

LOS ÍNDICES DE COLOMBIA

Santiago Serrano¹

^{1,2}Escuela de Ingeniería, Universidad de los Andes, s.serrano11@uniandes.edu.col

¹Instituto de altas investigaciones financieras, Banco del Parque,
delcurso@bp.com.col

05 de Julio de 2018

Abstract

El proposito principal de este trabajo es describir los procesos para partir una poblacion N-dimensional en partes de k tamaño en forma de una muestra. El proceso, que se denomina 'k-means' aparece para dar particiones que son razonablemente eficientes en el sentido de las variables dentro de las categorias. Eso es, que si p es la probabilidad de densidad para la poblacion S , $S=s_1, s_2, \dots, s_n$. La parte de E_n y de u_i siendo $i=1, 2, 3, \dots, k$ es el promedio condicional de p sobre S . Diremos de ahora en adelante en el documento 'tiende a ser bajo' para referirnos en principio a las consideraciones intuitivas y corroboradas del analisis matematico y practicas computacionales.

Introducción

The main purpose of this paper is to describe a process for partitioning an N-dimensional population into k sets on the basis of a sample. The process, which is called 'k-means', appears to give partitions which are reasonably efficient in the sense of within-class variance. That is, if p is the probability mass function for the population, $S = S_1, S_2, \dots, S_n$ is a partition of E_n . We say 'tends to be low' primarily because of intuitive considerations, corroborated to some extent by mathematical analysis and practical computational experience.

Comencemos viendo que hay en la sección 1 en la página 2.

1 Exploración Univariada

En esta sección exploro cada índice. En esta sección exploro cada índice. En esta sección exploro cada índice. En esta sección exploro cada índice. En esta sección exploro cada índice. En esta sección exploro cada índice. En esta sección exploro cada índice.

Para conocer el comportamiento de las variables se ha preparado donde se describe la distribución de las modalidades de cada variable. Los números representan la situación de algún país en ese indicador, donde el mayor valor numérico es la mejor situación.

Table 1: Medidas estadísticas

Statistic	Mean	St. Dev.	Max	Min	Median
IDH	0.802	0.042	0.879	0.691	0.804
Poblaci..n.Cabecera	1,196,730.000	1,982,287.000	10,070,801	13,090	717,197
Poblaci..n.Resto	360,590.300	331,887.600	1,428,858	21,926	268,111.5
Poblaci..n.Total	1,557,320.000	2,202,522.000	10,985,285	43,446	1,028,429

Como apreciamos en los países en la mejor situación son los menos, salvo en el caso del *índice de libertas mundial*¹

¹Nótese que esto se puede deber a la **menor** cantidad de categorías.

Para resaltar lo anterior, tenemos la Figura 1 en la página 3.

```

■barplots, echo=FALSE,fig=TRUE■= par(mfrow=c(2,2))
                                title='IDH' paleta='red'
hist(colbIDH, main = title, col = paleta, xlab = ' IDH', ylab = ' Frecuencia')
                                title='Poblaci<U+00F3>n de cabecera' paleta='red'
hist(colbpoblacion.cabecera, main = title, col = paleta, xlab = ' Poblaci <
                                U + 00F3 > n', ylab = " Frecuencia")
                                title='Poblaci<U+00F3>n resto' paleta='red' hist(colbPoblac, main =
                                title, col = paleta, xlab = ' Poblaci < U + 00F3 > n', ylab = " Frecuencia")
                                title='Poblaci<U+00F3>n total' paleta='red'
                                hist(colbpoblacion.total, main = title, col = paleta, xlab = ' Poblaci <
                                U + 00F3 > n', ylab = " Frecuencia")

```

Figure 1: Distribución de Indicadores

Además de la distribución de los variable, es importante saber el valor central. Como los valores son de naturaleza ordinal debemos pedir la **mediana** y otras medidas de posición (como los *cuartiles*, los que no pediremos pues son pocos valores). La mediana de cada variable la mostramos en la página.

■summary, fig=TRUE, echo=false■== dado el sesgo de las pobaciones, podriamos transformarla para que se acerque a la normalidad

```

par(mfrow=c(2,1))
colbcabeLog = log(colbpoblacion.cabecera) colbrestoLog = log(colbpoblacion.resto)
title='Logaritmo de la poblaci<U+00F3>n de cabecera' paleta='red' demoTableRelPlot=hist(colbcabeLog,
title, col = paleta, xlab = "", ylab = " Frecuencia")
title='Logaritmo de la poblaci<U+00F3>n restante' paleta='red' demoTableRelPlot=hist(colbrestoLog, m
title, col = paleta, xlab = "", ylab = " Frecuencia")
hist(colbcabeLog)hist(colbrestoLog)
Histogramas de los logaritmos

```

2 Exploración Bivariada

En este trabajo estamos interesados en el impacto de los otros índices en el nivel de Democracia. Veamos las relaciones bivariadas que tiene esta variable con todas las demás:

```
■corrDem, results=tex, echo=FALSE■=
  library(stargazer) explanans=names(colb)[c(7:8)] usando las logs corrDem=cor(colbIDH,colb[,explanans],
"na.or.complete")
  stargazer(corrDem )
  Veamos la correlación entre las variables independientes:
  ■corrTableX, results=tex, echo=FALSE■= y la correlaci<U+00F3>n entre
las variables independientes:
  corrTableX=round(cor(colb[,explanans], use = "na.or.complete"),2) corrTableX<sub>copy</sub> =
corrTableXcorrTableX[upper.tri(corrTableX)] < -""
  ver: corrTableX
  visualmente
  stargazer(corrTableX,title="Correlaci<U+00F3>n entre las variables inde-
pendientes")
```

Lo visto en la Tabla 1 se refuerza claramente en la Figura 2.

Figure 2: correlación entre predictores

3 Modelos de Regresión

Finalmente, vemos los modelos propuestos. Primero sin la libertad mundial como independiente, y luego con está. Los resultados se muestran en la Tabla ?? de la página ??.

```
■regresiones, echo=FALSE■= LinRegA = lm(IDH ~., data = colb[,c(1,7)])
LinRegB = lm(IDH ~., data = colb[,c(1,7:8)]) stargazer(LinRegA,LinRegB,title="Modelos
de Regresi<U+00F3>n",label="regresiones")
```

Como se vió en la Tabla ??, cuando está presente el *índice de libertad mundial*, el *índice de libertad de prensa* pierde significancia.

4 Exploración Espacial

Como acabamos de ver en la Tabla ?? en la página ??, si quisieras sintetizar la multidimensionalidad de nuestros indicadores, podríamos usar tres de las cuatro variables que tenemos (un par de las originales tiene demasiada correlación).

Así, propongo que calculemos conglomerados de países usando toda la información de tres de los indicadores. Como nuestras variables son ordinales utilizaremos un proceso de conglomeración donde las distancias serán calculadas usando la medida **gower** propuestas en [?], y para los enlazamientos usaremos la técnica de **medoides** según [?]. Los tres conglomerados se muestran en la Figura 3.

```
■getMap, echo=FALSE,results=hide■= Exploraci<U+00F3>n Espacial —
—
Calculemos conglomerados de regiones, usando toda la informaci<U+00F3>n
de las tres variables. usaremos la tecnica de k-means propuesta por MacQueen.
library(rgdal) folder='COLmaps' file = 'COLadm1.shp' mapaFile = file.path(folder, file) mapCol <
-rgdal :: readOGR(mapaFile, stringsAsFactors = F)
lo tenemos: plot(mapCol)
veamos que variables hay:
head(mapCol@data)
con esto hagamos el merge: subcolb = colb[,c(1 : 2, 7 : 8)] mapColidh =
merge(mapCol, subcolb, by.x = 'NAME1', by.y = 'Departamento', all.x = F)
cuántas regiones me quedaron luego del merge? nrow(mapColidh) todas!!...
preparacion para clusterizar:
que tengo?: names(mapColidh) nombredelavariablenqueusar < U+00E9 >:
dimensions = c("NAME1", "IDH", "cabeLog", "restoLog")
creo un nuevo data frame con esas: dataCluster=mapColidh@data[, c(dimensions)]
como la data es numerica la normalizo (menos la column 1): dataCluster[, -
1]=scale(dataCluster[, -1])
APLICANDO TECNICA KMEANS
calculo 3 clusters
resultado=kmeans(dataCluster[, -1], 3)
creo data frame con los clusters: clusters=as.data.frame(resultado$cluster)
a<U+00F1>ado columna con nombre de regiones clustersNAME1 = dataClusterNAME1names(clusters)
c('cluster', NAME1) hagoelmergehaciaelmapa : mapColidh = merge(mapColidh, clusters, by = '
NAME1', all.x = F)
lo tengo? names(mapColidh)
a pintar:
library(RColorBrewer) library(classInt)
variable a colorear varToPlot=mapColidh$cluster
decidir color: unique(varToPlot) aggregate(mapColidh@data[, c(10, 11, 12)], by =
list(mapColidh@data$cluster), FUN=mean)
```

Bibliograf<U+00ED>a

```

■plotMap1, echo=FALSE, fig=TRUE■=
  preparo colores numberOfClasses = length(unique(varToPLot))
colorForScale='Set2' paleta = brewer.pal(numberOfClasses, colorForScale)
  grafico mapa basico plot(mapCol,col='grey',border=0)
  grafico mapa cluster plot(mapCol,dh,col = paleta[varToPLot],border =
F,add = T)legend('left',legend = c("LOW","UP","MEDIUM"),fill =
paleta,cex = 0.6,bty = "n",title = "conglomerado")

```

Figure 3: Mapa del IDH