

Construcción de Retroalimentación: Módulo 1

Construcción de un modelo estadístico base (Portafolio Implementación)

Roel De la Rosa - A01197595
13/9/2022

Durante este reporte lo que se busca es analizar un conjunto de datos que contienen información sobre los niveles de contaminación que se encuentran en peces en distintos lagos. Se busca encontrar y justificar relaciones entre las distintas variables para poder encontrar, en este caso, cuáles son las variables que más afectan el nivel de contaminación de mercurio en la carne de los peces. Para hacerlo se utilizaron métodos estadísticos tales como las pruebas de normalidad de Mardia y de Anderson-Darling, para saber el comportamiento de las variables, además de análisis de correlación y de componentes principales para entender mejor las relaciones entre las variables. Se concluyó que la Alcalinidad, el PH, el Calcio y la Clorofila son los factores más relevantes.

Algunas de las preguntas que se quisieran responder son las siguientes: ¿Cómo se distribuyen las variables que se analizarán? ¿Hay relaciones entre las distintas variables? ¿Es necesario escalar los datos para obtener mejores análisis?

Considero que este tipo de estudios son necesarios, pues por ejemplo, el mercurio es uno de los metales más tóxicos con los que una persona puede ingerir. Es necesario que se hagan estudios para entender cómo es que ciertas variables pueden afectar el nivel de contaminación de mercurio en los animales que más se consumen para evitar que alguien consuma material contaminado y termine con efectos adversos a su salud.

```
## Warning: package 'ggplot2' was built under R version 4.1.3

## Warning: replacing previous import 'lifecycle::last_warnings' by
## 'rlang::last_warnings' when loading 'pillar'

## Warning: package 'reshape2' was built under R version 4.1.3

## Warning: package 'psych' was built under R version 4.1.3

##

## Attaching package: 'psych'

## The following objects are masked from 'package:ggplot2':
##
##   %%, alpha

## Warning: package 'Hmisc' was built under R version 4.1.3

##

## Attaching package: 'Hmisc.'

## The following object is masked from 'package:psych':
##
##   describe

## The following objects are masked from 'package:base':
##
##   format.pval, units

## Warning: package 'clusterSim' was built under R version 4.1.3

## Warning: package 'factextra' was built under R version 4.1.3

## Welcome! Want to learn more? See two factextra-related books at https://goo.gl/ve3Wba

## Warning: package 'mvnornalTest' was built under R version 4.1.3

##

## Attaching package: 'mvnornalTest'

## The following object is masked from 'package:psych':
##
##   mardia

## Warning: package 'MVN' was built under R version 4.1.3

## Warning: package 'RVAideMemoire' was built under R version 4.1.3

## *** Package RVAideMemoire v 0.9-81-2 ***
```

Leemos los datos

```
## X1      X2      X3      X4      X5      X6      X7      X8      X9      X10     X11     X12
## 1 1 Alligator  5.9 6.1 3.0 0.1 2.3 5 0.8 1.23 5 0.8 1.23 5 0.8 1.23 5 1
## 2 2 Annie      3.5 5.1 1.9 3.2 1.33 7 0.92 1.90 1.33 0
## 3 3 Apopka     116 0 9.1 44.1 128 0.3 4 8.1 5 8 0.4 0.4 0
## 4 4 Blue Cypress 39.4 6.9 16.4 3.5 0.44 12 0.13 0.8 0.4 0.4 0
## 5 5 Brick      2.5 4.6 2.9 1.0 1.20 12 0.69 1.50 1.33 1
## 6 6 Bryant     19.6 7.3 4.5 44.1 0.27 14 0.04 0.48 0.25 1
```

```
## 'data.frame': 53 obs. of 12 variables:
## $ X1: int 1 2 3 4 5 6 7 8 9 10 ...
## $ X2: chr "Alligator" "Annie" "Apopka" "Blue Cypress" ...
## $ X3: num 5.9 3.5 116 39.4 2.5 19.6 5.2 7.1 3.5 26.4 0.8 ...
## $ X4: num 6.1 5.1 0 1 6.9 4.6 7.3 5.4 8.1 5 8 0.4 0.4 ...
## $ X5: num 3.0 3.9 44.1 16.4 2.9 4.5 2.8 55.2 0.2 9.2 4.0 ...
## $ X6: num 0.7 3.2 128 3.3 0.5 1.8 ...
## $ X7: num 1.23 1.33 0.04 0.44 1.2 0.27 0.48 0.19 0.83 0.81 ...
## $ X8: int 5 7 6 12 12 14 10 12 24 12 ...
## $ X9: num 0.85 0.92 0.04 0.13 0.69 0.04 0.3 0.09 0.26 0.41 ...
## $ X10: num 1.43 1.9 0.06 0.84 1.5 0.48 0.72 0.38 1.4 1.47 ...
## $ X11: num 1.53 1.33 0.04 0.44 1.33 0.25 0.45 0.16 0.72 0.81 ...
## $ X12: int 1 0 0 0 1 1 1 1 1 1 ...
```

Transformación de los datos

Renombro las variables para poderlas interpretar de mejor manera.

X1 = número de identificación X2 = nombre del lago X3 = alcalinidad (mg/l de carbonato de calcio) X4 = PH X5 = calcio (mg/l) X6 = clorofila (mg/l) X7 = concentración media de mercurio (parte por millón) en el tejido muscular del grupo de peces estudiados en cada lago X8 = número de peces estudiados en el lago X9 = mínimo de la concentración de mercurio en cada grupo de peces X10 = máximo de la concentración de mercurio en cada grupo de peces X11 = estimación (mediante regresión) de la concentración de mercurio en el pez de 3 años (o promedio de mercurio cuando la edad no está disponible) X12 = indicador de la edad de los peces (0: jóvenes, 1: maduros)

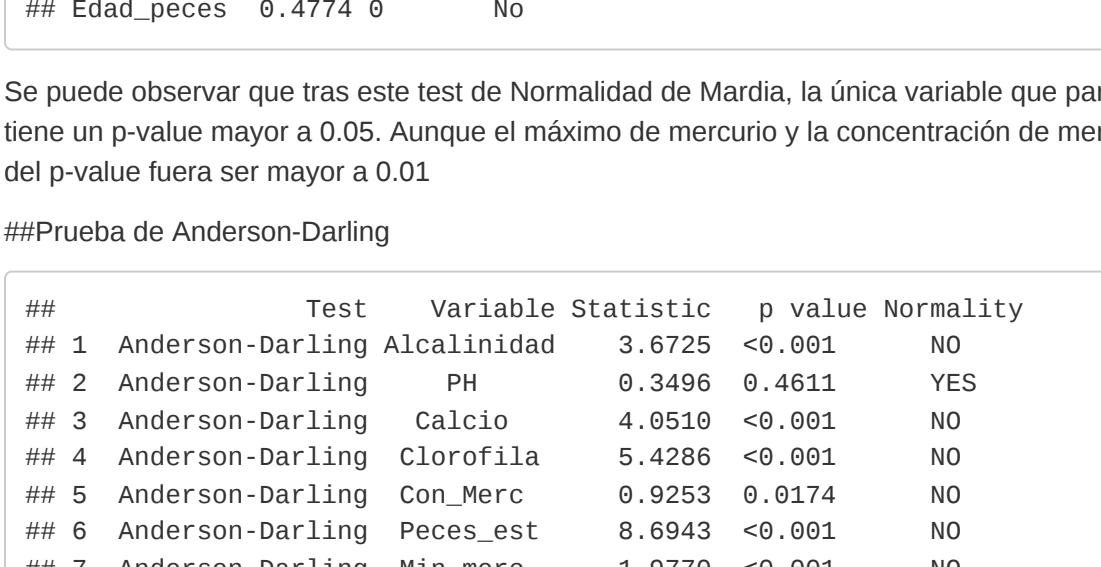
```
##          Lago Alcalinidad PH Calcio Clorofila Con_Merc Peces_est Min_merc
## 1 Alligator  5.9 6.1 3.0 0.1 2.3 5 0.8 1.23 5 0.8 1.23 5 0.85
## 2 Annie      3.5 5.1 1.9 3.2 1.33 7 0.92 1.90 1.33 0
## 3 Apopka     116 0 9.1 44.1 128 0.3 4 8.1 5 8 0.4 0.4 6 0.04
## 4 Blue Cypress 39.4 6.9 16.4 3.5 0.44 12 0.13 0.8 0.4 0.4 12 0.13
## 5 Brick      2.5 4.6 2.9 1.0 1.20 12 0.69 1.50 1.20 12 0.69
## 6 Bryant     19.6 7.3 4.5 44.1 0.27 14 0.04 0.27 14 0.04
## Max_merc Est_conc Edad_peces
## 1 1.43 0.53 1
## 2 1.90 1.33 0
## 3 0.06 0.04 0
## 4 0.84 0.44 0
## 5 1.50 1.33 1
## 6 0.48 0.25 1
```

```
## [1] "Alligator" "Annie" "Apopka"
## [4] "Blue Cypress" "Brick" "Bryant"
## [7] "Cherry" "Crescent" "Deer Point"
## [10] "Dias" "Dorr" "Down"
## [13] "Eatam" "East Tohopekaliga" "Fara-13"
## [16] "George" "Griffin" "Harney"
## [19] "Hart" "Hatchcreek" "Tamonia"
## [22] "Istokpoga" "Jackson" "Josephine"
## [25] "Kingsley" "Kissimmee" "Lochloosa"
## [28] "Louisia" "Mizockskuee" "Minimela"
## [31] "Monroe" "Newmans" "Ocean Pond"
## [34] "Ocheese Pond" "Okeechobee" "Orange"
## [37] "Panasafokee" "Parker" "Placid"
## [40] "Puzzle" "Rudman" "Rousseau"
## [43] "Swampson" "Snipp" "Talliquin"
## [46] "Tarpon" "Tohopekaliga" "Trafford"
## [49] "Trout" "Tsala Apopka" "Weir"
## [52] "Wildcat" "Yale"

## [1] 53
```

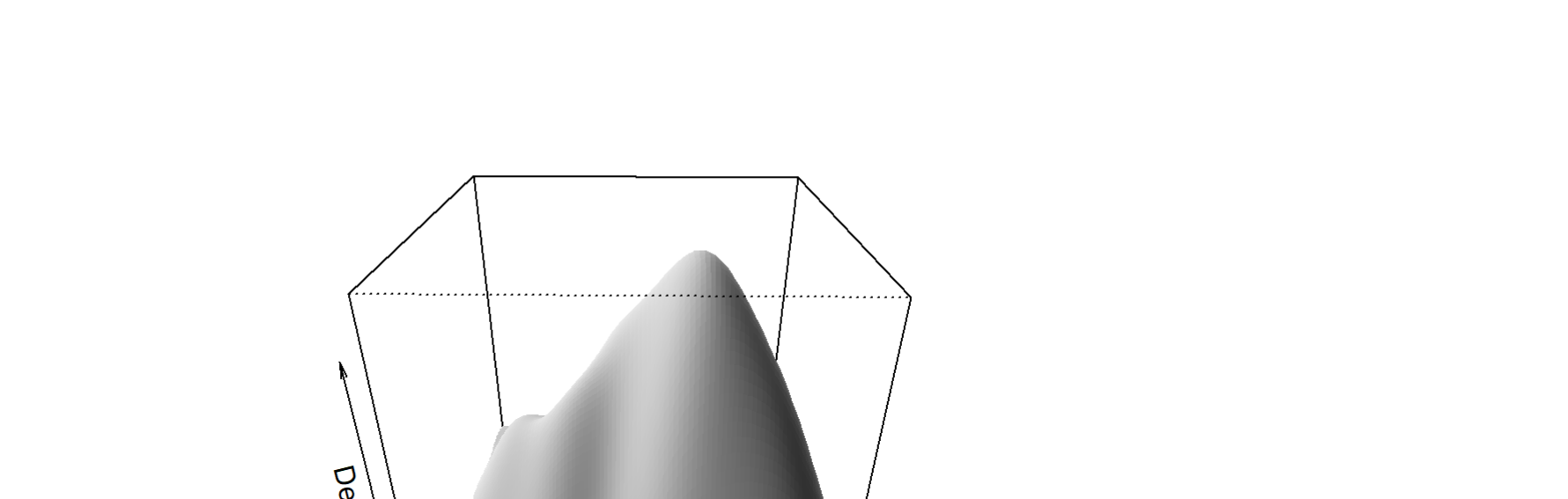
Hay 53 registros y se tienen 53 diferentes lagos, por lo que sabemos que esta columna realmente no nos aporta mucho.

Distribución de la concentración de mercurio



Se puede observar que hay bastantes casos en los que la concentración de mercurio por kg de pez es mayor a 0.5, lo cual incumple el Reglamento 34687-MAG y los reglamentos internacionales CE 1881/2006 y Codex Standard 193-1995.

De hecho, la mediana es 0.48 y la moda es de 0.527. Esto se encuentra incluso ya pasando los límites de lo establecido.



#Prueba de Normalidad de Mardia

```
## $mv.test
## Test Statistic p-value Result
## 1 Skewness 592.6673 0 NO
## 2 Kurtosis 5.8768 0 NO
## 3 MV Normality <NA> <NA> NO

## $uv.shapiro
## W p-value UV.Normality
## Alcalinidad 0.8203 0 NO
## PH 0.901 0.5552 YES
## Calcio 0.7913 0 NO
## Clorofila 0.6817 0 NO
## Con_Merc 0.9421 0.0125 NO
## Con_Merc 0.383 0 NO
## Min_merc 0.877 1e-04 NO
## Max_merc 0.9555 0.0467 NO
## Est_conc 0.9258 0.0028 NO
## Edad_peces 0.4774 0 NO
```

Se puede observar que al traza este test de Normalidad de Mardia, la única variable que parece que se comporta con normalidad es el PH, pues tiene un p-value mayor a 0.05. Aunque el máximo de mercurio y la concentración de mercurio podrían ser considerados normales si el criterio del p-value fuera ser mayor a 0.01

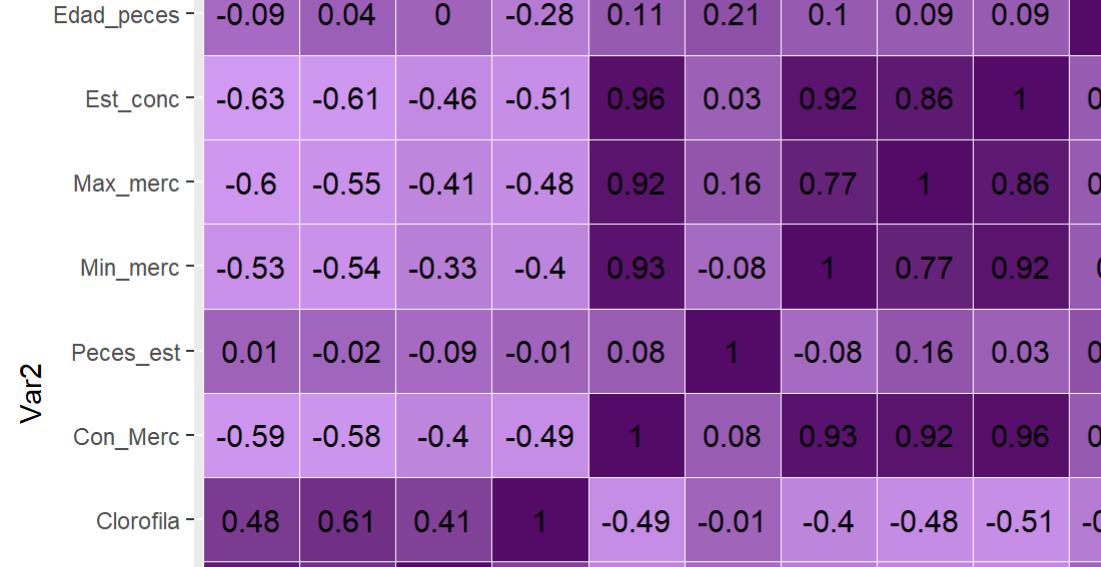
#Prueba de Anderson-Darling

```
## Test Variable Statistic p value Normality
## 1 Anderson-Darling Alcalinidad 3.6725 <0.001 NO
## 2 Anderson-Darling PH 0.3496 0.4611 YES
## 3 Anderson-Darling Calcio 4.0510 <0.001 NO
## 4 Anderson-Darling Clorofila 5.4286 <0.001 NO
## 5 Anderson-Darling Con_Merc 0.9253 0.0174 NO
## 6 Anderson-Darling Peces_est 8.6943 <0.001 NO
## 7 Anderson-Darling Min_merc 1.9770 <0.001 NO
## 8 Anderson-Darling Max_merc 0.6585 0.001 YES
## 9 Anderson-Darling Est_conc 1.0469 0.0086 NO
## 10 Anderson-Darling Edad_peces 14.3350 <0.001 NO
```

Aquí podemos observar que en tras la prueba de Anderson Darling, las únicas variables que parece que muestran normalidad son el PH (el cual viene en la prueba de Mardia) y el máximo de mercurio.

```
## Test Statistic p-value Result
## 1 Skewness 592.6673 0 NO
## 2 Kurtosis 5.8768 0 NO
## 3 MV Normality <NA> <NA> NO
```

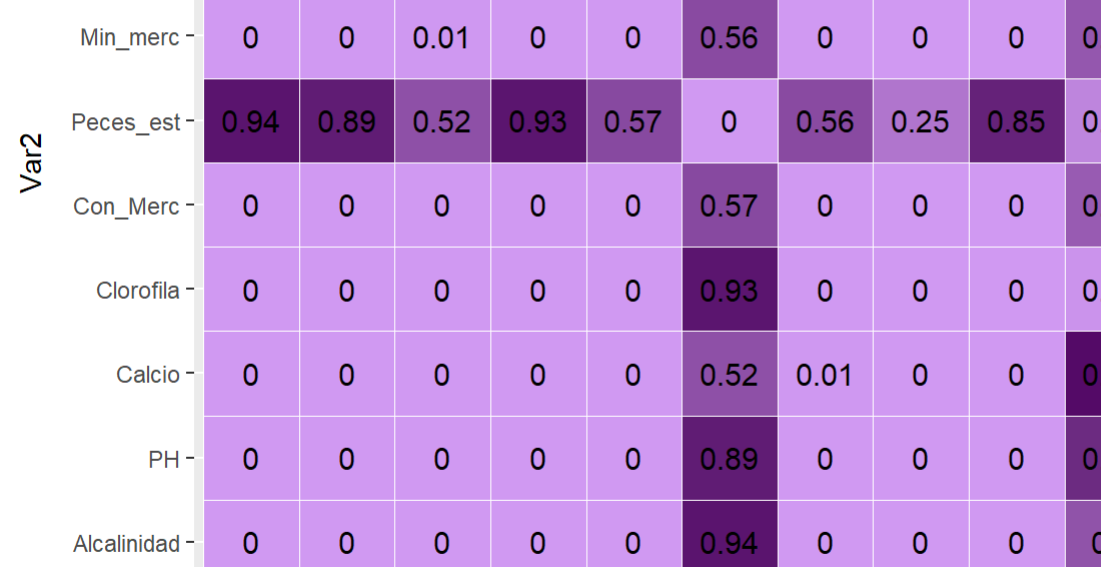
A partir de la prueba de Mardia, podemos saber que no se tiene normalidad multivariada pues los p-valores del sesgo y de la kurtosis son ambos prácticamente 0.



```
## $multivariateNormality
## Test Statistic p value Result
## 1 Mardia Skewness 6.53855430534145 0.162377362354508 YES
## 2 Mardia Kurtosis -0.889321233851276 0.373830462900813 YES
## 3 MVN <NA> <NA> YES

## $univariateNormality
## Test Variable Statistic p value Normality
## 1 Anderson-Darling PH 0.3496 0.4611 YES
## 2 Anderson-Darling Max_merc 0.6585 0.001 YES

## $descriptives
## n Mean Std.Dev Median Min Max 25th 75th Skew
## PH 53 6.5905660 1.2884493 6.80 3.60 9.10 5.80 7.40 -0.2458771
## Max_merc 53 0.8745283 0.5220469 0.84 0.06 2.04 0.48 1.33 0.4645925
## PH -0.6239638
## Max_merc -0.6692490
```

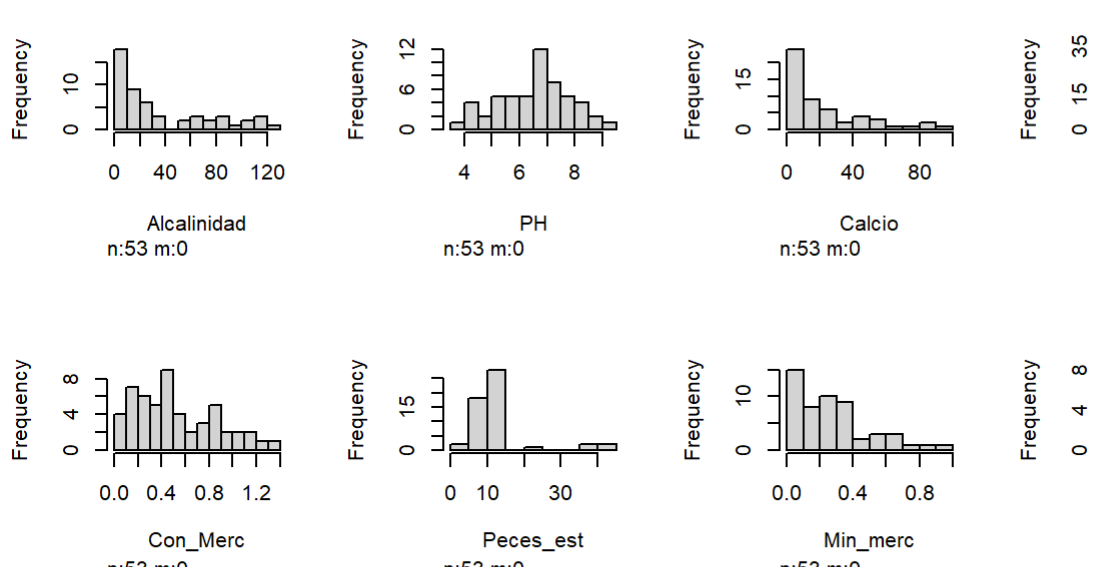


```
## $multivariateNormality
## Test Statistic p value Result
## 1 Mardia Skewness 6.53855430534145 0.162377362354508 YES
## 2 Mardia Kurtosis -0.889321233851276 0.373830462900813 YES
## 3 MVN <NA> <NA> YES

## $univariateNormality
## Test Variable Statistic p value Normality
## 1 Anderson-Darling PH 0.3496 0.4611 YES
## 2 Anderson-Darling Max_merc 0.6585 0.001 YES

## $descriptives
## n Mean Std.Dev Median Min Max 25th 75th Skew
## PH 53 6.5905660 1.2884493 6.80 3.60 9.10 5.80 7.40 -0.2458771
## Max_merc 53 0.8745283 0.5220469 0.84 0.06 2.04 0.48 1.33 0.4645925
## PH -0.6239638
## Max_merc -0.6692490
```

Outliers



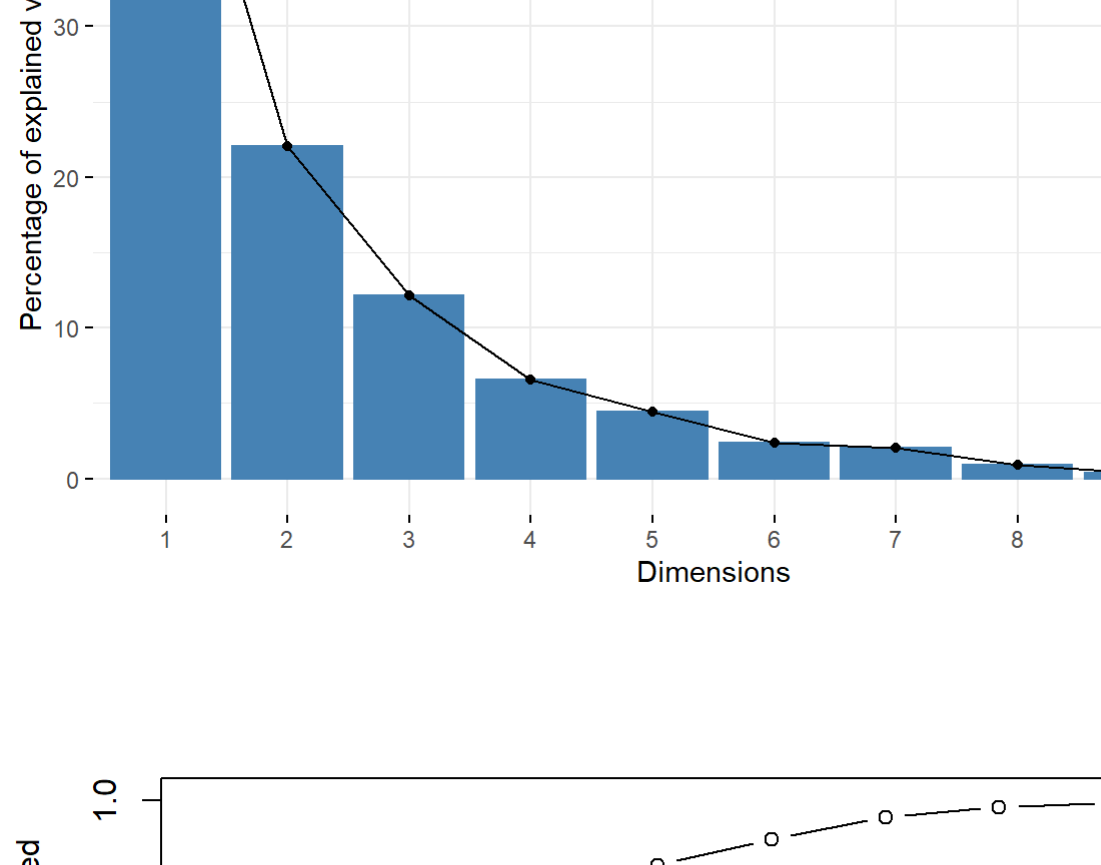
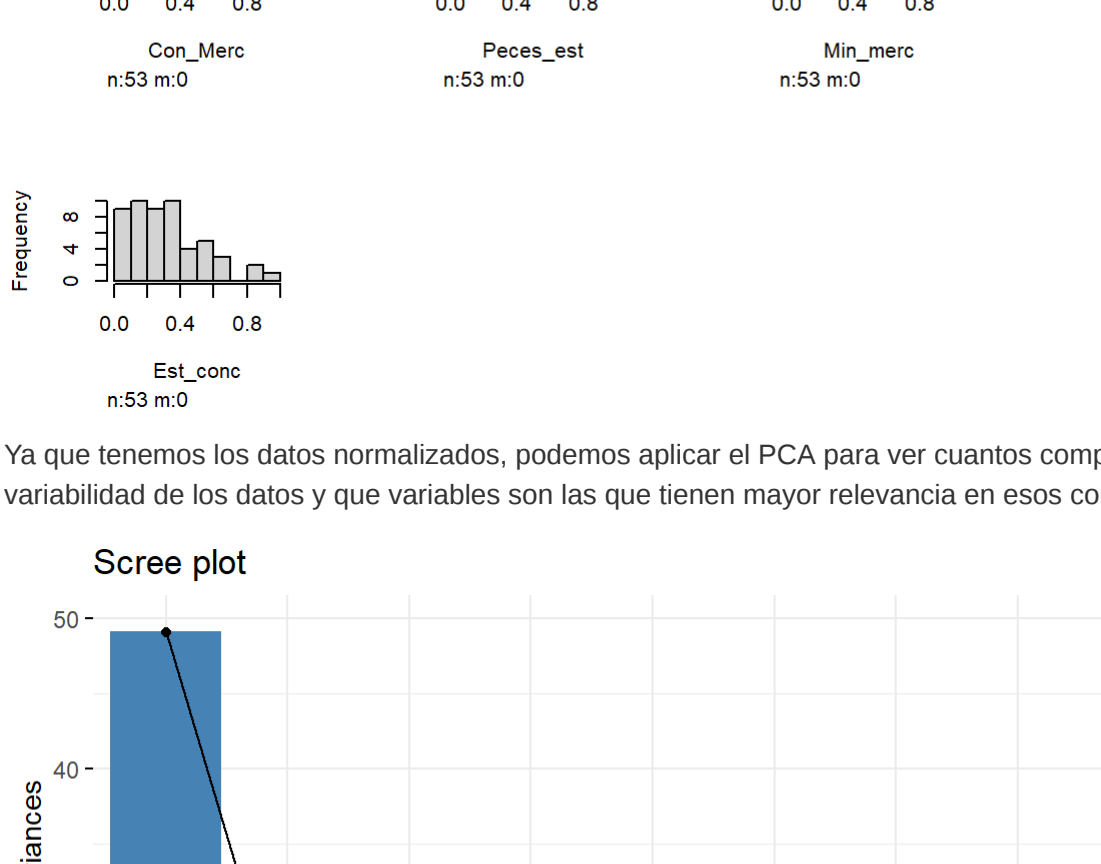
[1] 24 33

Análisis de Correlación

Para conocer los factores que más pueden influir hacemos un análisis de la correlación entre las variables a examinar.

Matriz de Correlación

En la siguiente figura se puede observar una matriz de correlación entre las variables. Se puede observar que, para la concentración de mercurio, se tiene una correlación negativa con la Clorofila, el Calcio, el PH y la Alcalinidad.



observar en el primer heatmap como es que las variables tienen correlación entre sí. Para poder entender mejor la relación entre las variables se han calculado los p-valores entre las correlaciones de todas las variables en el segundo heatmap. Después se han usado pruebas de hipótesis para ver si las variables tienen independencia o asociación entre sí.

$H_0: \rho = 0$ El caso de independencia entre variables $H_1: \rho \neq 0$ El caso de asociación entre variables

Regla de decisión: Rechazar H_0 si el p-value < 0.05.

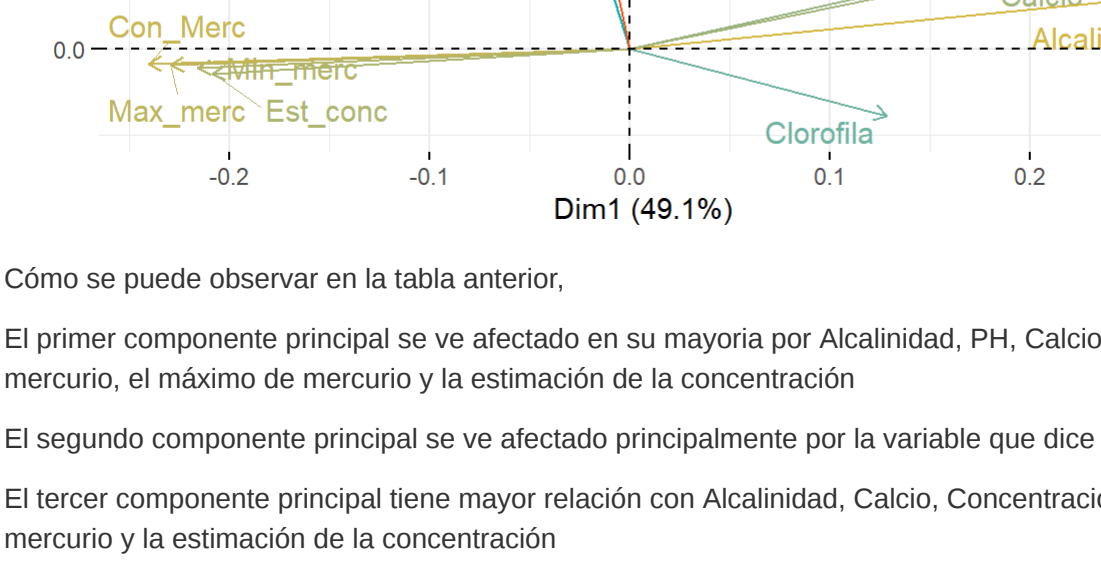
Dado que queremos ver que variables afectan a la concentración de mercurio, podemos observar que Alcalinidad, PH, Calcio, Clorofila, el mínimo de mercurio, el máximo de mercurio y la estimación de mercurio tienen un p-value menor a 0.05, por lo que se rechaza H_0 y se llega a la conclusión que estas variables son las que pueden afectar a la concentración de mercurio.

Análisis de Componentes Principales

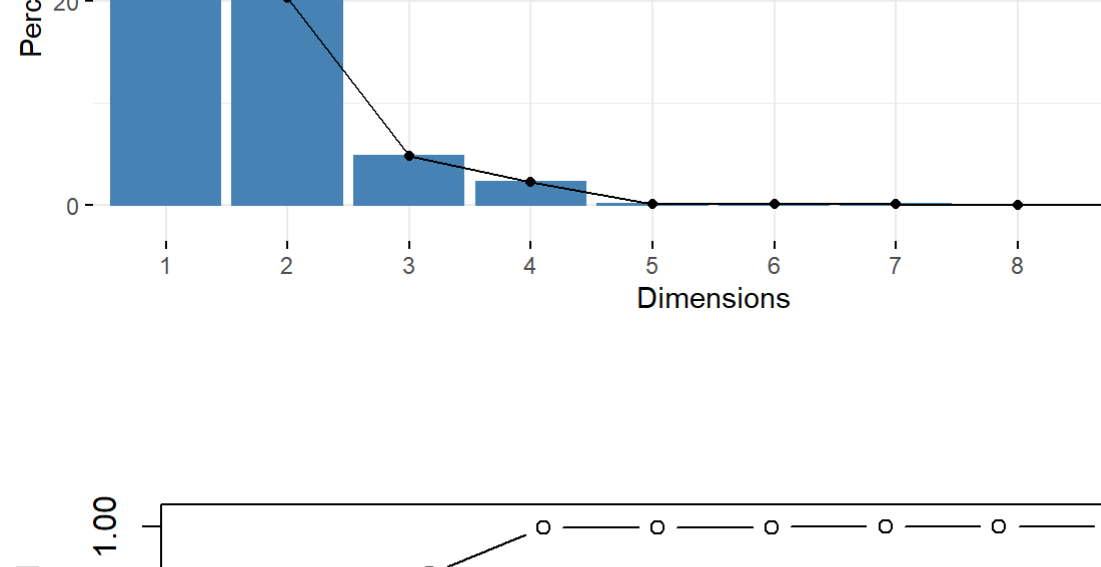
Antes de hacer el PCA debemos de saber si es necesario escalar los datos. Para poder saber esto primero vamos a ver las distribuciones de las variables.



Podemos observar que algunas variables que tienen valores pequeños, mientras que algunas otras tienen valores relativamente mucho más grandes, por ello vamos a buscar normalizar los datos para que cuando se realice el PCA no se tenga algún sesgo hacia las variables con valores más altos.



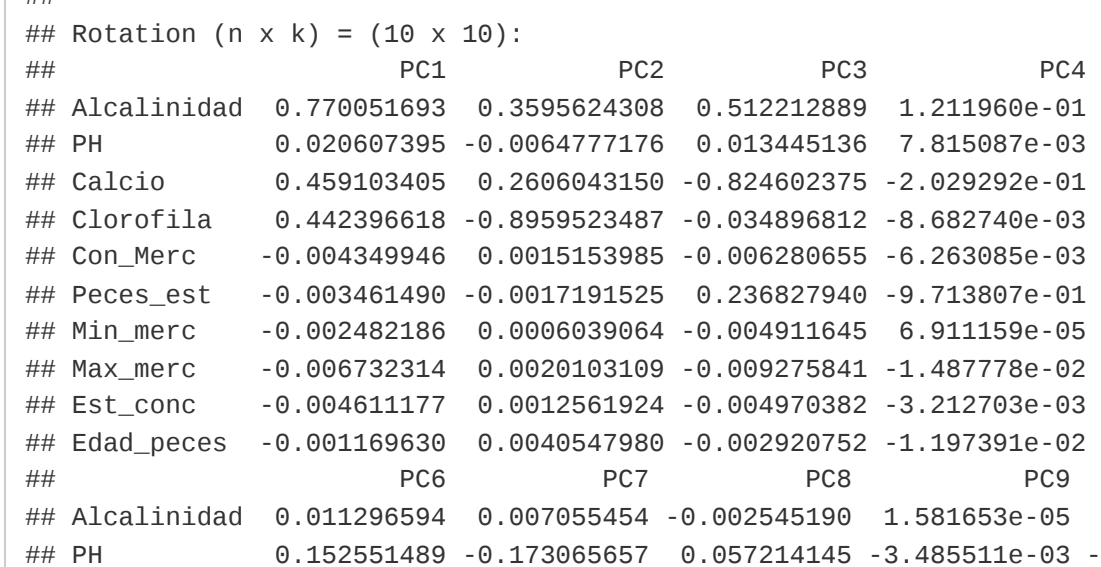
Ya que tenemos los datos normalizados, podemos aplicar el PCA para ver en cuáles componentes principales representan la mayoría de la variabilidad de los datos y que variables son las que tienen mayor relevancia en esos componentes.



Cómo se puede observar con 4 componentes principales se tiene cerca del 90% de la variabilidad de los datos. Ahora solo falta observar esos componentes principales.

```
## Standard deviations (1, ..., p)=10:
## [1] 0.59681084 0.39917379 0.29659562 0.21724577 0.17846501 0.13028089
## [7] 0.12054968 0.08034846 0.05387504 0.0352391

##
## Rotation (n x k) = (10 x 10):
## PC1 PC2 PC3 PC4 PC5
## Alcalinidad 0.2599910 0.0711931 -0.45380835 -0.092834387 0.22948361
## PH 0.30861238 0.1482087 -0.17878490 -0.069814950 0.55957988
## Calcio 0.31853493 0.11186139 0.52405512 0.09361573 0.28519334
## Clorofila 0.21582602 0.09564014 0.08869493 -0.14252997 -0.70838086
## Con_Merc -0.48348016 -0.0959523487 -0.034896812 -0.862740e-03 0.014101339
## Peces_est -0.02861474 -0.14242070 0.04076418 0.928389977 0.12533319
## Max_merc 0.003435490 0.001791525 0.236827940 0.7139897e-01 0.003620863
## Min_merc -0.002482186 0.006039964 0.004911645 6.91159e-05 0.06172060
## Max_merc -0.006703314 0.002010310 -0.009275841 -1.487778e-02 0.076192576
## Est_conc -0.004611177 0.0012561924 -0.004970382 3.212703e-03 0.079787785
## Edad_peces -0.001969330 0.0049547980 -0.002920752 -1.197391e-02 -0.188987340
```



Cómo se puede observar en la tabla anterior, El primer componente principal se ve afectado en su mayoría por Alcalinidad, PH, Calcio, Clorofila, la concentración de mercurio, el mínimo de mercurio, el máximo de mercurio y la estimación de la concentración

El segundo componente principal se ve afectado principalmente por la variable que dice si un pez es joven o adulto

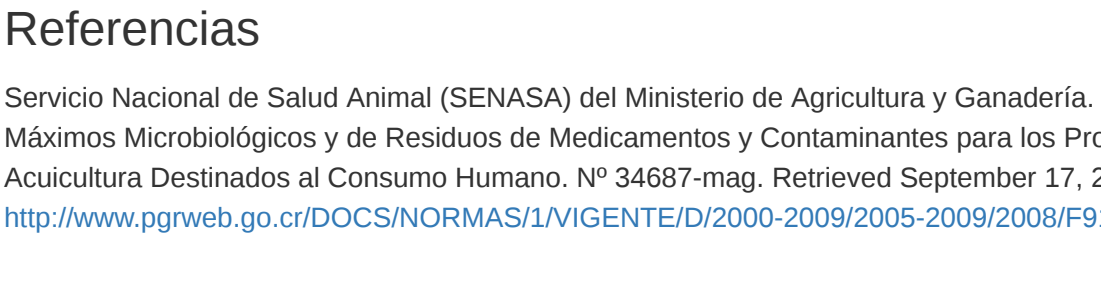
El tercer componente principal tiene mayor relación con Alcalinidad, Calcio, Concentración de mercurio, el mínimo de mercurio, el máximo de mercurio y la estimación de la concentración

El cuarto componente principal se ve afectado en su mayoría por la cantidad de peces estudiados.

Esto es como se puede ver gráficamente en el loading plot anterior. Se puede ver como es que esas variables se van agrupando.

Por su parte, haré el mismo análisis pero esta vez sin la normalización de los datos. Esto con objetivo de ver si se tienen resultados similares.

Ya que tenemos los datos normalizados, podemos aplicar el PCA para ver en cuáles componentes principales representan la mayoría de la variabilidad de los datos y que variables son las que tienen mayor relevancia en esos componentes.



Cómo se puede observar con 4 componentes principales se tiene cerca del 90% de la variabilidad de los datos. Ahora solo falta observar esos componentes principales.

```
## Standard deviations (1, ..., p)=10:
## [1] 0.33531760 0.16983350 0.06774799 0.04681655

##
## Rotation (n x k) = (10 x 10):
## PC1 PC2 PC3 PC4 PC5
## Alcalinidad 0.2599910 0.0711931 -0.45380835 -0.092834387 0.22948361
## PH 0.30861238 0.1482087 -0.17878490 -0.069814950 0.55957988
## Calcio 0.31853493 0.11186139 0.52405512 0.09361573 0.28519334
## Clorofila 0.21582602 0.09564014 0.08869493 -0.14252997 -0.70838086
## Con_Merc -0.48348016 -0.0959523487 -0.034896812 -0.862740e-03 0.014101339
## Peces_est -0.02861474 -0.14242070 0.04076418 0.928389977 0.12533319
## Max_merc 0.003435490 0.001791525 0.236827940 0.7139897e-01 0.003620863
## Min_merc -0.002482186 0.006039964 0.004911645 6.91159e-05 0.06172060
## Max_merc -0.006703314 0.002010310 -0.009275841 -1.487778e-02 0.076192576
## Est_conc -0.004611177 0.0012561924 -0.004970382 3.212703e-03 0.079787785
## Edad_peces -0.001969330 0.0049547980 -0.002920752 -1.197391e-02 -0.188987340
```


Cómo se puede observar en las gráficas y tabla anterior, parece que el pca si se ve muy sesgado por las variables que tienen mayores rangos, por lo que la normalización que realizamos fue un acierto.

Conclusiones

A partir del análisis de las correlaciones entre las variables y su significancia, además del análisis de componentes principales con los datos normalizados, se llega a la conclusión que los principales factores que influyen en el nivel de contaminación por mercurio en los peces de lagos de Florida son:

Alcalinidad, PH, Calcio, Clorofila, el mínimo de mercurio, el máximo de mercurio y la estimación de mercurio.

Algo que se debe de tomar en consideración es que el mínimo de mercurio, el máximo de mercurio y la estimación de mercurio, por la forma en la que fueron obtenidas y calculadas, tienen bastante relación con la concentración de mercurio, por lo que considero que la Alcalinidad, el PH, el Calcio y la Clorofila son los factores más relevantes.

Referencias

Servicio Nacional de Salud Animal (SENASA) del Ministerio de Agricultura y Ganadería. (2008). RTRC 409: 2008 Reglamento de Límites Máximos Microbiológicos y de Residuos de Medicamentos y Contaminantes para los Productos y Subproductos de la Pesca y de la Acuicultura Destinados al Consumo Humano. N° 34687-mag. Retrieved September 17, 2022, from <http://www.pgweb.gob.cr/DOCS/NORMAS/1/VIGENTE/DE/2000-2009/2005-2009/2008/R916BED07.HTML>