * DecYearAvg))
##                    -0.1321093284                   0.0302996603
## I(sin(2 * pi * 2 * DecYearAvg))
##                     0.0006763433
## Model M 3 Coefficients:
##                      (Intercept)                     Longitude
##                     2.4732735959                  -0.5466619191
##                         Latitude                     Elevation
##                    -0.1349785313                   0.0004045066
##                       DecYearAvg I(cos(2 * pi * 1 * DecYearAvg))
##                     0.0023844957                   0.1827996241
## I(sin(2 * pi * 1 * DecYearAvg)) I(cos(2 * pi * 2 * DecYearAvg))
##                    -0.1321214746                   0.0301466169
## I(sin(2 * pi * 2 * DecYearAvg)) I(cos(2 * pi * 3 * DecYearAvg))
##                     0.0005288376                  -0.0067885505
## I(sin(2 * pi * 3 * DecYearAvg))
##                     0.0285280669
## Model M 4 Coefficients:
##                      (Intercept)                     Longitude
##                     2.4720157747                  -0.5466639296
##                         Latitude                     Elevation
##                    -0.1349685854                   0.0004043331
##                       DecYearAvg I(cos(2 * pi * 1 * DecYearAvg))
##                     0.0023848791                   0.1828549245
## I(sin(2 * pi * 1 * DecYearAvg)) I(cos(2 * pi * 2 * DecYearAvg))
##                    -0.1320910037                   0.0301210681
## I(sin(2 * pi * 2 * DecYearAvg)) I(cos(2 * pi * 3 * DecYearAvg))
##                     0.0006026946                  -0.0067670436
## I(sin(2 * pi * 3 * DecYearAvg)) I(cos(2 * pi * 4 * DecYearAvg))
##                     0.0284716117                   0.0140612946
## I(sin(2 * pi * 4 * DecYearAvg))
##                     0.0122540024
```

In order assess the models, we can examine the $R^2$ values. An $R^2$ value ranges from 0 to 1 and gives an indication of the fraction of variance that is in the dependent variable. Generally speaking, the higher the $R^2$ value, the better the model fits the given data. This does however not always serve as a good indicator if, for example, there exists high variability of the noise.

The $R^2$ value for each model is presented below:

```
models_r_squared <- numeric(length(models_fitted))
for (i in 1:length(models_fitted)) {
  model <- models_fitted[[i]]
  models_r_squared[i] <- summary(model)$r.squared
}
model_names <- c("M0", "M1", "M2", "M3", "M4")
r_squared_df <- data.frame(Model = model_names, R_squared = models_r_squared)
knitr::kable(r_squared_df )
```
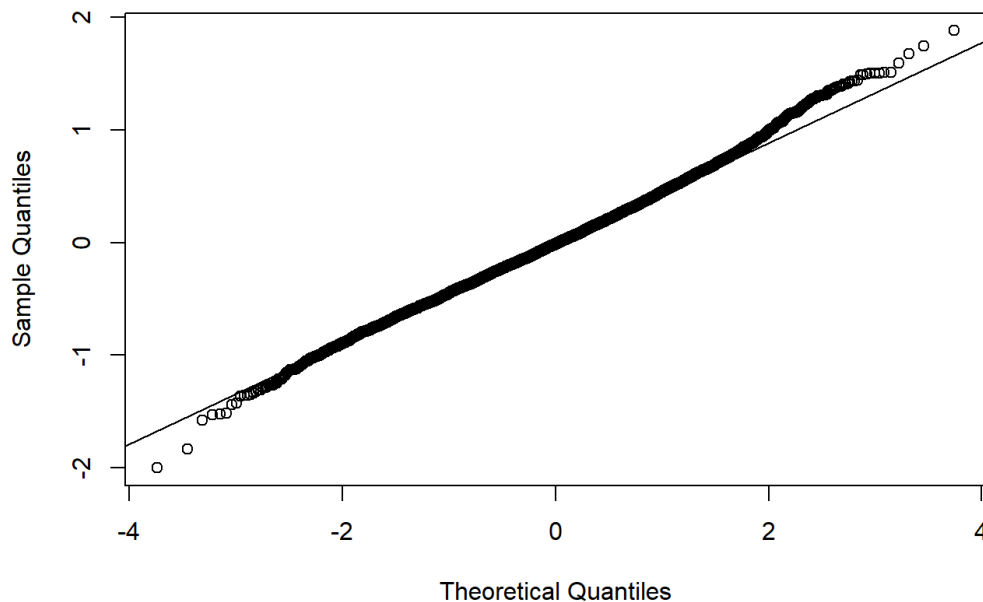
| Model | R_squared |
| --- | --- |
| M0 | 0.3423833 |

| Model | R_squared |
|-------|-----------|
| M1 | 0.4117222 |
| M2 | 0.4129767 |
| M3 | 0.4141532 |
| M4 | 0.4146312 |

We see that while there is a large difference between $M_0$ and $M_1$, the $R^2$ values for other models are broadly similar. As $M_4$ has the highest value, we can say there is a good chance it best fits the data. If we were to increase the value of $K$, we would expect the model to better fit the data, but we would have to be increasingly wary of overfitting, which would limit the model's ability to predict unseen values.

An important assumption of a liner regression model is that the residuals are normally distributed. In order to assess this, a Q-Q plot was conducted for the $M_4$ model, showing the residuals plotted against a diagonal line representing the theoretical normal distribution.

```
residuals <- resid(models_fitted[[5]])
qqnorm(residuals)
qqline(residuals)
```



We can see that although the tails are slightly heavy, the residuals follow a largely normal distribution, and this condition is indeed satisfied.

## Stratified Cross-Validation

In this section we test how well each model predicts precipitation for different locations. In order to do this, we will compute two score values for each precipitation measurement for each model, and aggregate to the 12 months of the year. The two computed scores are the standard error and Dawid-Sebastiani scores. We define the standard error scores as $\frac{\sum_{i=1}^{n}(y_i-\hat{y}_i)^2}{n}$ for sample size $n$, observed value $y_i$ and predicted value $\hat{y}_i$ and the Dawid-Sebastiani scores as $\sum_{i=1}^{n}\left(\frac{(y_i-\hat{y}_i)^2}{\sigma_i^2} + \log(\sigma_i^2)\right)$, for prediction standard deviation $\sigma_i$. While both scores are examples of proper scores as both satisfy the property $S(F,G) \geq S(G,G)$ for the true distribution $G$ and predicted distribution $F$, only the Dawid-Sebastiani score is strictly proper, as only for the Dawid-Sebastiani scores does a unique forecast $F$ maximize the expected scores. As the Dawid-Sebastiani also accounts for the prediction standard deviation, it is often the better method of scoring. In both cases, the lower the score, the better the forecasting model.
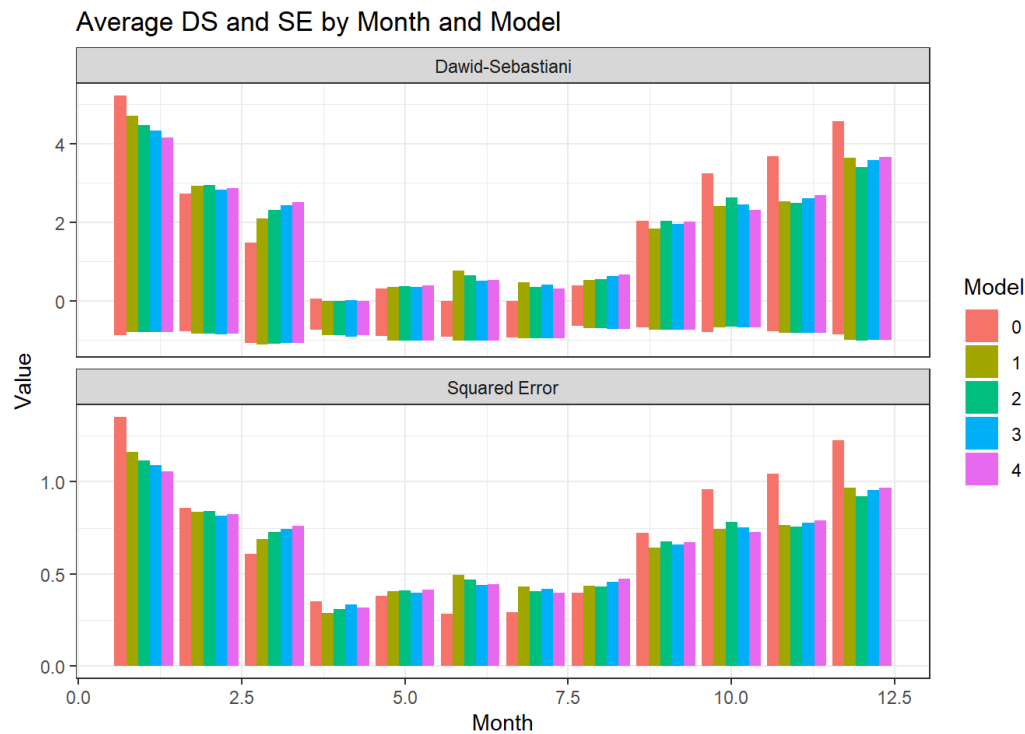
In order to compute these scores, we use the `cross_validated_scores` function. The function works by taking the given data and looping through all the station IDs, leaving one out each time to fit each model, and then using that model to predict the squared average precipitation values for measurements with the excluded ID. This is then used to compute the prediction standard deviations and expected values to compute the score values. A loop is created to compute these values for all models and this data frame is then aggregate to the 12 months of the year.

```
all_results <- data.frame()
for (k in 0:4) {
  model_formula <- Model_Generator(k)
  current_results <- cross_validated_scores(ghcnds_sr, model_formula)
  current_results$model <- k
  all_results <- rbind(all_results, current_results)
}
monthly_average <- aggregate(cbind(SE, DS) ~ Name + Month + model, data = all_results, FUN = mean)
```

In order to plot the results, we create a new data frame to index the score results as being either standard error or Dawid-Sebastiani scores.

```
monthly_average_indexed <- monthly_average %>%
  gather(key = key, value = value, DS, SE)
```
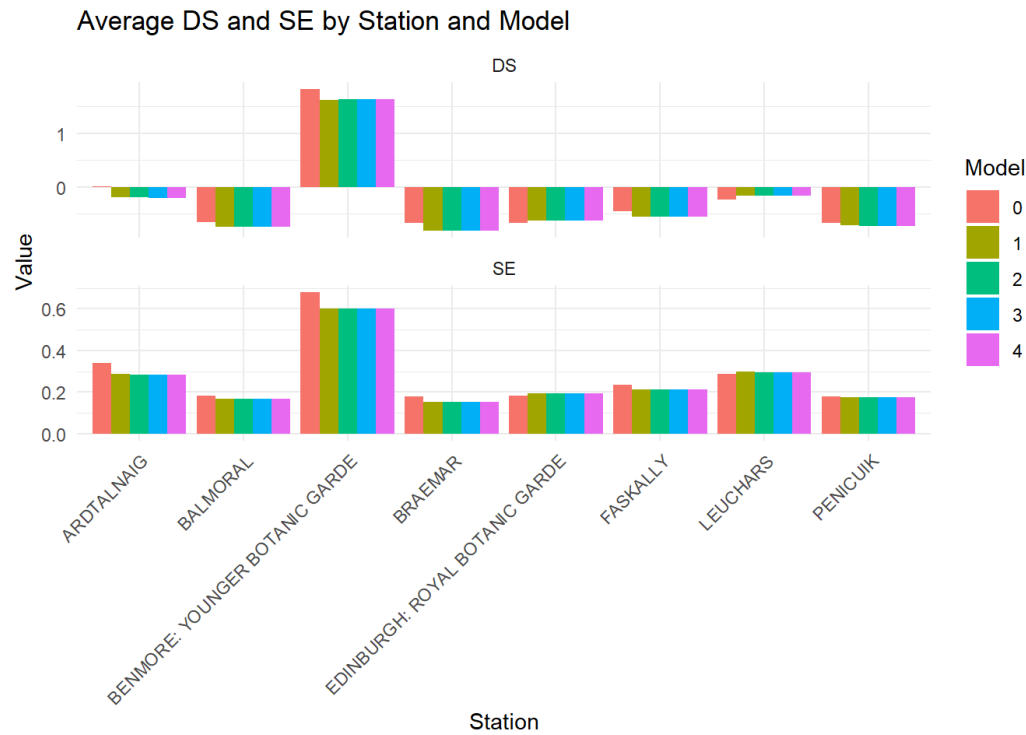
If we average over all stations and set the x axis to be in months, we get the following plot:



We see that in general the $M_0$ model appears to preform better in terms of the Dawid-Sebastiani and standard error scores in the summer months, while the other models tend to preform better in winter, with relatively little difference between $M_1$-$M_4$. This could be down to a number of reasons, including the fact that the weather is perhaps less variable in the summer months or that the models with more covariate terms suffer from overfitting.

To see if the models are more accurate at predicting weather data for certain stations, we average the data by weather station across all months, to produce the following bar chart. We again index the data frame to set the scores as being either standard error or Dawid-Sebastiani.

```
station_average <- aggregate(cbind(SE, DS) ~ Name + model, data = all_results, FUN = mean)

station_average_indexed <-station_average %>%
  gather(key = key, value = value, DS, SE)
```

## Average DS and SE by Station and Model



We see here that the scores for the $M_0$ model are generally worse than the others, with the exception of the Leuchars station, while the scores between the other models are broadly similar. We also see from the Dawid-Sebastiani and standard error scores that the models are worst at predicting the precipitation data for the Benmore: Younger Botamnic Gardens station. The models are almost equally accurate at predicting data for the the other stations, with the scores for the Ardtailnaig and Leuchars stations being slightly higher than the others. It can be seen that the stations which provide the worst scores tend to be relatively low in altitude, which is something that could be explored to improve the model further.

# Code appendix

# Function definitions

```r
library(tidyverse)
library(dbplyr)
library(grid)
library(StatCompLab)
# Stewart Ross (s2078228)

data(ghcnd_stations, package = "StatCompLab")
data(ghcnd_values, package = "StatCompLab")

#' Permutation Test
#'
#' @param station_id The ID of the station of which the approximated p value and
#' Monte Carlo Standard Deviation are to be computed
#' @param N The number of permutations
#'
#' @return A `data.frame` with the associated approximated p value and Monte Carlo Standard Deviation

permutation_test <- function(station_id, N) {
  station_data <- ghcnds %>%
    filter(ID == station_id, Element == "PRCP")

  # Approximate the winter average precipitation for the given station
  winter_average <- station_data %>%
    filter(Season == "Winter") %>%
    summarise(Winter_Average = mean(Value))$Winter_Average

  # Approximate the summer average precipitation for the given station
  summer_average <- station_data %>%
    filter(Season == "Summer") %>%
    summarise(Summer_Average = mean(Value))$Summer_Average

  # Compute the observed test statistic as the average of the winter and summer precipitation
  observed_difference <- abs(winter_average - summer_average)

  # Perform N permutations using the replicate function and calculate permuted differences
  permuted_differences <- replicate(N, {
    permuted_data <- station_data %>%
      mutate(Value = sample(Value, size = nrow(station_data), replace = FALSE)) %>%
      group_by(Season) %>%
      summarise(Permuted_Average = mean(Value), .groups = "drop")
    winter_average <- filter(permuted_data, Season == "Winter")$Permuted_Average
    summer_average <- filter(permuted_data, Season == "Summer")$Permuted_Average
    abs(winter_average - summer_average)})

  # Compute the approximated p value for the given station
  p_value <- mean(abs(permuted_differences) >= abs(observed_difference))
  SD_value <- sqrt(p_value * (1 - p_value) / N)
  return(data.frame(p_value, SD_value))
}

#' Compute the confidence regions
#'
#' @param p_values List of approximated p values for each station
#' @param SD_Values List of Monte Carlo Standard Deviations for each station
#'
#' @return A data frame with each station ID, Name, P value, Monte Carlo standard deviation and
#' upper and lower confidence limit boundaries

p_value_CI <- function(p_values, SD_values, N) {
  CI_lower <- numeric(length(p_values))
  CI_upper <- numeric(length(p_values))
  # Treat P=0 as special case
  for (i in 1:length(p_values)){
    if (p_values[i] == 0) {
      upper_limit <- 1 - (0.025^(1 / N))
      CI <- c(0, upper_limit)}
    else {
      # Compute lower and upper CI values
      lower_limit <- p_values[i] - qnorm(0.975) * SD_values[i]
```

```r
        upper_limit <- p_values[i] + qnorm(0.975) * SD_values[i]
        CI <- c(lower_limit, upper_limit)}
    CI_lower[i] <- CI[1]
    CI_upper[i] <- CI[2]}
  results_df <- data.frame(ID=ghcnd_stations$ID, Name = ghcnd_stations$Name,
                            P_values = p_values, SD_values = SD_values,
                            CI_lower = CI_lower, CI_upper = CI_upper)
  return(results_df)
}


# Question 2

#' Model Generator
#'
#' @param K The model formula to be computed i.e M_k
#'
#' @return The specified formula as a formula object which can be utilized by the lm function

Model_Generator <- function(K) {
  for (k in 0:K) {
    if (K == 0) {
      model_formula <- "Value_sqrt_avg ~ Longitude + Latitude + Elevation + DecYearAvg"}
    else {
      model_formula <- "Value_sqrt_avg ~ Longitude + Latitude + Elevation + DecYearAvg"
      for (k in 1:K) {
        model_formula <- paste0(model_formula, " + I(cos(2 * pi * ", k, " * DecYearAvg)) + I(sin(2 * pi * ", k, " * DecYea
rAvg))")
      }
    }
    model <- as.formula(model_formula)
  }
  return(model)
}


#' Cross-Validated Scores
#'
#' @param data The data frame to be analysed
#' @param formula The given M_K formula to be analysed
#'
#' @return A data frame with the same number of rows as the original data frame
#' and the columns being the given station name, year, month, standard error and Dawid-Sabastiani scores for each
#' data entry

cross_validated_scores <- function(data, formula) {
  result <- data.frame(
    Name=data$Name,
    Year=data$Year,
    Month = data$Month,
    SE = numeric(nrow(data)),
    DS = numeric(nrow(data)))
  for (id in unique(data$ID)) {
    fit <- lm(formula, data %>% filter(ID != id))
    pred <- predict(fit, newdata = data %>% filter(ID == id), se.fit = TRUE)
    residual_var <- sum(fit$residuals^2)/fit$df.residual
    sd_pred <- sqrt(pred$se.fit^2 + residual_var)
    result$SE[data$ID == id] <- (data$Value_sqrt_avg[data$ID == id] - pred$fit)^2
    result$DS[data$ID == id] <- ((data$Value_sqrt_avg[data$ID == id] - pred$fit)^2 / sd_pred^2) + log(sd_pred^2)
  }
  result
}


# Place your function definitions that may be needed in the report.Rmd, including function documentation.
# You can also include any needed library() calls here
```