

Examen, problema 1

EYP2435 - Análisis Multivariado

Sebastián Celaya

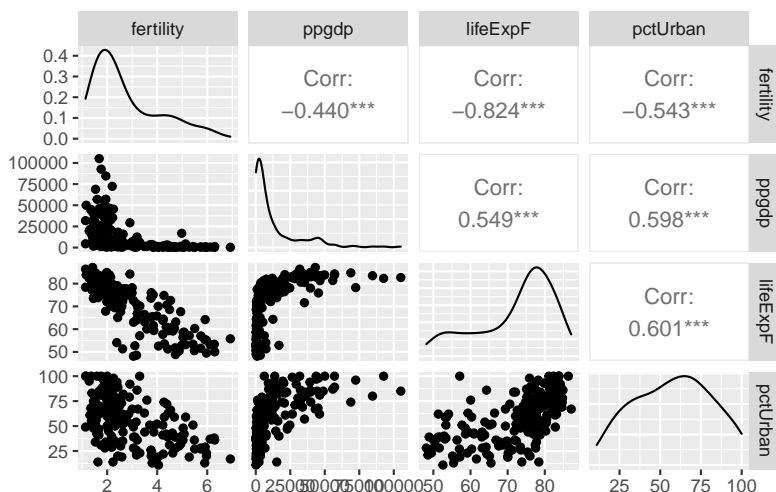
Camila Echeverría

Datos de naciones Unidas, UN. Los datos en el archivo `UNData`, en `canvas` contiene algunas variables, entre ellas `ppgdp`, el producto nacional bruto por persona de 2009 en dólares estadounidenses, `fertility`, la tasa de natalidad por cada 100 mujeres en la población en el año 2009, `lifeExp` esperanza de vida y `pctUrban` porcentaje de población urbana; además de dos variables cualitativas, región y grupo. Los datos son para 199 localidades, en su mayoría Países miembros de la ONU, pero también para otras áreas como Hong Kong que no son países independientes.

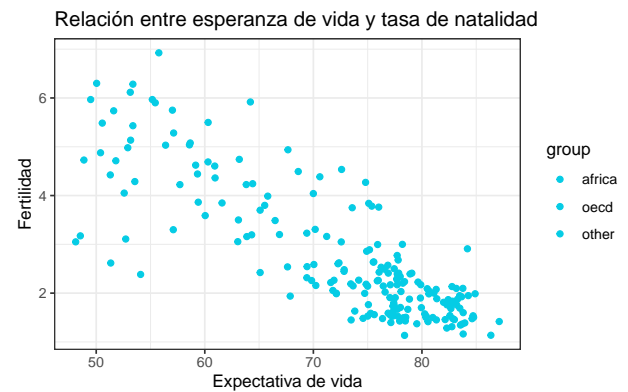
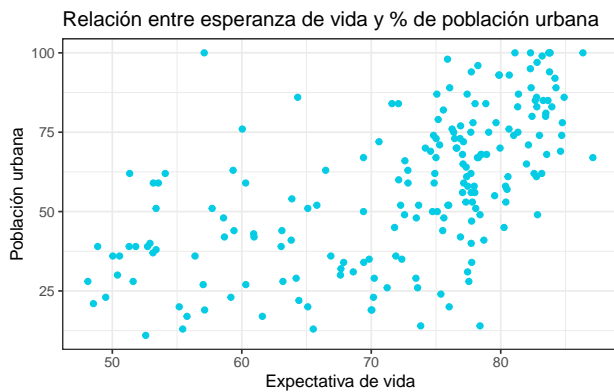
- Haga un análisis exploratorio de los datos. Escriba sus conclusiones.
- Suponga que se desea analizar el efecto de las otras variables disponibles en (`fertility`, `lifeEXP`). Para esto, proponga un modelo estadístico y verifique si es adecuado. Usando algunas herramientas de inferencia estadística, responda si existe efecto de las otras variables sobre (`fertility`, `lifeEXP`). Escriba sus conclusiones.

Solución a)

En primer lugar, veremos cómo se comportan las variables continuas entre sí. Para esto, realizaremos un gráfico de correlación.

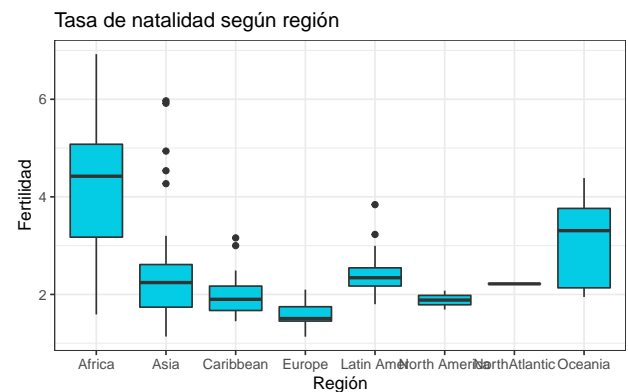
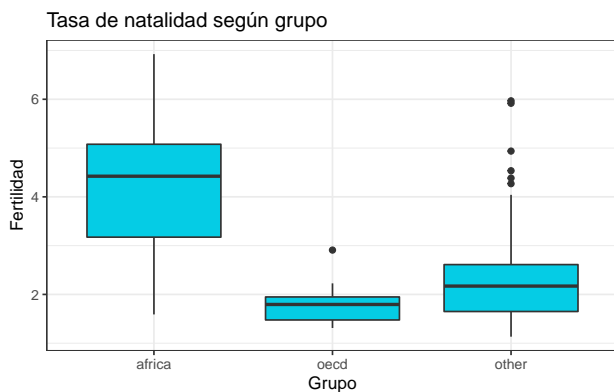


Las dos relaciones más fuertes que encontramos son entre `lifeExpF` y `fertility` y `lifeExpF` y `pctUrban`. En los siguientes gráficos podemos ver más de cerca cómo se relacionan estas variables.



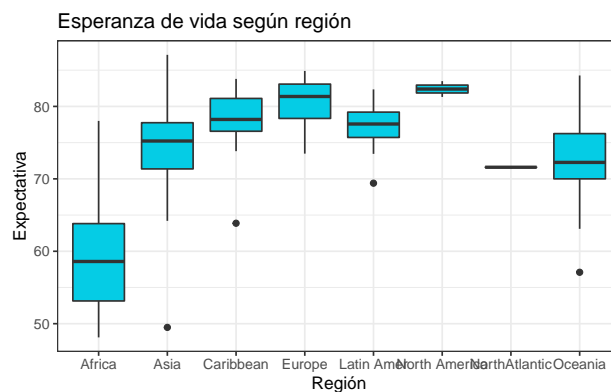
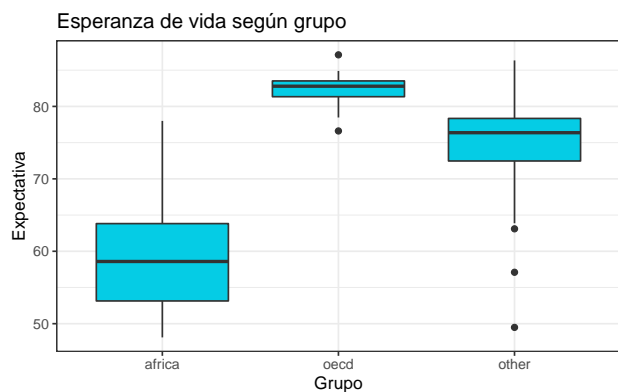
Por un lado, mientras la correlación es de 0.601, el gráfico no muestra una relación completamente lineal entre la esperanza de vida y el porcentaje de población urbana. De hecho, en algún punto hasta pareciera una simple nube de puntos. Sin embargo, la relación con la tasa de natalidad sí se ve bastante fuerte y bastante marcada, lo que era de esperarse de una correlación de -0.824.

Luego, debemos ver cómo se comportan las variables numéricas si las agrupamos por las variables categóricas. Para esto, comenzaremos realizando boxplots de **fertility**.

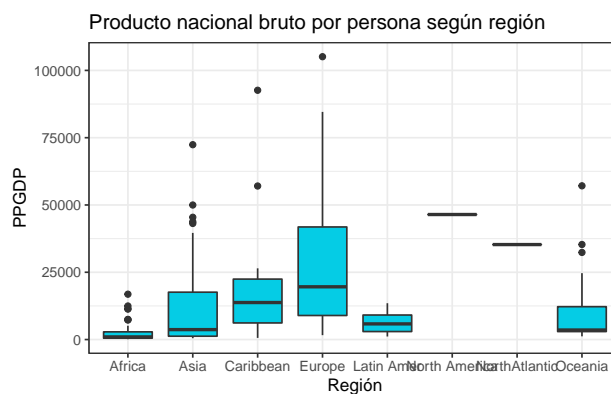
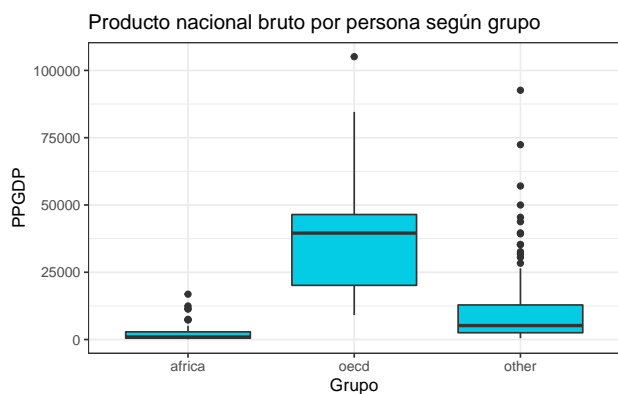


Lo primero que salta a la luz es que África, tanto como región como grupo, posee las tasas de natalidad más altas, siendo su mediana superior al tercer cuartil del resto de los grupos o regiones, aunque presenta un rango intercuartílico bastante amplio. Además las regiones pertenecientes a los grupos **oecd** u otros presentan tasas bastante bajas en comparación.

Un detalle que igual podría llamar la atención es la nula dispersión de la fertilidad de la región del Atlántico Norte. Esto se debe a que, según los datos, esta región está compuesta de un único país: Groenlandia.

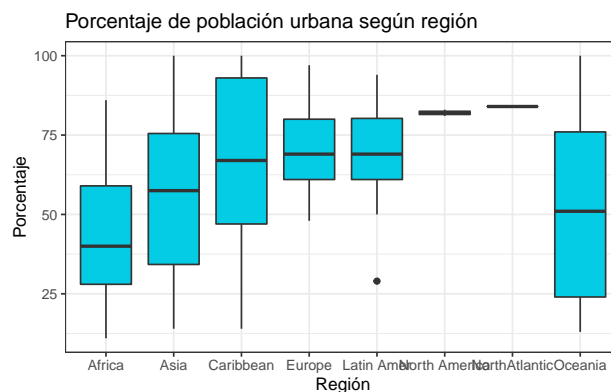
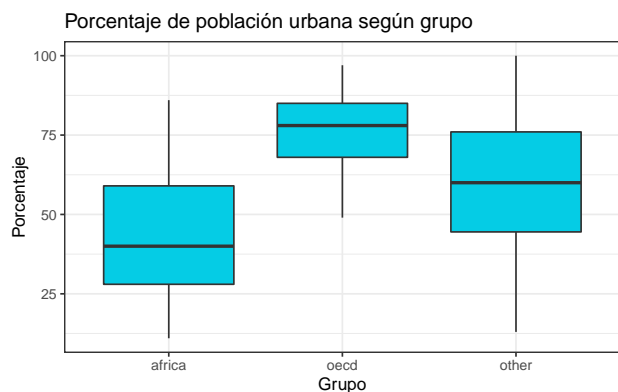


Luego, al realizar los gráficos de `lifeExpF`, podemos ver que la situación se revierte drásticamente. El grupo de la `oecd` es el que tiene mayor esperanza de vida, mientras que África es el que tiene la menor. En cuanto a las regiones, Norte América es la que lidera, seguida desde cerca por Europa y el Caribe.



Con respecto al producto nacional, se repite el patrón de la esperanza de vida: África posee el menor, mientras que el grupo de la `oecd` posee el mayor. Además, las medianas de América del Norte y Groenlandia son las más altas de entre las regiones.

Finalmente, podemos ver que el patrón se vuelve a repetir para el porcentaje de población urbana.



De esta manera, a priori, pareciera que el grupo y/o región tiene un gran impacto en las distintas variables cuantitativas de la base de datos, lo que, a su vez, determina cómo se comportarán estas.

Solución b)