

# Analítica de datos aplicada a estudios sobre desarrollo

## Segundo examen

**Giovany Babativa-Márquez, PhD**

**Andrea González Peña, PhD**

Noviembre 03, 2023

## Instrucciones

El trabajo puede ser desarrollado como máximo en parejas y se debe entregar el script reproducible o archivo de R utilizado para generar los resultados, así como un documento de máximo 13 páginas que contenga como mínimo:

- Resumen
- Introducción
- Materiales y métodos
- Resultados
- Conclusiones
- Bibliografía

El contenido de este documento debe incluir el desarrollo de los ejercicios propuestos más adelante, lo cual forma parte de su estructura retórica. Es decir, no debe entregar por separado el desarrollo de cada ejercicio, sino que todos estos elementos deben servir para construir el documento que debe entregar.

La fecha límite para la entrega será el 18 de noviembre y se debe cargar por la plataforma de Bloque Neón. Los trabajos enviados por correo electrónico no serán considerados para su evaluación.

Tenga en cuenta que se hará una comparación entre los trabajos y de encontrar una alta similitud entonces la nota será dividida entre el número de personas involucradas. Si el reporte tiene más de las páginas permitidas (sin contar la bibliografía) tendrá una penalidad de 5 décimas por hoja adicional.

## 2 Ejercicios propuestos

El Proyecto de Opinión Pública Latinoamericana (LAPOP por sus siglas en inglés), es una encuesta realizada desde el año 2004 que busca medir los valores, actitudes y comportamientos democráticos. Para el año 2015 ya se realizaba en 28 países con un tamaño de muestra de más de 50.000 encuestas y considera un diseño muestral probabilístico en cada país. La ronda de 2021, corresponde al último estudio realizado y se llevó a cabo en 22 países con más de 64.000 encuestas.

Para este examen se han descargado los datos de 2993 encuestas realizadas en Colombia en la ronda 2021 y el cuestionario utilizado, los cuales fueron descargados de la página de LAPOP en este [enlace](#). Asimismo, los estudiantes pueden acceder a los materiales de forma directa desde el repositorio de GitHub de este curso:

- *Ficha técnica*: haga clic [aquí](#)
- *Cuestionario*: haga clic [aquí](#)
- *Conjunto de datos en formato stata*: haga clic [aquí](#)

### 2.1 Ejercicio 1

En este punto se hará un trabajo conceptual, imagine que se le ha solicitado construir tres (3) índices:

- Índice de confianza en instituciones
- Índice de democracia
- Índice de valores antidemocráticos

Para ello debe descargar el cuestionario, y para cada índice defina el concepto que desea medir e identifique, desde el cuestionario, las variables que usaría en cada caso.

### 2.2 Descripción de los datos

El siguiente paso consiste en construir la base de datos. Para ello, el conjunto de datos original fue preprocesado para importar, ordenar y transformar las variables que serán relevantes en nuestro análisis, el script utilizado puede ser descargado de [acá](#). Este conjunto de datos contiene 16 variables construidas desde las preguntas de la encuesta LAPOP, en el que finalmente se tuvieron en cuenta 2971 encuestas y puede ser descargado [acá](#). A continuación se describen las variables utilizadas.

Nombre	Descripción	Preguntas
<b>just_golpe</b>	Circunstancias en que se justificaría que los militares de este país tomen el poder por un golpe de Estado.	<i>jc13, jc13covid</i>

Nombre	Descripción	Preguntas
<b>just_cierre_cong</b>	Justificación para que el presidente del país cierre el Congreso y gobierne sin Congreso	<i>jc15a</i>
<b>conf_gobierno_nal</b>	Mide en una escala del 1 al 4 la confianza en que el gobierno nacional hace lo correcto.	<i>anestg</i> en escala invertida
<b>conf_instit</b>	Mide en una escala del 1 al 7 el nivel de respeto a las instituciones políticas en el país del encuestado.	<i>b2</i>
<b>conf_alcaldia</b>	Mide en una escala de 1 a 7 el nivel de confianza en la Alcaldía.	<i>b32</i>
<b>conf_elecciones</b>	Mide en una escala de 1 a 7 el nivel de confianza en las elecciones	<i>b47a.</i>
<b>conf_policia</b>	Mide en una escala de 1 a 7 el nivel de confianza en la policía	<i>b18</i>
<b>conf_medios</b>	Mide en una escala del 1 al 7 el nivel de confianza en los medios de comunicación	<i>b37</i>
<b>conf_fuerzas_mil</b>	Mide en una escala del 1 al 7 el nivel de confianza en las fuerzas militares.	<i>b12</i>
<b>sat_democracia</b>	Variable dicotómica que mide el nivel de satisfacción con la democracia de los encuestados.	<i>pn4</i> en escala invertida
<b>prot_derechos</b>	Mide en una escala del 1 al 7 la protección de los derechos básicos del ciudadano por parte del sistema político colombiano	<i>b3</i>
<b>orgullo_sistema</b>	Mide en una escala del 1 al 7 qué tanto se siente usted orgulloso de vivir bajo el sistema político colombiano	<i>d4</i>

Recuerde cargar en R los paquetes necesarios y el conjunto de datos, así:

```
rm(list = ls())

options(scipen = 999)
library(pacman)

p_load(tidyverse, janitor, corrplot, haven,
       devtools, FactoMineR, factoextra,
       ggcorrplot, GGally)

url <- "https://github.com/jgbabativam/AnaDatos/raw/main/datos/df_colombia.dta"

datos <- read_dta(url)
```

A partir de la base de datos construida desarrolle los siguientes ejercicios:

2. Analice la relación lineal entre las variables, para ello calcule la matriz de correlación entre las variables cuantitativas y represente el resultado con un diagrama usando el paquete `corrplot` o `ggcorrplot`. Concluya sobre estos resultados.

Realice un Análisis de Componentes Principales sobre las variables cuantitativas:

3. Discuta sobre la cantidad de información del conjunto de los datos es explicado por las dos primeras dimensiones. Presente un gráfico que lo ilustre. Ayuda use la función `fviz_screplot()` del paquete `factoextra`.
4. Haga una gráfica del plano factorial generado por las dimensiones 1 y 2 para las variables. Concluya sobre la asociación de las variables indicando cuáles presentan fuertes correlaciones y apuntan en la misma dirección representando el mismo concepto, indique también si considera que hay variables que no se representen bien en el plano de las primeras dos dimensiones. Ayuda use la función `fviz_pca_var()` del paquete `factoextra`.
5. Use el siguiente comando para visualizar de forma simultánea a los encuestados y las variables

```
fviz_pca_biplot(res, repel = F, col.var = "black", col.ind = "gray")
```

Tenga en cuenta que `res` se refiere al objeto que contiene los resultados del PCA. De forma visual, considerando la densidad de punto grises (encuestados), ¿podría conjeturar que hay una buena confianza en las instituciones? explique.

6. A partir de los resultados, explique que representa un puntaje alto en la primera componente principal. Asigne un nombre apropiado al índice que se obtendría desde esta dimensión. Apoye su conclusión en el gráfico que resulta del siguiente comando (tenga en cuenta que `res` se refiere al objeto que contiene los resultados del PCA):

```
fviz_contrib(res, choice = "var", axes = 1, top = 10)
```

7. A partir de los resultados, explique que representa un puntaje alto en la segunda componente principal. Asigne un nombre al índice que se obtendría desde esta dimensión. Apoye su conclusión en el gráfico que resulta del siguiente comando (tenga en cuenta que `res` se refiere al objeto que contiene los resultados del PCA):

```
fviz_contrib(res, choice = "var", axes = 2, top = 10)
```

8. A partir de los resultados, explique que representa un puntaje alto en la tercera componente principal. Asigne un nombre al índice que se obtendría desde esta dimensión. Apoye su conclusión en el gráfico que resulta del

siguiente comando (tenga en cuenta que `res` se refiere al objeto que contiene los resultados del PCA):

```
fviz_contrib(res, choice = "var", axes = 3, top = 10)
```

9. Construya los índices formados por las tres (3) primeras componentes principales, debe hacer uno por cada componente principal y agréguelos al conjunto de datos. Ayuda: use el siguiente código para hacerlo:

#### ##### Indices

```
index1 <- as.data.frame(res$ind$coord[,1]) |>
  rename(score = `res$ind$coord[, 1]`) |>
  mutate(Indice1 = round(GGally::rescale01(score)*100, 1)) |>
  select(-score)

index2 <- as.data.frame(res$ind$coord[,2]) |>
  rename(score = `res$ind$coord[, 2]`) |>
  mutate(Indice2 = round(GGally::rescale01(score)*100, 1)) |>
  select(-score)

index3 <- as.data.frame(res$ind$coord[,3]) |>
  rename(score = `res$ind$coord[, 3]`) |>
  mutate(Indice3 = round(GGally::rescale01(score)*100, 1)) |>
  select(-score)

df_index <- bind_cols(datos, index1, index2, index3)
```

10. Realice los análisis descriptivos del valor de los índices por región, sexo, edad y nivel educativo. Para ello puede calcular el promedio o usar gráficas apropiadas que le permitan concluir si existe una mayor confianza en determinadas regiones, rangos de edad o niveles educativos.
11. Genere las conclusiones generales del ejercicio.