

# **ANÁLISIS MULTIVARIANTE PARA COMPARAR EL ESTADO ECOLÓGICO DE LAS COMUNIDADES DE MACROINVERTEBRADOS EN RÍOS DEL PÁRAMO ECUATORIANO Y SU IMPACTO EN ÁREAS DE CONSERVACIÓN\***

**Maestría en Biodiversidad y Cambio Climático**

**Karina Hernández**

**16 de abril de 2024**

En el presente Notebook de RMarkdown, se desarrolla todo el proceso estadístico multivariante que busca Analizar las Componentes Principales (ACP) y los Clúster que determinan la taxonomía de variables que identifican de mejor forma a las comunidades de macroinvertebrados en ríos del páramo ecuatoriano. Primero, se aplicará un ACP para reducir la dimensionalidad de los datos y examinar patrones de variación en la composición de las variables físico-químicas y ecológicas entre diferentes ríos. Se realizará una evaluación de la adecuación del modelo mediante el estadístico de Kaiser-Meyer-Olkin (KMO) con su nivel significancia respecto al Test de Bartlett, el cual nos proporcionará información sobre la idoneidad de los datos para la aplicación del componentes principales. Además, se empleará un algoritmo de Clúster para agrupar variables en función de su similitud, esto permitirá identificar patrones para definir grupo de variables y xplorar posibles diferencias. Al integrar estos enfoques analíticos, este estudio busca proporcionar una comprensión más profunda del estado ecológico de los ríos del páramo ecuatoriano y de la efectividad de las áreas de conservación en la protección de las comunidades de macroinvertebrados. Los hallazgos de este análisis contribuirán a informar futuras estrategias de conservación y gestión de estos ecosistemas críticos.

# CONTENIDOS

## Listado de figuras

# Listado de Tablas

```
1 # Librerías y directorio de trabajo:
2 library(tidyverse)
3 library(DT)
4 library(highcharter)
5 library(psych)
6 library(gridExtra)
7 library(glue)
8 library(e1071)
9 library(reshape2)
10 library(factoextra)
11 library(FactoMineR)
12 library(corrplot)
13 library(knitr)
14 library(kableExtra)
15 library(RColorBrewer) # Librería para la paleta de colores
16
17 knitr::opts_knit$set(root.dir = "C:/Users/marcelochavez/Documents/TESIS/KARINA/")
```

# 1 Análisis Exploratorio de los Índices Físico-Químicos y Ecológicos

Empezaremos con mostrar un Boxplot comparativo para visualizar si existen valores atípicos y como se distribuyen las variables en los años **1995 y 2021** según el formato o estructura de cada tabla con escalas de medición originales en cada variable:

## 1.1 Dataset de los Índices Físico Químicos:

```
1 # Carga del DataSet y transformaciones -----
2
3 indices_fq <- readxl::read_excel("DATASETS/DatosFisicoQuimicosEcologicos.xlsx") %>%
4   mutate_if(is.numeric, ~ round(., 2)) %>%
5   rename_all(toupper)
6
7 df_largo <- indices_fq %>%
8   pivot_longer(cols = -RIOS,
9               names_to = c(".value", "ANIO"),
10              names_pattern = "(.*)_(\\d{4})")
11
12 indices_fq_ddc <- df_largo %>%
13   mutate(ID = glue("{RIOS}_{ANIO}")) %>%
14   column_to_rownames(var = "ID") %>%
15   select(-RIOS,
16         -ANIO)
17
18 var_scale <- as.data.frame(lapply(indices_fq_ddc, function(x) if(is.numeric(x)) scale
19
20 df_fq_1995 <- pivot_longer(df_largo,
```

```

21         cols = c(ANCHO,
22                 PROFUNDIDAD,
23                 VELOCIDAD,
24                 DESCARGA,
25                 PH,
26                 TEMPERATURA,
27                 OXIGENO_DISUELTO,
28                 CONDUCTIVIDAD,
29                 RIQUEZA,
30                 ABUNDANCIA,
31                 DIVERSIDAD,
32                 ABI,
33                 BMWP,
34                 IBF),
35         names_to = "VARIABLE",
36         values_to = "VALOR") %>%
37     filter(ANIO=="1995")
38
39 df_fq_2021 <- pivot_longer(df_largo,
40                             cols = c(ANCHO,
41                                     PROFUNDIDAD,
42                                     VELOCIDAD,
43                                     DESCARGA,
44                                     PH,
45                                     TEMPERATURA,
46                                     OXIGENO_DISUELTO,
47                                     CONDUCTIVIDAD,
48                                     RIQUEZA,
49                                     ABUNDANCIA,
50                                     DIVERSIDAD,
51                                     ABI,
52                                     BMWP,
53                                     IBF),
54                             names_to = "VARIABLE",
55                             values_to = "VALOR") %>%

```

```

56     filter(ANIO=="2021")
57
58
59 # Data para BoxPlot:
60
61 df_bx_95 <- data_to_boxplot(
62     df_fq_1995,
63     variable = VALOR,
64     group_var = VARIABLE,
65     # group_var2 = codigo,
66     add_outliers = T
67 )
68
69 df_bx_21 <- data_to_boxplot(
70     df_fq_2021,
71     variable = VALOR,
72     group_var = VARIABLE,
73     # group_var2 = codigo,
74     add_outliers = T
75 )
76
77 # Visualización del Dataset
78
79 DT::datatable(df_largo,
80               class = 'cell-border stripe',
81               # filter = 'top',
82               caption = htmltools::tags$caption(
83                 style = 'caption-side: bottom; text-align: left;',
84                 'Tabla 1: ', htmltools::em('Indicadores Físico-Químicos y Ecológicos pa
85                 extensions = c('Buttons','Scroller'),
86                 options = list(scrollX = TRUE,
87                               initComplete = JS(
88 "function(settings, json) {",
89 "$ (this.api().table().header()).css({'background-color': '#000', 'color': '#fff'}
90 "})"),

```



```
91     dom = 'Bfrtip',
92     buttons = c('excel'),
93     deferRender = TRUE,
94     scrollY = 500,
95     scroller = TRUE))
```

## 1.2 Boxplot de los indicadores físico-químicos y ecológicos:

```

1 highchart() %>%
2   hc_xAxis(type = "category",
3     labels = list(
4       rotation = 270
5     )) %>%
6   hc_add_series_list(df_bx_95) %>%
7   hc_xAxis(title = list(text = "VARIABLES")) %>%
8   hc_yAxis(title = list(text = "VALORES")) %>%
9   hc_title(text = "BOXPLOT COMPARATIVO ENTRE ÍNDICES FÍSICO QUÍMICOS y ECOLÓGICOS -
10  hc_subtitle(text = "AÑO 1995", align="left") %>%
11  hc_caption(text = "ELABORADO POR: Karina Hernández") %>%
12  hc_legend(enabled = F, title = list(text = "<b>ESTACIONES DE MONITOREO:</b>")) %>%
13  hc_tooltip(formatter = JS("function() {
14    var tooltip = '';
15    if (this.point.low && this.point.q1 && this.point.median && this.point.q3 && this
16      tooltip += '<b>Variable: </b>' + this.point.name + '<br/>' +
17        '<b>Mínimo: </b>' + this.point.low.toFixed(2) + '</b><br/>' +
18        '<b>Q1: </b>' + this.point.q1.toFixed(2) + '</b><br/>' +
19        '<b>Mediana: </b>' + this.point.median.toFixed(2) + '</b><br/>' +
20        '<b>Q3: </b>' + this.point.q3.toFixed(2) + '</b><br/>' +
21        '<b>Máximo: </b>' + this.point.high.toFixed(2) + '</b><br/>';
22    }
23    if (this.point.y && !(this.point.low && this.point.q1 && this.point.median && thi
24      tooltip += '<b>Valor Atípico: </b>' + this.point.y.toFixed(2) + '<br/>';
25    }
26    return tooltip;})) %>%
27  hc_add_theme(hc_theme_gridlight())

```

La variable **Abundancia** presenta una alta dispersión para el año 1995 en relación a todas las variables de estudio

```

1 highchart() %>%
2   hc_xAxis(type = "category",
3     labels = list(
4       rotation = 270
5     )) %>%

```

```

6 hc_add_series_list(df_bx_21) %>%
7 hc_xAxis(title = list(text = "VARIABLES")) %>%
8 hc_yAxis(title = list(text = "VALORES")) %>%
9 hc_title(text = "BOXPLOT COMPARATIVO ENTRE ÍNDICES FÍSICO QUÍMICOS y ECOLÓGICOS -")
10 hc_subtitle(text = "AÑO 2021", align="left") %>%
11 hc_caption(text = "ELABORADO POR: Karina Hernández") %>%
12 hc_legend(enabled = F, title = list(text = "<b>ESTACIONES DE MONITOREO:</b>")) %>%
13 hc_tooltip(formatter = JS("function() {
14 var tooltip = '';
15 if (this.point.low && this.point.q1 && this.point.median && this.point.q3 && this
16     tooltip += '<b>Variable: </b>' + this.point.name + '<br/>' +
17         '<b>Mínimo: </b>' + this.point.low.toFixed(2) + '</b><br/>' +
18         '<b>Q1: </b>' + this.point.q1.toFixed(2) + '</b><br/>' +
19         '<b>Mediana: </b>' + this.point.median.toFixed(2) + '</b><br/>' +
20         '<b>Q3: </b>' + this.point.q3.toFixed(2) + '</b><br/>' +
21         '<b>Máximo: </b>' + this.point.high.toFixed(2) + '</b><br/>';
22 }
23 if (this.point.y && !(this.point.low && this.point.q1 && this.point.median && thi
24     tooltip += '<b>Valor Atípico: </b>' + this.point.y.toFixed(2) + '<br/>';
25 }
26 return tooltip;}")) %>%
27 hc_add_theme(hc_theme_gridlight())

```

La variable **Diversidad** presenta una alta dispersión para el año 2021 en relación a todas las variables de estudio

### 1.3 Exploratorio de los datos:

```

1 exploratorio <- function(df) {
2   tipos <- sapply(df, class)
3   valores_min <- sapply(df, function(x) ifelse(is.numeric(x), min(x[complete.cases(
4   valores_max <- sapply(df, function(x) ifelse(is.numeric(x), max(x[complete.cases(
5   coeficientes_asimetria <- sapply(df, function(x) ifelse(is.numeric(x), round(e107
6   curtosis <- sapply(df, function(x) ifelse(is.numeric(x), round(kurtosis(x[complet
7   promedio <- sapply(df, function(x) ifelse(is.numeric(x), round(mean(x[complete.ca

```

```

8 medianas <- sapply(df, function(x) ifelse(is.numeric(x), round(median(x[complete.
9 modas <- sapply(df, function(x) ifelse(is.numeric(x), {
10     tab <- table(x[complete.cases(x)])
11     as.numeric(names(tab)[tab == max(tab)]), NA))
12 rangos <- sapply(df, function(x) ifelse(is.numeric(x), round(max(x[complete.cases
13 varianzas <- sapply(df, function(x) ifelse(is.numeric(x), round(var(x[complete.ca
14 desviaciones <- sapply(df, function(x) ifelse(is.numeric(x), round(sd(x[complete.
15 coeficientes_variacion <- sapply(df, function(x) ifelse(is.numeric(x), round(sd(x
16
17 resumen <- data.frame(Variable = names(df),
18                       Tipo = tipos,
19                       Mínimo = valores_min,
20                       Máximo = valores_max,
21                       Rango = rangos,
22                       Promedio = promedio,
23                       Mediana = medianas,
24                       Desviación_Estándar = desviaciones,
25                       Coeficiente_Variación = coeficientes_variacion,
26                       Moda = modas,
27                       Varianza = varianzas,
28                       Coeficiente_Asimetría = coeficientes_asimetria,
29                       Curtosis = curtosis,
30                       stringsAsFactors = FALSE)
31
32 resumen$Tipo <- ifelse(resumen$Tipo == "factor", "Categoría", resumen$Tipo)
33 resumen$Tipo <- ifelse(resumen$Tipo == "POSIXct", "Fecha", resumen$Tipo)
34 resumen$Tipo <- ifelse(resumen$Tipo == "logical", "Booleana", resumen$Tipo)
35 resumen$Tipo <- ifelse(resumen$Tipo == "numeric", "Numérica", resumen$Tipo)
36
37 # Elimina los nombres de las etiquetas de las filas
38 rownames(resumen) <- NULL
39
40 return(resumen)
41 }
42

```

```

43 DT::datatable(exploratorio(df_largo),
44               class = 'cell-border stripe',
45               # filter = 'top',
46               caption = htmltools::tags$caption(
47                 style = 'caption-side: bottom; text-align: left;',
48                 'Tabla 2: ', htmltools::em('Estimadores de Tendencia Central, Dispersión
49               extensions = c('Buttons','Scroller'),
50               options = list(scrollX = TRUE,
51                             initComplete = JS(
52                 "function(settings, json) {",
53                 "$(this.api().table().header()).css({'background-color': '#000', 'color': '#fff'}",
54                 "})",
55                 dom = 'Bfrtip',
56                 buttons = c('excel'),
57                 deferRender = TRUE,
58                 scrollY = 500,
59                 scroller = TRUE))

```

De las medidas del Análisis Exploratorio tanto para tendencia central, como de dispersión se puede deducir lo siguiente:

- Las variables que presentan mayor variabilidad son Abundancia y Descarga con **Coefficientes de Variación** por encima del 100%, esto lo puede corroborar la desviación estándar en comparación con el promedio de cada variable respectivamente.
- La variable Diversidad tiene un **Rango** ( $ValorMximo - ValorMnimo$ ) considerable, esto nos puede alertar más adelante sobre el comportamiento de dicha variable.
- El **Coefficiente de Asimetría** en la variable Temperatura es cercano a **0** lo que podría sugerir que la variable se acerca a una **Distribución Normal Estándar**
- Respecto a la **Curtosis** (Mayor a 3) en las variables Velocidad y Diversidad son distribuciones con colas más largas, lo que puede sugerir la presencia de datos atípicos.
- En las variables Descarga y Ancho en **Curtosis** presentan valores cercanos a **0**, lo que puede sugerir variables con una Distribución cercana a la Normal.

- Es importante mencionar que el **Análisis Exploratorio** no es determinante para analizar la presencia de normalidad en los datos. Por consiguiente, es importante a través de un **Test Estadístico** realizar la comprobación correspondiente.
- Finalmente, en un contexto de análisis de datos previo a la aplicación de cualquier técnica estadística, una vez revisado los **estimadores en conjunto con los Boxplot** podemos confirmar que es necesario **realizar un escalamiento a la matriz de datos**, para evitar los sesgos que pueden causar los valores atípicos y su unidad de medida, en el conjunto completo de los indicadores físico-químicos y ecológicos.

## 2 Análisis de la normalidad de los datos:

Con base a los resultados obtenidos del Análisis Exploratorio, a continuación analizaremos el tipo de distribución de cada variable, es evidente que si consideramos a un año de estudio (1995 o 2021) con 11 estaciones de monitoreo, no podemos hablar de datos cercanos a una **Distribución Normal o Gaussiana** por el principio de la **Ley de los Grandes Números**. Por lo cual se ha considerado los datos de los dos años.

En cada variable se aplica el *Test de Normalidad* y su significancia estadística para determinar si obedece o no a muestra proveniente de una Distribución Normal Estándar  $N \approx (\bar{x} = 0, \sigma^2 = \text{constante})$ .

El **Test de Shapiro-Wilk (W)** (Gutiérrez, 2008), es uno de las pruebas estadísticas más robustas, y se obtiene de la siguiente forma:

$$W = \frac{\left(\sum_{i=1}^n a_i x_{(i)}\right)^2}{\sum_{i=1}^n (x_{(i)} - \bar{x})^2}$$

**Donde:**

- $n$  es el tamaño de la muestra.
- $x_i$  es el  $i$  — simo valor más pequeño de la muestra ordenada.
- $\bar{x}$  es la media de la muestra.
- $a_i$  son los coeficientes que se utilizan para calcular el estadístico de Shapiro-Wilks, y se calculan de la siguiente forma:

$$\begin{aligned} a_i &= \frac{m_i}{\sum_{i=1}^n m_i^2} \\ m_i &= \sum_{j=1}^n (x_{(j)} - \bar{x}) s_{(j)} \\ s_{(j)} &= \sum_{k=1}^j (x_{(k)} - \bar{x})^2 \end{aligned}$$

Los coeficientes  $a_i$  y los estadísticos auxiliares  $m_i$  y  $x_j$  se calculan a partir de las observaciones de la muestra y se utilizan en la fórmula del estadístico de Shapiro-Wilk. Estos coeficientes y estadísticos se definen matemáticamente de la siguiente manera:

**Donde:**

- $m_i$  es la media de los elementos de las siguientes dos sumas:
- La suma en el numerador se calcula sobre todos los  $n$  valores de la muestra.
- La suma en el denominador se calcula sobre todos los  $n$  valores de la muestra.
- $s_j$  representa la suma de cuadrados de las diferencias entre cada observación y la media, considerando solo las primeras  $j$  observaciones en orden ascendente. Esta cantidad se utiliza para calcular  $m_i$ , que a su vez se utiliza para calcular  $a_i$ , y finalmente en el cálculo del estadístico  $W$ .

**Nota técnica:** El valor del estadístico de Shapiro-Wilk ( $W$ ) se compara con los valores críticos de la distribución para determinar si la muestra sigue una distribución normal. Si el valor de  $W$  es significativamente menor que el valor crítico correspondiente, entonces hay evidencia para rechazar la hipótesis nula de normalidad.

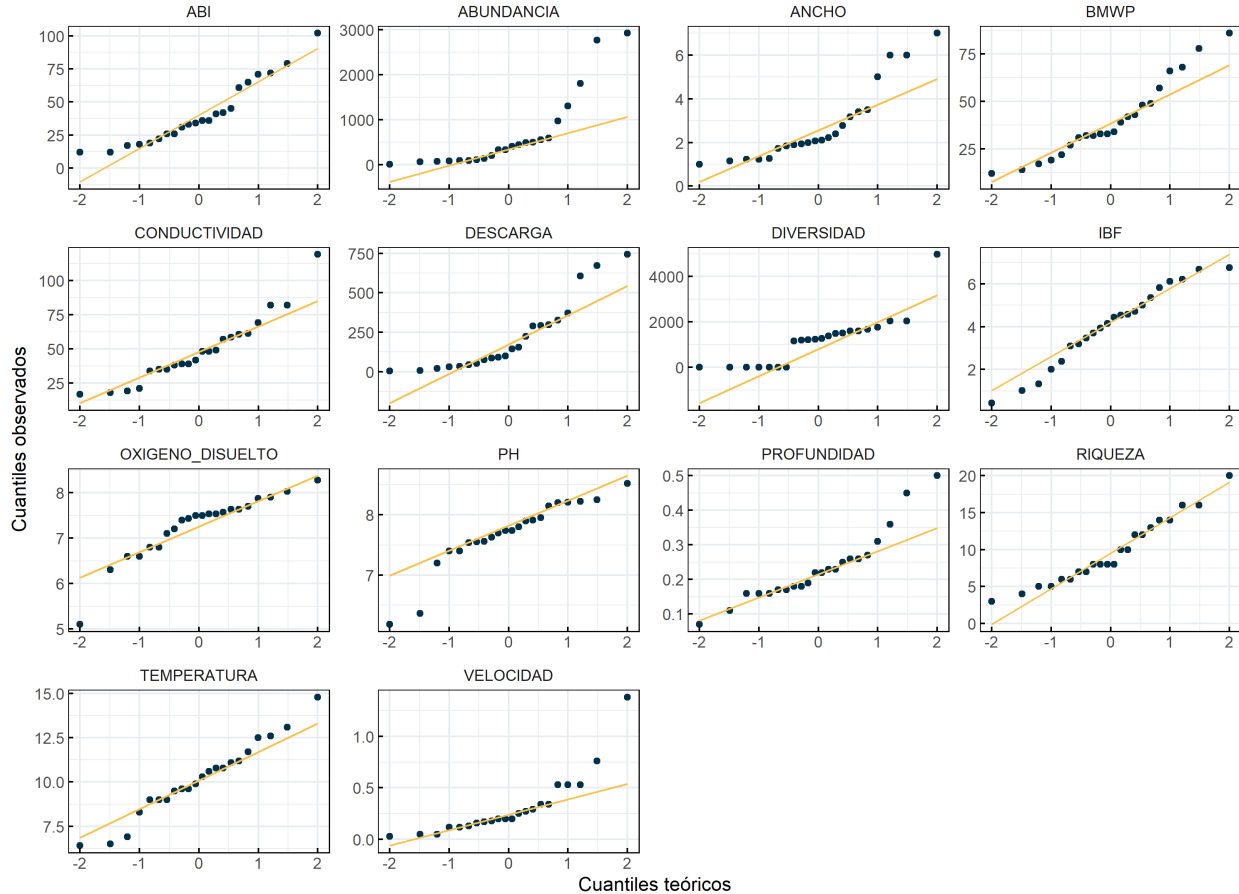
```

1  numeric_vars <- indices_fq_ddc[sapply(indices_fq_ddc, is.numeric)]
2
3  # Convertir el dataframe en formato largo (long format) para usar facet_wrap
4  df_long <- tidyr::pivot_longer(numeric_vars, cols = everything())
5
6  # Crear el Q-Q plot con facet_wrap
7  ggplot(df_long, aes(sample = value)) +
8    geom_qq(color = "#003049") +
9    stat_qq_line(color = "#fcbf49") +
10   labs(
11     title = "Gráfico de los Cuantiles (Q-Q Plot) para los Índices Físico-Químicos",
12     subtitle = "Estaciones de Monitoreo (Ríos) para los años 1995 y 2021",
13     x = "Cuantiles teóricos",
14     y = "Cuantiles observados",
15     caption = "Elaborado por: Karina Hernández"
16   ) +
17   facet_wrap(~ name, scales = "free") +
18   theme_minimal() +
19   theme(
20     panel.grid = element_line(color = "#e9ecef", size = 0.5), # Añadir grid
21     panel.border = element_rect(color = "black", fill = NA), # Añadir recuadro
22     axis.ticks = element_line(color = "black"), # Añadir ticks en los ejes
23     axis.line = element_blank(), # Eliminar líneas de ejes interiores

```



Gráfico de los Cuantiles (Q-Q Plot) para los Índices Físico-Químicos y Ecológicos  
Estaciones de Monitoreo (Ríos) para los años 1995 y 2021



```

24 plot.title = element_text(hjust = 0.5), # Centrar el título
25 plot.caption = element_text(hjust = 0), # Alinear el footnote a la izquierda
26 axis.ticks.length = unit(-0.1, "cm") # Ajustar longitud de los ticks hacia a
27 )

```

## 2.1 Test No paramétrico de Shapiro-Wilk para evaluar la normalidad en los Índices Físico-Químicos y Ecológicos:

```

1 # Definir una función para aplicar el test de Shapiro-Wilk a todas las variables numéricas
2 shapiro_test_all <- function(df) {
3
4     # Filtrar solo las variables numéricas

```

```

5  numeric_vars <- df[sapply(df, is.numeric)]
6
7  # Crear una lista para almacenar los resultados del test
8  test_results <- list()
9
10 # Iterar sobre cada variable numérica y aplicar el test de Shapiro-Wilk
11 for (var_name in names(numeric_vars)) {
12
13     result <- shapiro.test(numeric_vars[[var_name]])
14
15     # Redondear el estadístico W y el valor p a dos decimales
16     result$statistic <- round(result$statistic, 2)
17
18     result$p.value <- round(result$p.value, 2)
19
20     test_results[[var_name]] <- result
21 }
22
23 # Retornar los resultados como una lista
24 return(test_results)
25 }
26
27 # Aplicar la función al dataframe indices_fq_ddc
28 shapiro_results <- shapiro_test_all(indices_fq_ddc)
29
30 # Convertir los resultados en un tibble
31 results_df <- tibble::enframe(shapiro_results,
32                               name = "Variable",
33                               value = "Resultados")
34
35 # Extraer los estadísticos W y los valores p en columnas separadas
36 results_df <- results_df %>%
37   mutate(W = sapply(Resultados,
38                     "[[",
39                     "statistic"),

```

```

40
41     p_value = sapply(Resultados,
42                       "[[",
43                       "p.value"))
44
45 # Seleccionar las columnas relevantes
46 results_df <- results_df %>%
47   select(Variable, W, p_value)
48
49 # Definir los niveles de significancia
50 significance_levels <- c(0.05, 0.01, 0.10)
51
52 # Función para determinar si el valor p es mayor o menor que un nivel de significancia
53 check_significance <- function(p_value,
54                                significance_level) {
55
56   if (p_value < significance_level) {
57     return("Menor al  $\alpha$ ")
58   } else {
59     return("Mayor o igual  $\alpha$ ")
60   }
61 }
62
63 # Aplicar la función a cada valor p y nivel de significancia
64 for (level in significance_levels) {
65
66   results_df <- results_df %>%
67     mutate(
68       !!paste0("Nivel de Significancia al ",
69               level) := sapply(p_value,
70                                check_significance,
71                                significance_level = level),
72       !!paste0("Rechazas H0 al ",
73               level) := ifelse(p_value < level,
74                                "Sí",

```

```

75         "No")
76     )
77 }
78
79 # Mostrar la tabla actualizada
80 knitr::kable(results_df, format = "markdown")

```

| Variable         | W        | p_value | Nivel de Significancia al 0.05 | Rechazas H0 al 0.05 | Nivel de Significancia al 0.01 | Rechazas H0 al 0.01 | Nivel de Significancia al 0.1 | Rechazas H0 al 0.1 |
|------------------|----------|---------|--------------------------------|---------------------|--------------------------------|---------------------|-------------------------------|--------------------|
| ANCHO            | 0.830.00 | 0.000   | Menor al $\alpha$              | Sí                  | Menor al $\alpha$              | Sí                  | Menor al $\alpha$             | Sí                 |
| PROFUNDIDAD      | 0.940.03 | 0.000   | Menor al $\alpha$              | Sí                  | Mayor o igual $\alpha$         | No                  | Menor al $\alpha$             | Sí                 |
| VELOCIDAD        | 0.750.00 | 0.000   | Menor al $\alpha$              | Sí                  | Menor al $\alpha$              | Sí                  | Menor al $\alpha$             | Sí                 |
| DESCARGA         | 0.830.00 | 0.000   | Menor al $\alpha$              | Sí                  | Menor al $\alpha$              | Sí                  | Menor al $\alpha$             | Sí                 |
| PH               | 0.880.01 | 0.000   | Menor al $\alpha$              | Sí                  | Mayor o igual $\alpha$         | No                  | Menor al $\alpha$             | Sí                 |
| TEMPERATURA      | 0.980.84 | 0.000   | Mayor o igual $\alpha$         | No                  | Mayor o igual $\alpha$         | No                  | Mayor o igual $\alpha$        | No                 |
| OXIGENO DISUELTO | 0.850.01 | 0.000   | Menor al $\alpha$              | Sí                  | Mayor o igual $\alpha$         | No                  | Menor al $\alpha$             | Sí                 |
| CONDUCTIVIDAD    | 0.920.07 | 0.000   | Mayor o igual $\alpha$         | No                  | Mayor o igual $\alpha$         | No                  | Menor al $\alpha$             | Sí                 |
| RIQUEZA          | 0.950.33 | 0.000   | Mayor o igual $\alpha$         | No                  | Mayor o igual $\alpha$         | No                  | Mayor o igual $\alpha$        | No                 |
| ABUNDANCIA       | 0.710.00 | 0.000   | Menor al $\alpha$              | Sí                  | Menor al $\alpha$              | Sí                  | Menor al $\alpha$             | Sí                 |
| DIVERSIDAD       | 0.790.00 | 0.000   | Menor al $\alpha$              | Sí                  | Menor al $\alpha$              | Sí                  | Menor al $\alpha$             | Sí                 |
| ABI              | 0.910.04 | 0.000   | Menor al $\alpha$              | Sí                  | Mayor o igual $\alpha$         | No                  | Menor al $\alpha$             | Sí                 |
| BMWP             | 0.940.16 | 0.000   | Mayor o igual $\alpha$         | No                  | Mayor o igual $\alpha$         | No                  | Mayor o igual $\alpha$        | No                 |
| IBF              | 0.960.60 | 0.000   | Mayor o igual $\alpha$         | No                  | Mayor o igual $\alpha$         | No                  | Mayor o igual $\alpha$        | No                 |

## 2.2 Comprobación de hipótesis para análisis normalidad:

Para la comprobación de normalidad se ha planteado la siguiente prueba de hipótesis:

$H_0$  : Las variables provienen de una Distribución Normal

$H_1$  : Las variables NO provienen de una Distribución Normal

**Análisis:** De igual forma en un contexto general los indicadores físico-químicos y ecológicos, con un nivel de significancia del 5%. De las 14 variables, 5 presentan normalidad y el resto de variables no lo son. Este resultado nos conduce a generalizar que para todo el conjunto de datos es necesario reducir la dimensionalidad y/o aplicar un **método No Paramétrico** para descubrir clústers o taxonomías entre las unidades de estudio.

## 3 Análisis Factorial por Componentes Principales:

El **análisis factorial** en un contexto de rigurosidad estadística presente los siguientes pasos:

- Matriz de correlaciones y su significancia estadística
- Determinante de la matriz de correlaciones
- Prueba de esfericidad de Bartlett
- Comprobación de hipótesis para el análisis factorial por componentes principales
- Gráfico de Sedimentación
- Proporción de Varianza Acumulada
- Componentes Principales
- Autovalores y Varianzas
- Círculo de las Correlaciones
- Variables versus Dimensiones
- Biplot (variables versus las estaciones de monitoreo)
- Índice Multivariante

### 3.1 Matriz de Correlaciones:

```
1 # Calcular correlaciones y p-values
2 correlation_matrix <- cor(indices_fq_ddc)
3 p_values <- cor.mtest(indices_fq_ddc)$p
4
5 # Convertir p-values en asteriscos de acuerdo con el nivel de significancia
6 significant_level <- ifelse(p_values < 0.001, "***",
7                             ifelse(p_values < 0.01, "**",
8                                     ifelse(p_values < 0.05, "*", "")))
```

```

9
10 # Crear gráfico
11 ggplot(data = melt(round(correlation_matrix, 2)), aes(x = Var1, y = Var2, fill = value)) +
12   geom_tile() +
13   geom_text(aes(label = paste(round(value, 2), significant_level)), color = "black")
14   labs(
15     title = "Mapa de calor de correlaciones",
16     x = "",
17     y = "",
18     fill = "Nivel de Correlación",
19     caption = "Elaborado por: Karina Hernández"
20   ) +
21   scale_fill_gradient(low = "#caf0f8", high = "#2a9d8f") +
22   theme_minimal() +
23   theme(
24     axis.text.x = element_text(angle = 90, vjust = 0.5, hjust = 1),
25     axis.title.x = element_text(angle = 90),
26     axis.title.y = element_text(angle = 0),
27     panel.grid = element_line(color = "#e9ecef", linewidth = 0.5), # Cambio a 1
28     panel.border = element_rect(color = "black", fill = NA),
29     axis.ticks = element_line(color = "black"),
30     axis.line = element_blank(),
31     plot.title = element_text(hjust = 0.5),
32     plot.caption = element_text(hjust = 0),
33     axis.ticks.length = unit(-0.1, "cm")
34   )

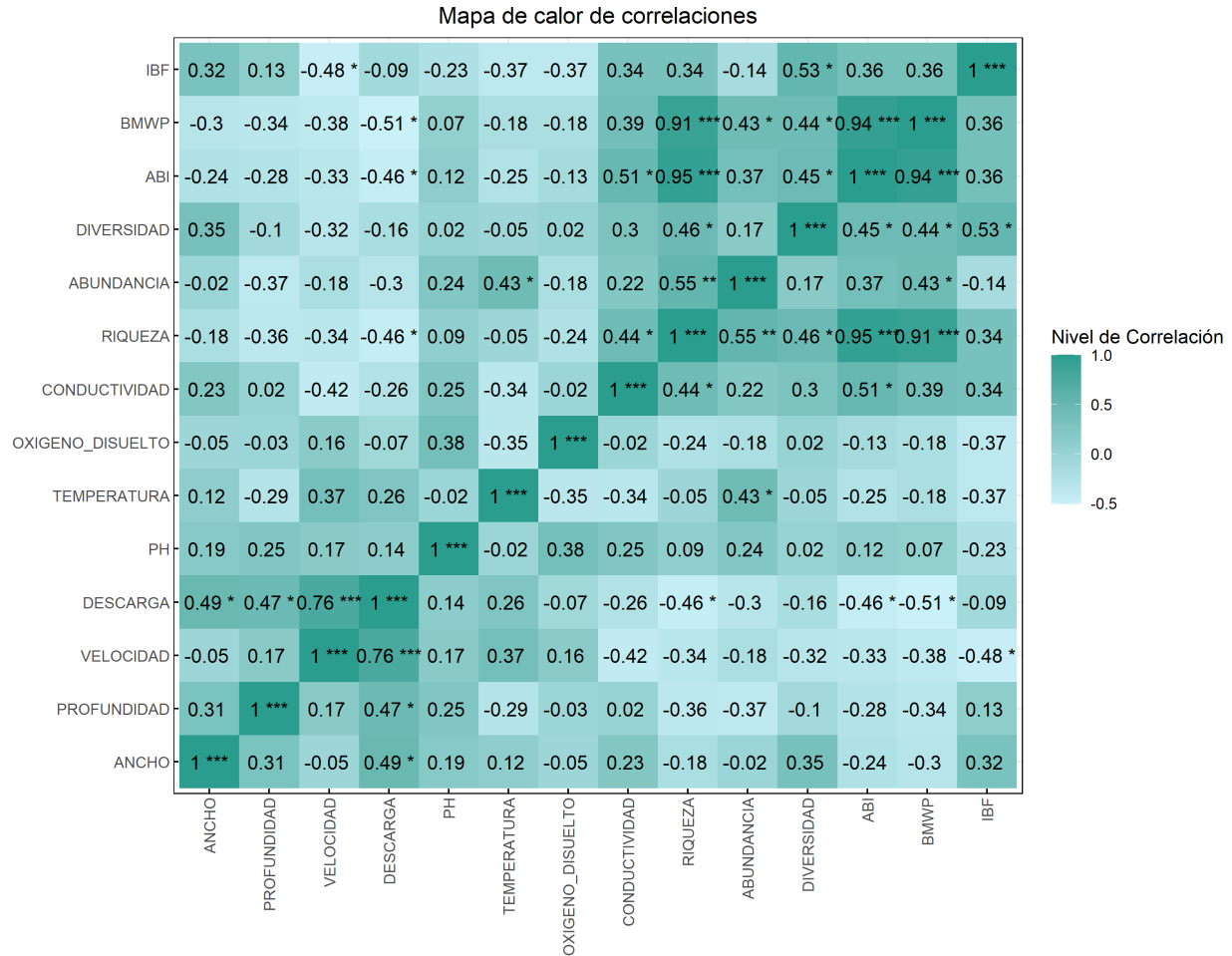
```

Es importante conocer que el nivel de significancia de la matriz de correlaciones para definir variables que a más de tener un buen nivel de correlación, es significativo estadísticamente.

$p - value < 0.001(***)$  : Es significativa la correlación al 99%

$p - value < 0.05(**)$  : Es significativa la correlación al 95%

$p - value < 0.01(*)$  : Es significativa la correlación al 90%



## 3.2 Determinante de la matriz de correlaciones

```
1 det(cor(indices_fq_ddc))
```

```
[1] 2.706342e-07
```

Si están correlacionadas, el determinantes de la matriz de correlaciones será próximo a cero. De lo contrario el **Análisis Factorial** no es viable.

## 3.3 Prueba de esfericidad de Bartlett:

```
1 KMO(indices_fq_ddc)
```

Kaiser-Meyer-Olkin factor adequacy

Call: KMO(r = indices\_fq\_ddc)



Overall MSA = 0.44

MSA for each item =

|         |             |                  |               |
|---------|-------------|------------------|---------------|
| ANCHO   | PROFUNDIDAD | VELOCIDAD        | DESCARGA      |
| 0.26    | 0.65        | 0.35             | 0.41          |
| PH      | TEMPERATURA | OXIGENO_DISUELTO | CONDUCTIVIDAD |
| 0.36    | 0.26        | 0.17             | 0.61          |
| RIQUEZA | ABUNDANCIA  | DIVERSIDAD       | ABI           |
| 0.60    | 0.42        | 0.61             | 0.48          |
| BMWP    | IBF         |                  |               |
| 0.58    | 0.55        |                  |               |

```
1 cortest.bartlett(cor(indices_fq_ddc))
```

```
$chisq
```

```
[1] 1413.954
```

```
$p.value
```

```
[1] 3.472263e-236
```

```
$df
```

```
[1] 91
```

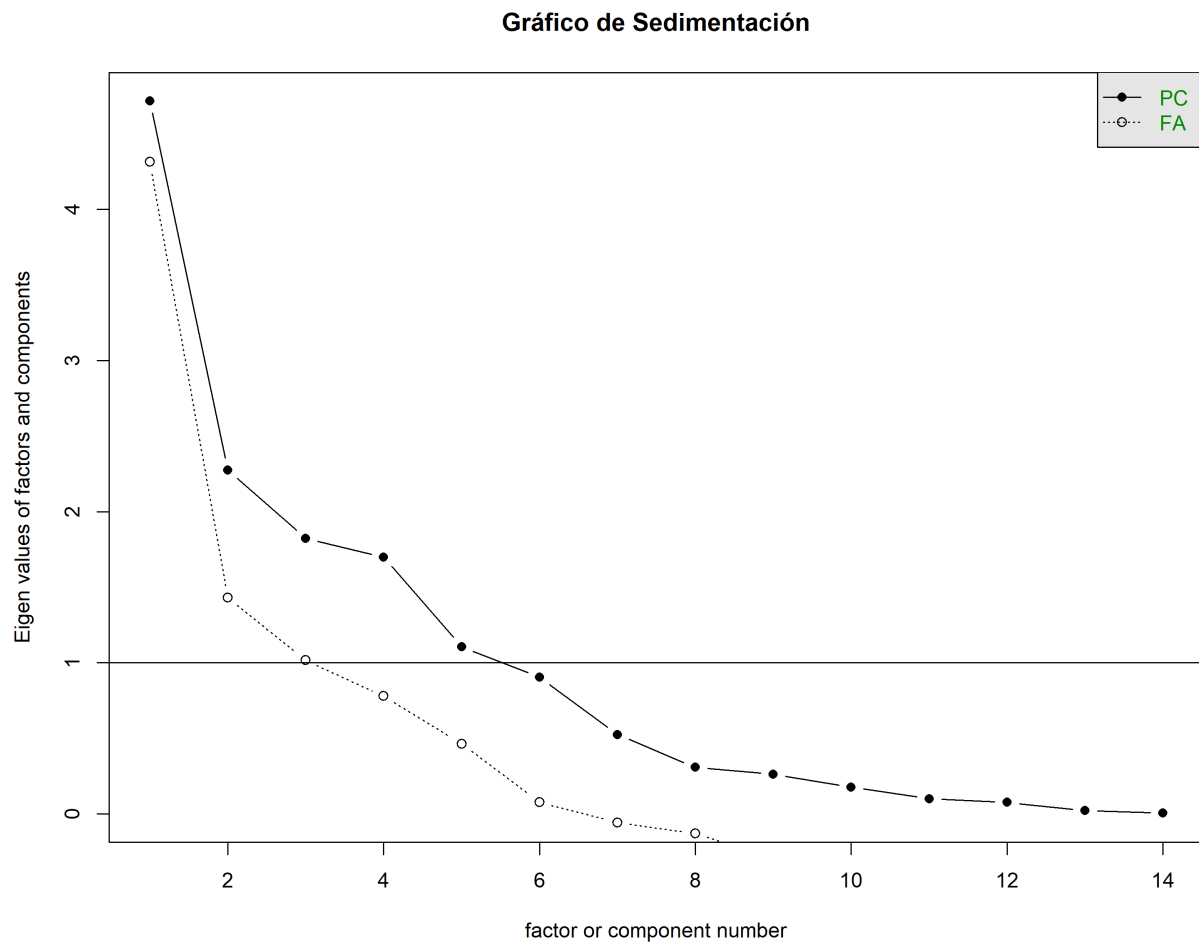
Un KMO de 0.44 es considerado bajo y podría sugerir que un análisis factorial puede no ser adecuado para estos datos. Sin embargo, la prueba de Bartlett tiene un p-valor menor al nivel de significancia del 5%, lo que indica que las variables están correlacionadas y que, por lo tanto, podría justificarse aplicar un método multivariante para reducción de dimensionalidad.

La hipótesis nula es que la matriz de correlaciones es una matriz identidad, lo que significa que no existirían correlaciones significativas, por lo que el modelo factorial no sería pertinente.

### 3.4 Comprobación de hipótesis para el análisis factorial por componentes principales

$H_0$  : Existe esfericidad

$H_1$  : No existe esfericidad



La esfericidad implica que todas las variables tienen la misma varianza y que no hay correlaciones entre ellas.

De aquí que el valor del  $KMO$  de 0.44 y su p-valor es menor al 95% de confianza. Por lo tanto: Rechazamos  $H_0$ , y concluimos que es factible estadísticamente realizar la reducción de dimensiones a través de un **Análisis Factorial por Componentes Principales**

### 3.5 Gráfico de Sedimentación:

```
1 scree(indices_fq_ddc,
2       main ="Gráfico de Sedimentación")
```

### 3.6 Proporción de varianza acumulada:

```
1 # Cálculo de los autovalores y autovectores:
2 value_pro<-eigen(cor(indices_fq_ddc))
3
4 Acumulado <- cumsum(value_pro$values)
5
6 Prop.acumulado <- Acumulado/sum(value_pro$values)
7
8 val.prop.porce <- data.frame(value_pro$values,
9                               Acumulado,
10                              Prop.acumulado)
11
12 row.names(val.prop.porce) = c(expression(lambda[1]),
13                                expression(lambda[2]),
14                                expression(lambda[3]),
15                                expression(lambda[4]),
16                                expression(lambda[5]),
17                                expression(lambda[6]),
18                                expression(lambda[7]),
19                                expression(lambda[8]),
20                                expression(lambda[9]),
21                                expression(lambda[10]),
22                                expression(lambda[11]),
23                                expression(lambda[12]),
24                                expression(lambda[13]),
25                                expression(lambda[14]))
26
27 colnames(val.prop.porce) <- c("Valor Propio",
28                               "Acumulado",
29                               "Prop. Acumulado")
30 kable(val.prop.porce,
31       caption = "Valores propios desde la matriz de correlaciones",
32       digits = 2,
33       format.args = list(decimal.mark = ",")) %>%
34   row_spec(row = 5,
```

35  
36

```
background="#a7c957",
bold=T)
```

Tabla 3.1: Valores propios desde la matriz de correlaciones

|                  | Valor Propio | Acumulado    | Prop. Acumulado |
|------------------|--------------|--------------|-----------------|
| lambda[1]        | 4,72         | 4,72         | 0,34            |
| lambda[2]        | 2,27         | 6,99         | 0,50            |
| lambda[3]        | 1,82         | 8,82         | 0,63            |
| lambda[4]        | 1,70         | 10,51        | 0,75            |
| <b>lambda[5]</b> | <b>1,11</b>  | <b>11,62</b> | <b>0,83</b>     |
| lambda[6]        | 0,90         | 12,52        | 0,89            |
| lambda[7]        | 0,52         | 13,05        | 0,93            |
| lambda[8]        | 0,31         | 13,36        | 0,95            |
| lambda[9]        | 0,26         | 13,62        | 0,97            |
| lambda[10]       | 0,18         | 13,80        | 0,99            |
| lambda[11]       | 0,10         | 13,90        | 0,99            |
| lambda[12]       | 0,08         | 13,97        | 1,00            |
| lambda[13]       | 0,02         | 13,99        | 1,00            |
| lambda[14]       | 0,01         | 14,00        | 1,00            |

El **Valor Propio** de la matriz de correlaciones cuando desciende hasta 1 sugiere que el número de componentes o factores que pueden explicar la variabilidad de todo el conjunto de datos está en un  $\lambda = 5$  y lo corrobora el **Gráfico de Sedimentación**

### 3.7 Componentes Principales:

```
1 pca <- prcomp(indices_fq_ddc,
2               scale=TRUE,
3               center = TRUE,
4               method = "varimax")
5 summary(pca)
```

Importance of components:

PC1      PC2      PC3      PC4      PC5      PC6      PC7

|                        |       |        |        |        |         |         |         |
|------------------------|-------|--------|--------|--------|---------|---------|---------|
| Standard deviation     | 2.172 | 1.5080 | 1.3503 | 1.3036 | 1.05141 | 0.95050 | 0.72437 |
| Proportion of Variance | 0.337 | 0.1624 | 0.1302 | 0.1214 | 0.07896 | 0.06453 | 0.03748 |
| Cumulative Proportion  | 0.337 | 0.4994 | 0.6297 | 0.7511 | 0.83002 | 0.89456 | 0.93204 |

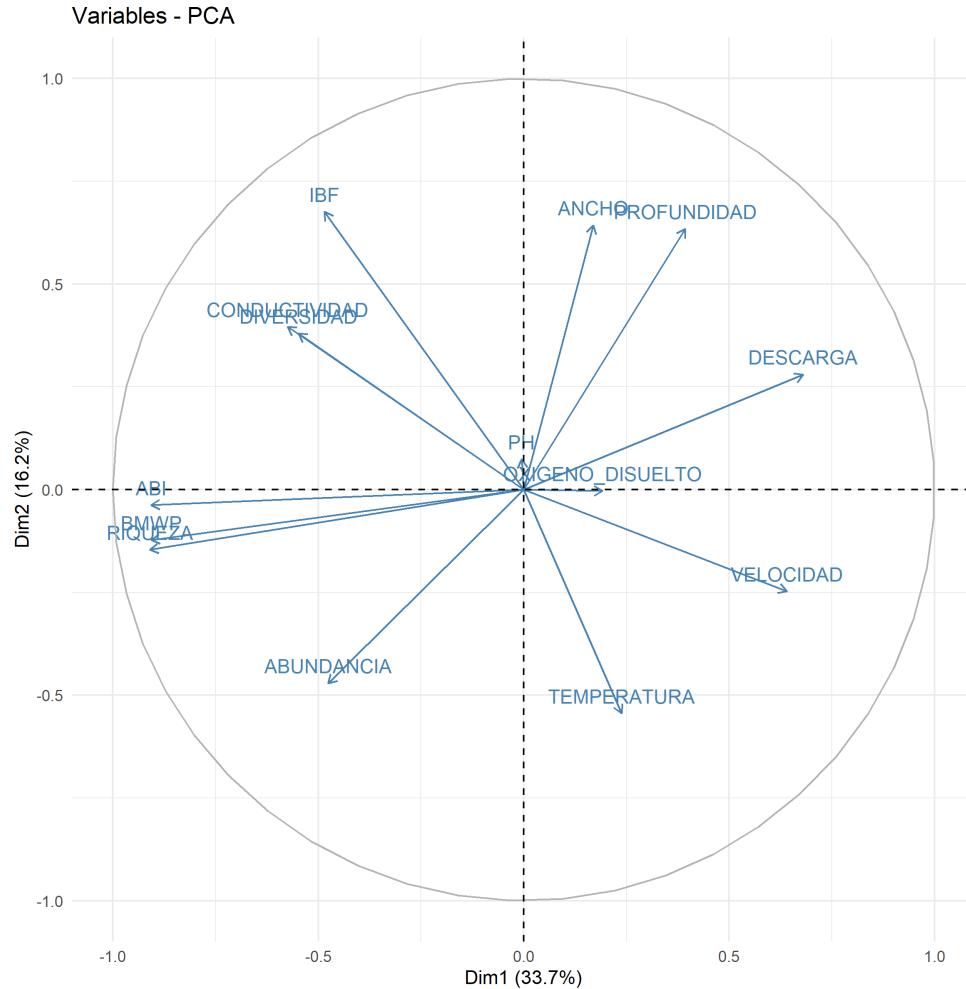
|                        |         |         |        |        |         |         |         |
|------------------------|---------|---------|--------|--------|---------|---------|---------|
|                        | PC8     | PC9     | PC10   | PC11   | PC12    | PC13    | PC14    |
| Standard deviation     | 0.55491 | 0.51189 | 0.4216 | 0.3153 | 0.27573 | 0.14732 | 0.08124 |
| Proportion of Variance | 0.02199 | 0.01872 | 0.0127 | 0.0071 | 0.00543 | 0.00155 | 0.00047 |
| Cumulative Proportion  | 0.95403 | 0.97275 | 0.9855 | 0.9926 | 0.99798 | 0.99953 | 1.00000 |

### 3.8 Autovalores y varianzas:

```
1 get_eig(pca)
```

|        | eigenvalue  | variance.percent | cumulative.variance.percent |
|--------|-------------|------------------|-----------------------------|
| Dim.1  | 4.717836203 | 33.69883002      | 33.69883                    |
| Dim.2  | 2.274213490 | 16.24438207      | 49.94321                    |
| Dim.3  | 1.823386036 | 13.02418597      | 62.96740                    |
| Dim.4  | 1.699435601 | 12.13882572      | 75.10622                    |
| Dim.5  | 1.105467492 | 7.89619637       | 83.00242                    |
| Dim.6  | 0.903448292 | 6.45320209       | 89.45562                    |
| Dim.7  | 0.524706456 | 3.74790326       | 93.20353                    |
| Dim.8  | 0.307925627 | 2.19946876       | 95.40299                    |
| Dim.9  | 0.262032510 | 1.87166079       | 97.27466                    |
| Dim.10 | 0.177786579 | 1.26990413       | 98.54456                    |
| Dim.11 | 0.099431722 | 0.71022659       | 99.25479                    |
| Dim.12 | 0.076026874 | 0.54304910       | 99.79783                    |
| Dim.13 | 0.021703083 | 0.15502202       | 99.95286                    |
| Dim.14 | 0.006600033 | 0.04714309       | 100.00000                   |

Al analizar los componentes principales podemos observar que en el **Componente 5** se explica el **83%** de la varianza de todas las dimensiones



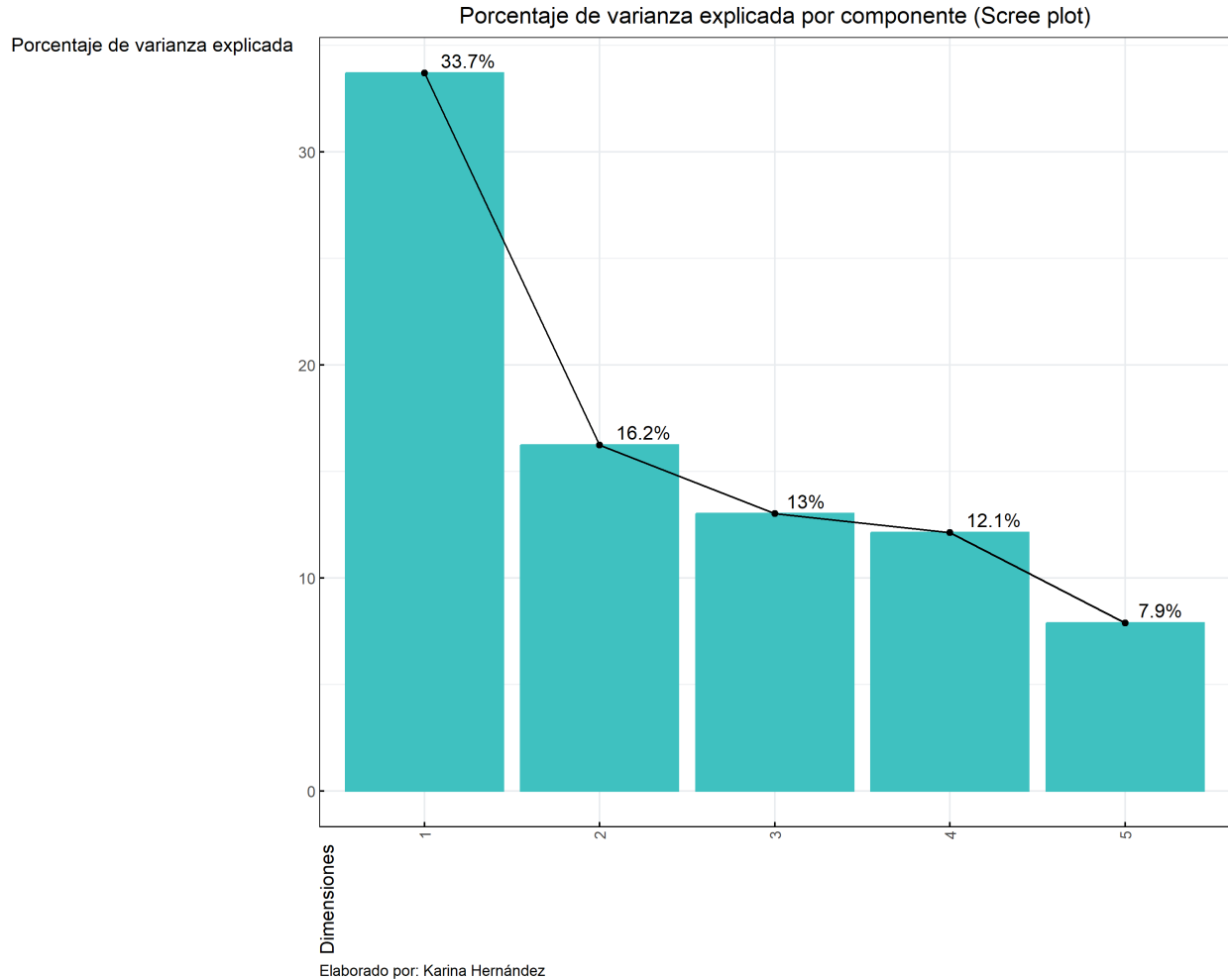
### 3.9 Círculo de correlaciones:

```
1 fviz_pca_var(pca,
2               col.var="steelblue") +
3 theme_minimal()
```

En el círculo de las correlaciones es una representación gráfica de las correlaciones entre las variables originales y los componentes principales extraídos durante el análisis factorial. Cada variable original se representa como un vector en el espacio de los componentes principales, donde la longitud del vector indica la correlación entre esa variable y el componente principal correspondiente, y la dirección del vector indica la dirección en la que esa variable contribuye al componente principal. Por ejemplo, al examinar el círculo de correlaciones, podemos observar que las variables **ABI**, **BMWP** y **RIQUEZA** están inversamente correlacionadas con la

Dimensión 2. Por otro lado, las variables **ANCHO y PROFUNDIDAD** muestran una correlación más fuerte con la Dimensión 1. En resumen, el círculo de correlaciones proporciona una representación visual que ayuda a entender cómo las variables originales están relacionadas y contribuyen a la estructura de los componentes principales identificados durante el análisis factorial

```
1 fviz_screepplot(pca,
2                 ncp = 5,
3                 barfill = "#3fc1c0",
4                 barcolor = "#3fc1c0",
5                 addlabels = TRUE,
6                 ylab = "Porcentaje de varianza explicada",
7                 xlab = "Dimensiones",
8                 main = "Porcentaje de varianza explicada por componente (Scree plot)",
9                 hjust = -0.3) +
10 theme(plot.title = element_text(size = 14, face = "bold"),
11        plot.caption = element_text(size = 8, face = "italic")) +
12 labs(caption = "Elaborado por: Karina Hernández") +
13 theme_minimal() +
14 theme(
15     axis.text.x = element_text(angle = 90, vjust = 0.5, hjust = 1),
16     axis.title.x = element_text(angle = 90),
17     axis.title.y = element_text(angle = 0),
18     panel.grid = element_line(color = "#e9ecef", linewidth = 0.5), # Cambio a li
19     panel.border = element_rect(color = "black", fill = NA),
20     axis.ticks = element_line(color = "black"),
21     axis.line = element_blank(),
22     plot.title = element_text(hjust = 0.5),
23     plot.caption = element_text(hjust = 0),
24     axis.ticks.length = unit(-0.1, "cm")
25 )
```

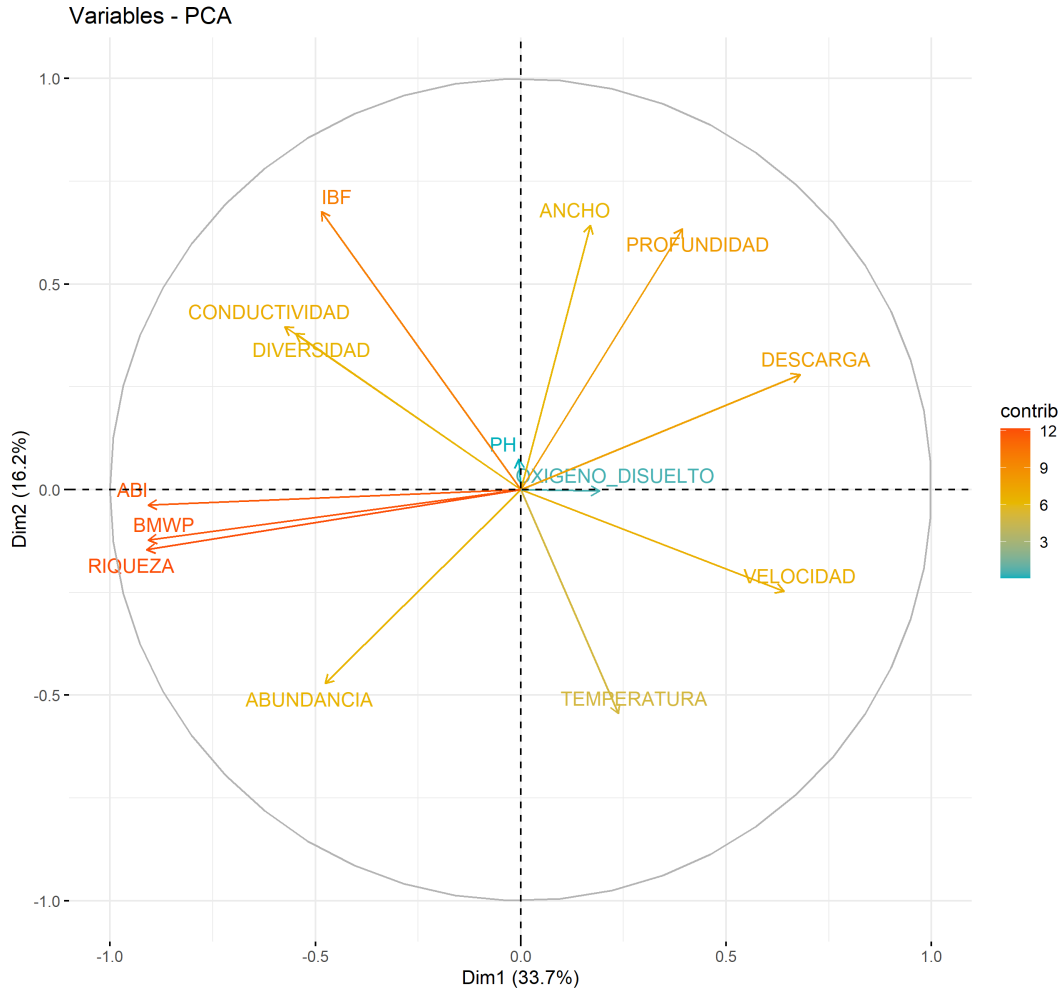


### 3.10 Variables versus Dimensiones:

```
1 fviz_pca_var(pca,  
2             col.var="contrib",  
3             gradient.cols = c("#00AFBB", "#E7B800", "#FC4E07"),  
4             repel = TRUE)
```

Es fundamental examinar el nivel de contribución de las variables a cada uno de los componentes principales. En este sentido, observamos que variables como **ABI**, **BMWP** y **RIQUEZA** destacan por su significativa contribución a la **Dimensión 2**. Sin embargo, al considerar **OXÍGENO** y **PH**, a pesar de ubicarse en los mismos componentes, la longitud de sus vectores indica una contribución mínima al modelo factorial en construcción.





### 3.11 Biplot (variables versus las estaciones de monitoreo)

```

1 # Biplot of individuals and variables
2 fviz_pca_biplot(pca, repel = TRUE)

```

Variables ecológicas como **ABI**, **BMWP** y **RIQUEZA** exhiben una alta similitud en la estación de monitoreo 8 durante el año 1995. De manera similar: **CONDUCTIVIDAD**, **DIVERSIDAD** e **IBF** comparten características destacadas en las estaciones de monitoreo 7 y 11, también en el año 1995. Por otro lado, un conjunto distinto de variables como: **ANCHO**, **PROFUNDIDAD** y **DESCARGA** para el año 2021, evidencian patrones similares en las estaciones de monitoreo 1, 7, 8 y 9. Sin embargo, variables como: **ABUNDANCIA**, **TEMPERATURA** y **VELOCIDAD**, a pesar de su comportamiento de no similitud entre ellas, aportan de forma significativa al modelo factorial individualmente

