

Los Conteos Rápidos del INE y la Topología del Símplex Unitario p -dimensional

Presenta:

M. en C. Arnoldo Daniel Miranda Fournier

Asesor:

Dr. Gabriel Nuñez Antonio

Laboratorio de Ciencia de Datos

Una travesía por la Ciencia de Datos para matemáticos y no matemáticos

17 de Septiembre de 2020

Contenido

- 1 Datos Composicionales
- 2 Modelo Propuesto
- 3 Aplicaciones
- 4 Conclusiones

Objetivo General

El objetivo de la tesis fue proponer una metodología para contribuir al análisis, descripción y modelación de datos composicionales desde un enfoque bayesiano.

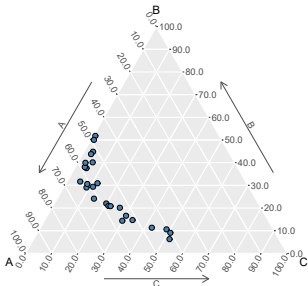
Introducción

¿Qué son los Datos Composicionales?

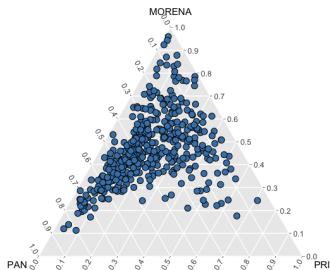
- Composición Geoquímica de Rocas
- Composición Geoquímica de Sedimentos
- Composición Proteica de la Leche
- Procesos Electorales

Ejemplos

- Composición Geoquímica de Rocas



- Datos Electorales (Puebla 2019)



- Geología, Política, Economía, Medicina, etc.

Su naturaleza y problemática

El simplex como espacio muestral

$$S^p = \{ \mathbf{x} \in \mathbb{R}^p : x_i \geq 0, \sum_{i=1}^p x_i = 1 \}$$

$$S^2 \subset \mathbb{R}^2$$

Segmento



$$S^3 \subset \mathbb{R}^3$$

Diagrama Ternario



$$S^4 \subset \mathbb{R}^4$$

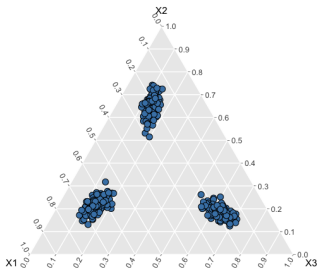
Tetraedro



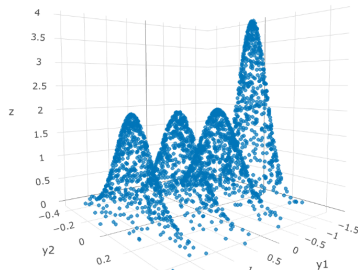
- Aitchison (1982) estudia las características y operaciones de los datos composicionales.

Objetivo específico

- Patrones más complejos

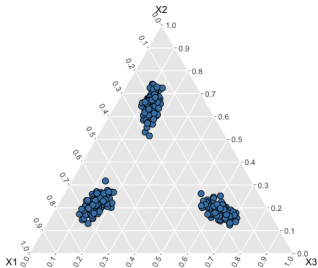


- Mezcla de Distribuciones Normales Multivariadas

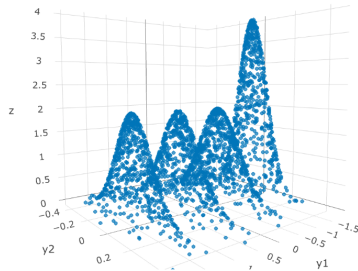


Objetivo específico

- Patrones más complejos



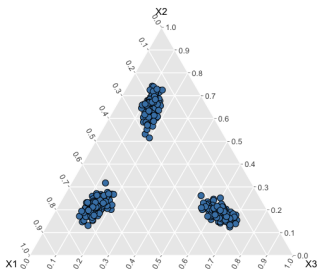
- Mezcla de Distribuciones Normales Multivariadas



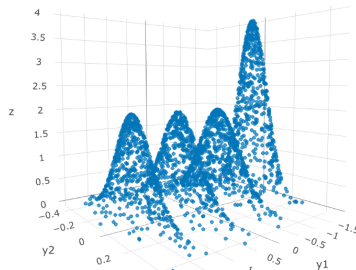
- Mediante la conocida transformación log-cociente *ilr*, definida por Egozcue et al. (2003).

Objetivo específico

- Patrones más complejos



- Mezcla de Distribuciones Normales Multivariadas



- Mediante la conocida transformación log-cociente *ilr*, definida por Egozcue et al. (2003).
- Los modelos de mezclas ofrecen una gran versatilidad en la descripción de datos, su flexibilidad permite aproximar problemas desde un enfoque no paramétrico, son ampliamente utilizados en Estadística.

Un modelo no-paramétrico

Sea $\mathbf{y} \in \mathbb{S}^{p+1}$. Entonces mediante la transformación ilr :

$$ilr(\mathbf{y}) = \mathbf{x}, \quad \mathbf{x} \in \mathbb{R}^p$$

Un modelo no-paramétrico

Sea $\mathbf{y} \in \mathbb{S}^{p+1}$. Entonces mediante la transformación ilr :

$$ilr(\mathbf{y}) = \mathbf{x}, \quad \mathbf{x} \in \mathbb{R}^p$$

El modelo no-paramétrico propuesto para modelar a \mathbf{x} es:

$$\mathbf{x} \sim f(\mathbf{x} | \boldsymbol{\omega}, \boldsymbol{\mu}, \boldsymbol{\lambda})$$

Donde

$$f(\mathbf{x} | \boldsymbol{\omega}, \boldsymbol{\mu}, \boldsymbol{\lambda}) = \sum_{i=1}^{\infty} \omega_i N_p(\mathbf{x} | \boldsymbol{\mu}_i, \boldsymbol{\lambda}_i) \quad \text{donde} \quad \sum_{i=1}^{\infty} \omega_i = 1, \quad (1)$$

Un modelo no-paramétrico

Sea $\mathbf{y} \in \mathbb{S}^{p+1}$. Entonces mediante la transformación *ilr*:

$$\text{ilr}(\mathbf{y}) = \mathbf{x}, \quad \mathbf{x} \in \mathbb{R}^p$$

El modelo no-paramétrico propuesto para modelar a \mathbf{x} es:

$$\mathbf{x} \sim f(\mathbf{x} | \boldsymbol{\omega}, \boldsymbol{\mu}, \boldsymbol{\lambda})$$

Donde

$$f(\mathbf{x} | \boldsymbol{\omega}, \boldsymbol{\mu}, \boldsymbol{\lambda}) = \sum_{i=1}^{\infty} \omega_i N_p(\mathbf{x} | \boldsymbol{\mu}_i, \boldsymbol{\lambda}_i) \quad \text{donde} \quad \sum_{i=1}^{\infty} \omega_i = 1, \quad (1)$$

Con N_p una densidad *Normal Multivariada* con media $\boldsymbol{\mu}_i \in \mathbb{R}^p$ y matriz de precisión $\boldsymbol{\lambda}_i$ simétrica y positiva definida, Bernardo and Smith (2009):

$$N_p(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\lambda}) = |\boldsymbol{\lambda}|^{\frac{1}{2}} (2\pi)^{-\frac{p}{2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^t \boldsymbol{\lambda} (\mathbf{x} - \boldsymbol{\mu}) \right\}$$

Un modelo no-paramétrico

Sea $\mathbf{y} \in \mathbb{S}^{p+1}$. Entonces mediante la transformación *ilr*:

$$\text{ilr}(\mathbf{y}) = \mathbf{x}, \quad \mathbf{x} \in \mathbb{R}^p$$

El modelo no-paramétrico propuesto para modelar a \mathbf{x} es:

$$\mathbf{x} \sim f(\mathbf{x} | \boldsymbol{\omega}, \boldsymbol{\mu}, \boldsymbol{\lambda})$$

Donde

$$f(\mathbf{x} | \boldsymbol{\omega}, \boldsymbol{\mu}, \boldsymbol{\lambda}) = \sum_{i=1}^{\infty} \omega_i N_p(\mathbf{x} | \boldsymbol{\mu}_i, \boldsymbol{\lambda}_i) \quad \text{donde} \quad \sum_{i=1}^{\infty} \omega_i = 1, \quad (1)$$

Con N_p una densidad *Normal Multivariada* con media $\boldsymbol{\mu}_i \in \mathbb{R}^p$ y matriz de precisión $\boldsymbol{\lambda}_i$ simétrica y positiva definida, Bernardo and Smith (2009):

$$N_p(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\lambda}) = |\boldsymbol{\lambda}|^{\frac{1}{2}} (2\pi)^{-\frac{p}{2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^t \boldsymbol{\lambda} (\mathbf{x} - \boldsymbol{\mu}) \right\}$$

Los parámetros del modelo son las colecciones infinitas de $\{\boldsymbol{\omega}, \boldsymbol{\mu}, \boldsymbol{\lambda}\}$.

El enfoque Bayesiano de la Estadística

- $X \sim f(x | \theta)$: Modelo de Probabilidad.
- $\pi(\theta)$: Distribución Inicial.
- Muestra aleatoria $\mathbf{x} = \{x_1, x_2, \dots, x_n\}$ del modelo $f(x | \theta)$
Verosimilitud: $f(\mathbf{x} | \theta) = f(x_1, x_2, \dots, x_n | \theta) = \prod_{i=1}^n f(x_i | \theta)$
- Teorema de Bayes...Distribución Final:

$$\begin{aligned}\pi(\theta | \mathbf{x}) &= \frac{f(\mathbf{x} | \theta)\pi(\theta)}{f(\mathbf{x})} \\ &= \frac{f(\mathbf{x} | \theta)\pi(\theta)}{\int_{\Theta} f(\mathbf{x} | \theta)\pi(\theta)d\theta} \\ &= \frac{\prod_{i=1}^n f(x_i | \theta)\pi(\theta)}{\int_{\Theta} \prod_{i=1}^n f(x_i | \theta)\pi(\theta)d\theta}\end{aligned}$$

- El enfoque Bayesiano permite transformar la estructura compleja de un modelo de mezclas en un conjunto de estructuras simples.

Metodología

- Inferencia desde el enfoque Bayesiano:
 - Procesos de Dirichlet.
 - Modelo de Mezclas de Procesos Dirichlet.

Metodología

- Inferencia desde el enfoque Bayesiano:
 - Procesos de Dirichlet.
 - Modelo de Mezclas de Procesos Dirichlet.
- Se utilizó el algoritmo dado por Walker (2007):
 - Gibbs Sampling.
 - Slice Sampling.
 - Variables latentes.
- Se complemento con la propuesta realizada por Kalli et al. (2011):
 - Slice Sampling más general.

Metodología

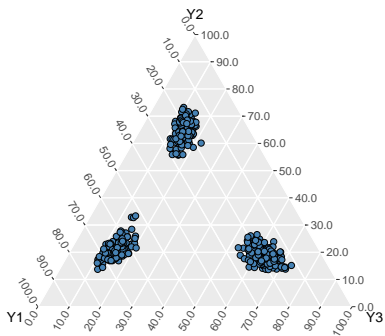
- Inferencia desde el enfoque Bayesiano:
 - Procesos de Dirichlet.
 - Modelo de Mezclas de Procesos Dirichlet.
- Se utilizó el algoritmo dado por Walker (2007):
 - Gibbs Sampling.
 - Slice Sampling.
 - Variables latentes.
- Se complementó con la propuesta realizada por Kalli et al. (2011):
 - Slice Sampling más general.
- La combinación y adaptación de ambos algoritmos se implementaron en el lenguaje estadístico **R**.

Metodología

- Inferencia desde el enfoque Bayesiano:
 - Procesos de Dirichlet.
 - Modelo de Mezclas de Procesos Dirichlet.
- Se utilizó el algoritmo dado por Walker (2007):
 - Gibbs Sampling.
 - Slice Sampling.
 - Variables latentes.
- Se complementó con la propuesta realizada por Kalli et al. (2011):
 - Slice Sampling más general.
- La combinación y adaptación de ambos algoritmos se implementaron en el lenguaje estadístico **R**.
- El algoritmo regresa una muestra de la distribución predictiva final en \mathbb{R}^P que es mapeada de regreso al Simplex a través de la transformación ilr^{-1} .

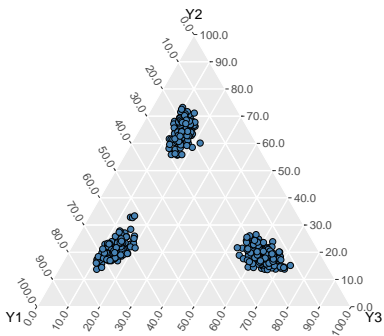
Datos Simulados

Datos Simulados 1



Datos Simulados

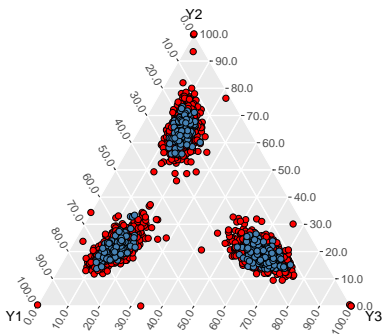
Datos Simulados 1



Intervalos Poblacionales Marginales		
Componente	Inferior	Superior
Y1	14.68 %	70.38 %
Y2	14.73 %	70.39 %
Y3	10.51 %	68.41 %

Datos Simulados

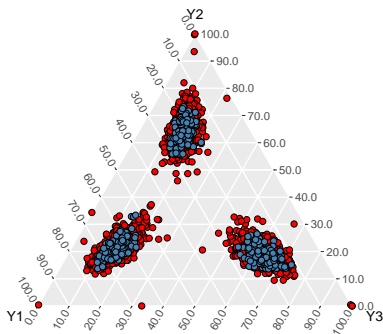
Datos Simulados 1

**Intervalos Poblacionales Marginales**

Componente	Inferior	Superior
Y1	14.68 %	70.38 %
Y2	14.73 %	70.39 %
Y3	10.51 %	68.41 %

Datos Simulados

Datos Simulados 1



Intervalos Poblacionales Marginales

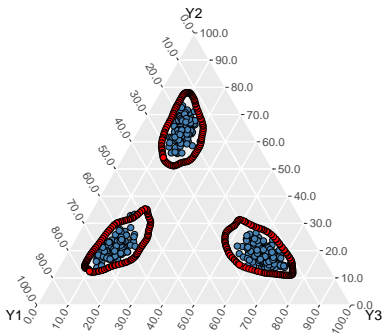
Componente	Inferior	Superior
Y1	14.68 %	70.38 %
Y2	14.73 %	70.39 %
Y3	10.51 %	68.41 %

Intervalos de Probabilidad Marginales

Componente	Inferior	Superior
Y1	13.42 %	72.70 %
Y2	13.71 %	72.45 %
Y3	09.34 %	69.74 %

Datos Simulados

Datos Simulados 1

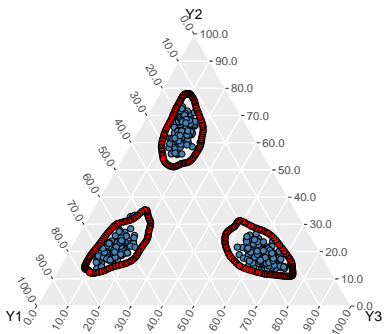


Intervalos Poblacionales Marginales		
Componente	Inferior	Superior
Y1	14.68 %	70.38 %
Y2	14.73 %	70.39 %
Y3	10.51 %	68.41 %

Intervalos de Probabilidad Marginales		
Componente	Inferior	Superior
Y1	13.42 %	72.70 %
Y2	13.71 %	72.45 %
Y3	09.34 %	69.74 %

Datos Simulados

Datos Simulados 1



Intervalos Poblacionales Marginales

Componente	Inferior	Superior
Y1	14.68 %	70.38 %
Y2	14.73 %	70.39 %
Y3	10.51 %	68.41 %

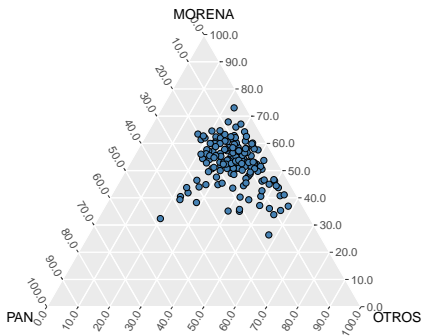
Intervalos de Probabilidad Marginales

Componente	Inferior	Superior
Y1	13.42 %	72.70 %
Y2	13.71 %	72.45 %
Y3	09.34 %	69.74 %

- Se propuso establecer una especificación inicial no informativa en el Símplex.

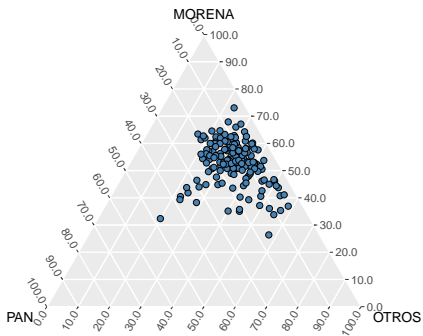
Datos Reales

Gubernatura del estado de Morelos 2018



Datos Reales

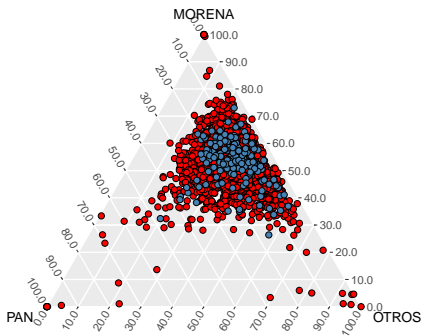
Gubernatura del estado de Morelos 2018



Cómputo Distrital		
Coalición		
MORENA	PAN	OTROS
52.59 %	14.05 %	33.36 %

Datos Reales

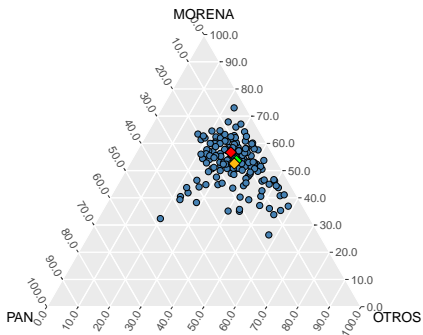
Gubernatura del estado de Morelos 2018



Cómputo Distrital		
Coalición		
MORENA	PAN	OTROS
52.59 %	14.05 %	33.36 %

Datos Reales

Gubernatura del estado de Morelos 2018

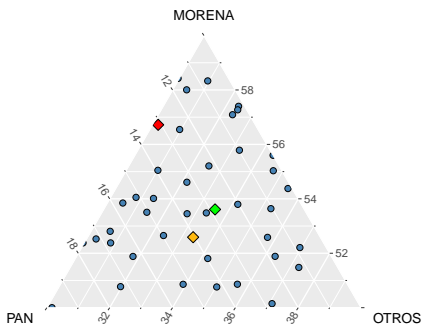


Cómputo Distrital ♦		
Coalición		
MORENA	PAN	OTROS
52.59 %	14.05 %	33.36 %

Estimadores Puntuales			
Estimador	Coalición		
	MORENA	PAN	OTROS
Media ♦	56.72 %	13.10 %	30.18 %
Mediana ♦	53.61 %	12.84 %	33.55 %

Datos Reales

Gubernatura del estado de Morelos 2018

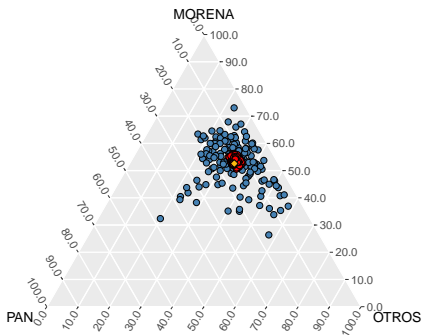


Cómputo Distrital ◆		
Coalición		
MORENA	PAN	OTROS
52.59 %	14.05 %	33.36 %

Estimadores Puntuales			
Estimador	Coalición		
	MORENA	PAN	OTROS
Media ◆	56.72 %	13.10 %	30.18 %
Mediana ◆	53.61 %	12.84 %	33.55 %

Datos Reales

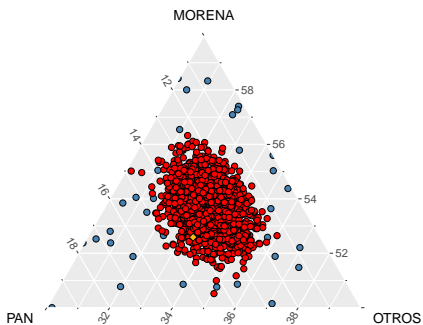
Gubernatura del estado de Morelos 2018



Cómputo Distrital ♦		
Coalición		
MORENA	PAN	OTROS
52.59 %	14.05 %	33.36 %

Datos Reales

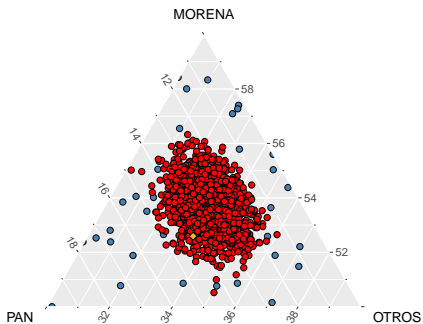
Gubernatura del estado de Morelos 2018



Cómputo Distrital ♦		
Coalición		
MORENA	PAN	OTROS
52.59 %	14.05 %	33.36 %

Datos Reales

Gubernatura del estado de Morelos 2018

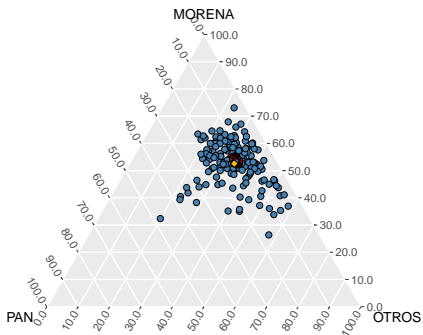


Cómputo Distrital ♦		
Coalición		
MORENA	PAN	OTROS
52.59 %	14.05 %	33.36 %

Intervalos de Probabilidad Marginales		
Coalición	Inferior	Superior
MORENA	52.20 %	55.31 %
PAN	11.80 %	14.22 %
OTROS	31.70 %	34.96 %

Datos Reales

Gubernatura del estado de Morelos 2018

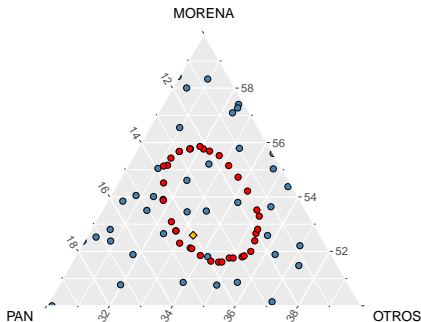


Cómputo Distrital ♦		
Coalición		
MORENA	PAN	OTROS
52.59 %	14.05 %	33.36 %

Intervalos de Probabilidad Marginales		
Coalición	Inferior	Superior
MORENA	52.20 %	55.31 %
PAN	11.80 %	14.22 %
OTROS	31.70 %	34.96 %

Datos Reales

Gubernatura del estado de Morelos 2018



Cómputo Distrital ◆

Coalición		
MORENA	PAN	OTROS
52.59 %	14.05 %	33.36 %

Intervalos de Probabilidad Marginales

Coalición	Inferior	Superior	Tamaño
MORENA	52.20 %	55.31 %	3.11 %
PAN	11.80 %	14.22 %	2.43 %
OTROS	31.70 %	34.96 %	3.26 %

Intervalos del INE

Coalición	Inferior	Superior	Tamaño
MORENA	51.00 %	53.80 %	2.8 %
PAN	13.40 %	16.10 %	2.7 %

Comparando con los resultados del INE nuestros intervalos exceden en 0.31 % a la componente MORENA y disminuye en 0.27 % a la componente PAN.

Comparación entre Modelos

Cómputo Distrital		
Coalición		
MORENA	PAN	OTROS
52.59 %	14.05 %	33.36 %

Modelo de Mezclas de Normales Medianas

Intervalos de Probabilidad Marginales			
Coalición	Inferior	Superior	Tamaño
MORENA	52.20 %	55.31 %	3.11 %
PAN	11.80 %	14.22 %	2.43 %
OTROS	31.70 %	34.96 %	3.26 %

Modelo de Mezclas de Gammas Medias

Intervalos de Probabilidad Marginales			
Coalición	Inferior	Superior	Tamaño
MORENA	51.99 %	54.57 %	2.58 %
PAN	12.25 %	14.46 %	2.21 %
OTROS	31.88 %	34.91 %	3.03 %

Intervalos del INE			
Coalición	Inferior	Superior	Tamaño
MORENA	51.00 %	53.80 %	2.8 %
PAN	13.40 %	16.10 %	2.7 %

Notas: Comparando con los resultados del INE el Modelo de Mezclas de Normales con la Mediana excede en 0.31 % a la componente MORENA y disminuye en 0.27 % a la componente PAN.

Comparando con los resultados del INE el Modelo de Mezclas de Gammas con la Media disminuye en 0.22 % a la componente MORENA y disminuye en 0.49 % a la componente PAN.

Por lo tanto, para el estado de Morelos aparentemente es más adecuado el Modelo de Mezclas de Gammas.

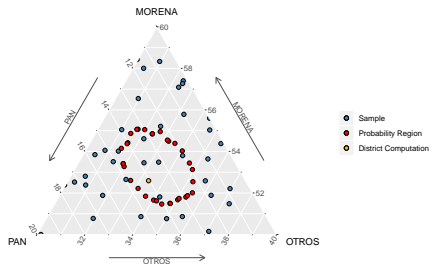
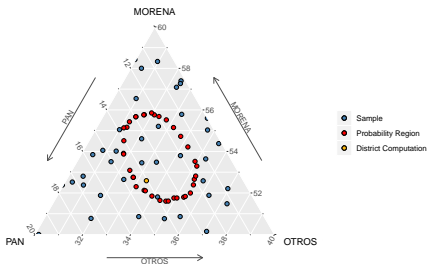
Comparación entre Modelos

Modelo de Mezclas de Normales Medianas

Modelo de Mezclas de Gammas Medias

Gubernatura Morelos 2018

Gubernatura Morelos 2018



Temas abarcados

- Datos Composicionales
- Estadística Bayesiana
- Modelos de mezclas
- Procesos Dirichlet
- Métodos de simulación
- Programación

Conclusiones

- Se propone un modelo que toma en cuenta la naturaleza de los datos (datos composicionales).
- Con esta metodología se esta en condiciones de ofrecer estimadores puntuales que cumplan con la restricción de que la suma de sus componentes sea uno.
- De igual manera se pueden ofrecer intervalos de probabilidad marginales, para cada componente involucrada en el vector composicional.
- La variabilidad de la proporción se mide al correr, cada vez el programa completo, obteniendo la mediana.

Conclusiones

- El modelo no-parámtrico (basado en mezclas infinitas de distribuciones normales multivariadas) fue viable gracias al uso de la transformación log-cociente *ilr*.
- Se podrían considerar modelos de mezclas con distribuciones definidas en ortantes positivos, con el propósito de omitir cualquier tipo de transformación log-cociente.
- El modelo puede mejorarse optimizando los algoritmos utilizados para la inferencia, por ejemplo, se podrían paralelizar algunos procesos.
- Aunque quedan cosas por hacer, esta propuesta de tesis sienta las bases para el análisis bayesiano de datos composicionales.

Referencias

- Aitchison, J. (1982). The statistical analysis of compositional data. *Journal of the Royal Statistical Society: Series B (Methodological)*, 44(2):139–160.
- Bernardo, J. M. and Smith, A. F. (2009). *Bayesian theory*, volume 405. John Wiley & Sons.
- Egozcue, J. J., Pawlowsky-Glahn, V., Mateu-Figueras, G., and Barcelo-Vidal, C. (2003). Isometric logratio transformations for compositional data analysis. *Mathematical Geology*, 35(3):279–300.
- Kalli, M., Griffin, J. E., and Walker, S. G. (2011). Slice sampling mixture models. *Statistics and computing*, 21(1):93–105.
- Walker, S. G. (2007). Sampling the dirichlet mixture model with slices. *Communications in Statistics—Simulation and Computation*®, 36(1):45–54.

Mi correo:
danielfourrier270790@gmail.com

¡Gracias!