

EST-46114 Métodos Multivariados y Datos Categóricos

Sesion 03 - Inferencia en la Distribucion Gaussiana Multivariada - Parte 2

Juan Carlos Martinez-Ovando

Maestria en Ciencia de Datos



Objetivos

- Estudiaremos soluciones de problemas inferenciales comunes asociados con la clase de modelos gaussianos multidimensionales.

Prreguntas inferenciales II

Preguntas para datos tipo swiss, e.g.:

P1. ¿Es el tasa Agriculture al menos **9 puntos porcentuales** mayor a la de Catholic?

P2. ¿Es la tasa de Agriculture al menos **1.25 veces mayor** que la de Catholic?

P3. Otras preguntas relacionadas con la media y dispersion. . . .

El diccionario de datos swiss esta disponible en <https://stat.ethz.ch/R-manual/R-devel/library/datasets/html/swiss.html>.

P1. Agriculture vs Catholic

La pregunta en cuestion puede plantearse de dos formas:

P1.a. Directamente sobre los valores de la tasa, i.e.

$$X_{Agriculture} > X_{Catholic} + 9.$$

referido a comunidades suizas, en general.

P1.b. Referidos a los valores de la tasa de una comunidad promedio, i.e.

$$\mu_{Agriculture} > \mu_{Catholic} + 9.$$

P1.b. Solucion

- La solucion para P1.b. es la mas sencilla, pues se refiere a los parametros del modelo y no a las variables aleatorias.
- Siendo referido a los parametros, puede a su vez resolverse de manera lineal o no lineal (refiriendo por linealidad a la propiedad y operacion sobre los parametros).

P1.b. Solucion lineal

La solucion lineal puede plantearse como un problema de decision bajo incertidumbre, considerando dos posibles hipotesis:

H_0 : Hipotesis 0

$$\mu_{Agriculture} - \mu_{Catholic} > 9.$$

H_1 : Hipotesis 1

$$\mu_{Agriculture} - \mu_{Catholic} \leq 9.$$

P1.b. Solucion lineal

Elementos del problema de decision:

1. Espacio de decisiones

$$\mathcal{A} = \mathbb{R} \times \mathbb{R} = \{(\mu_{Agriculture}, \mu_{Catholic}) : \mu_{Agriculture}, \mu_{Catholic} \text{ son reales}\}.$$

2. Espacio de incertidumbre:

$$\Theta = \mathbb{R} \times \mathbb{R} = \{(\mu_{Agriculture}, \mu_{Catholic}) : \mu_{Agriculture}, \mu_{Catholic} \text{ son reales}\}.$$

3. Funcion de utilidad: $U : \mathcal{A} \times \Theta \rightarrow \{0, 1\}$, siendo $U = 1$ cuando al eleccion/decision dle parametro coincide con el verdadero valor del paraemto, con $U = 0$ en el caso contrario.
4. Cuantificacion de incertidumbre: $\mathbb{P}(\mu_{Agriculture}, \mu_{Catholic} | \text{datos, complemento})$, probabilidad adecuada a la informacion de los datos.

Noten que el resto de los parametros no esta considerdo en las hipotesis, por lo que no es referido en los elementos del problema de decision.

P1.b. Solucion lineal

Noten que las hipotesis H_0 y H_1 inducen una particion sobre $\Theta = \Theta_0 \cup \Theta_1$, con

$$\Theta_0 = \{(\mu_{Agriculture}, \mu_{Catholic}) : \mu_{Agriculture} - \mu_{Catholic} > 0\},$$

y

$$\Theta_1 = \Theta \setminus \Theta_0.$$

Regla de decision:

La regla de decision consistira en elegir la hipotesis H con mayor probabilidad, i.e.

$$H^* = \arg \max_{H \in \{H_0, H_1\}} \mathbb{P}(H|\text{datos}),$$

siendo

$$\mathbb{P}(H_j|\text{datos}) = \mathbb{P}((\mu_{Agriculture}, \mu_{Catholic}) \in \Theta_j|\text{datos}).$$

Aprendizaje bayesiano

Retomemos de la sesion anterior los resultados del aprendizaje bayesiano de la clase de modelos gaussianos con los datos `swiss` usando la funcion `gaussian.posterior`:

El resultado es un objeto lista, `output`, con cuatro objetos

- `output[[1]]` - parametro μ_n
- `output[[2]]` - parametro s_n
- `output[[3]]` - parametro a_n
- `output[[4]]` - parametro \mathbf{B}_n

Aprendizaje bayesiano

```
output[{1}]
```

```
## [[1]]  
##           [,1]  
## Fertility    68.68125  
## Agriculture  49.60417  
## Examination  16.14583  
## Education    10.75000  
## Catholic     40.28667  
## Infant.Mortality 19.52708
```

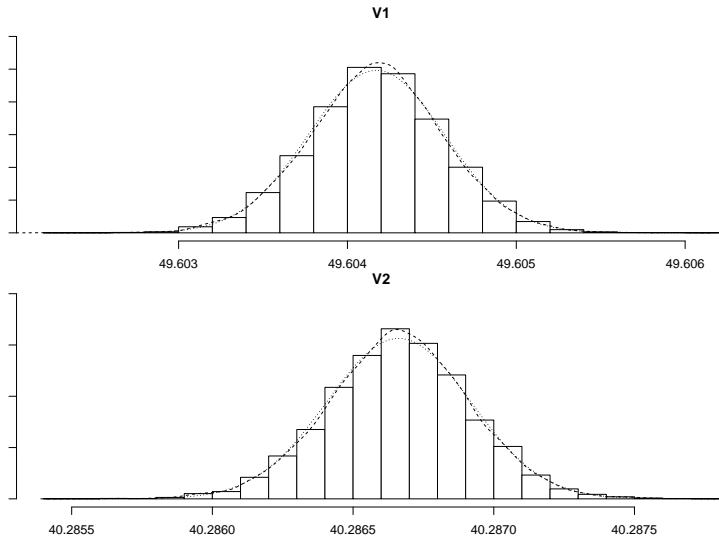
En primera instancia, pareciera ser que H_0 se cumple con estos datos. Sin embargo, $\mu_{Agriculture,n}$ y $\mu_{Catholic,n}$ son estimadores de cantidades desconocidas, por lo que la relación referida a H_0 dependerá de la variabilidad epistémica de estos parámetros, en función de a_n , s_n y B_n .

Incertidumbre epistémica

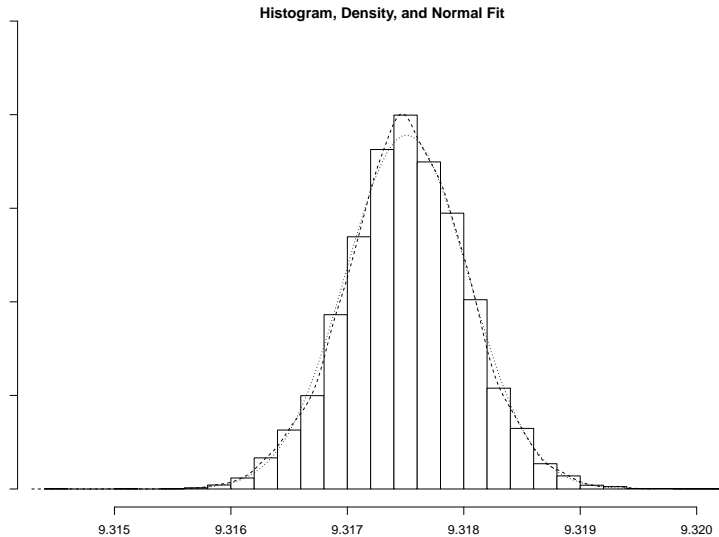
La **incertidumbre epistémica** se refiere a la cuantificación del desconocimiento acerca de $(\mu_{Agriculture,n}, \mu_{Catholic,n})$, en este caso. tal desconocimiento, marginal, esta referido a:

- Distribuciones marginales de $\mu_{Agriculture}$ y $\mu_{Catholic}$
- Interaccion / comovimiento / covarianza entre $\mu_{Agriculture}$ y $\mu_{Catholic}$

Visualizacion I



Visualizacion II



Regla de decision (solucion lineal)

La par probabilidades para H_0 y H_1 (aproximadas por Monte Carlo) son

```
prob_lin_h0 <- length(which(mu.sim.2[,1]-mu.sim.2[,2] > 9))/M  
prob_lin_h0
```

```
## [1] 1
```

y

```
prob_lin_h1 <- length(which(mu.sim.2[,1]-mu.sim.2[,2] <= 9))/M  
prob_lin_h1
```

```
## [1] 0
```

En este caso, la evidencia es contundente. **Consideren que este NO es el caso, en general.**

P1.b. Solucion lineal exacta

La solucion del problema, que hemos resuelto via simulacion se pudo haber obtenido analiticamente empleando el resultado de cerradura bajo linealidad, considerando la aplicacion del vector

$$\mathbf{c}_{Hipotesis} = (0, 1, 0, 0, -1, 0),$$

a

$$\boldsymbol{\mu},$$

considerando que

$$(\boldsymbol{\mu}, \boldsymbol{\Lambda}) \sim \text{N-Wi}(\boldsymbol{\mu}, \boldsymbol{\Lambda} | \mathbf{m}_n, \mathbf{s}_n, a_n, \mathbf{B}_n).$$

- Se sigue que $\mathbf{c}'_{Hipotesis} \boldsymbol{\mu}$ sigue una distribucion t -univariada.
- Los datos que simulamos, son una muestra 'iid' de esta distribucion marginal.
- Solo problemas asociados con transformaciones afines de \mathbf{X} o $\boldsymbol{\mu}$ pueden resolverse de manera analitica cerrada. En general, sera util recurrir a herramientas de simulacion estocastica.

P2. Solucion no lineal

Ahora, la solucion no lineal a la pregunta P2 puede plantearse como un problema de decision bajo incertidumbre, considerando dos posibles hipotesis:

H_0 : Hipotesis 0

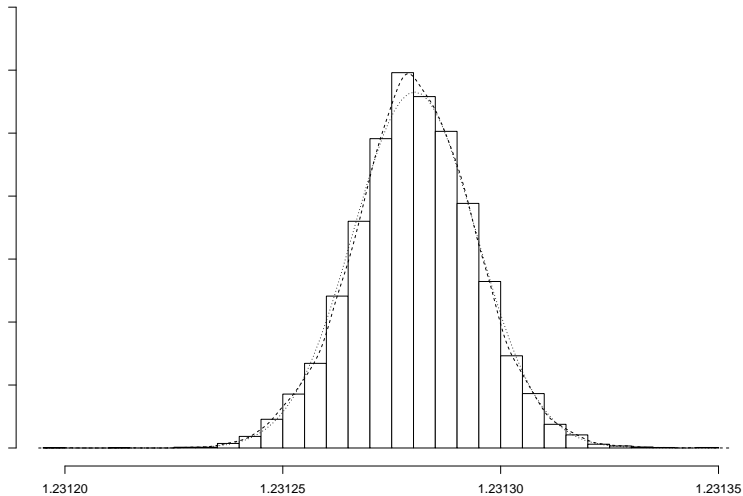
$$\mu_{Agriculture} / \mu_{Catholic} > 1.25.$$

H_1 : Hipotesis 1

$$\mu_{Agriculture} / \mu_{Catholic} \leq 1.25.$$

Visualizacion III

Histogram, Density, and Normal Fit



Regla de decision (solucion lineal)

La par probabilidades para H_0 y H_1 (aproximadas por Monte Carlo) son

```
prob_nolin_h0 <- length(which(mu.sim.2[,1]/mu.sim.2[,2] > 1.25))/M  
prob_nolin_h0
```

```
## [1] 0
```

y

```
prob_nolin_h1 <- length(which(mu.sim.2[,1]/mu.sim.2[,2] <= 1.25))/M  
prob_nolin_h1
```

```
## [1] 1
```

En este caso, la evidencia es contundente. **Consideren que este NO es el caso, en general.**

Ejercicio

Desarrollen las soluciones al problema `lineal` y no `lineal` deferido ahora a las variables

$$X_{Agriculture} \text{ y } X_{Catholic},$$

en lugar de referirlo a $\mu_{Agriculture}$ y $\mu_{Catholic}$.

Desarrollen las soluciones via simulacion, considerando lo siguiente:

$$\begin{aligned}\Lambda | \text{datos} &\sim \text{Wi}(\Lambda | \mathbf{a}_n, \mathbf{B}_n) \\ \mu | \Lambda, \text{datos} &\sim \text{N}(\mu | \mathbf{m}_n, s_n \Lambda) \\ \mathbf{X} | \mu, \Lambda, \text{datos} &\sim \text{N}(\mathbf{x} | \mu, \Lambda).\end{aligned}$$

En la siguiente sesion

Revisaremos como hacer:

1. Descomposicion singular de matrices de covarianzas
2. Analisis de componentes principales
3. Analisis inferencial de componentes principales

Lecturas complementarias

- Hothorn et al (2018) "On multivariate t and Gauss probabilities in R".
est46114_s04_suplemento1.pdf
- Martinez-Ovando (2016) "Paradigma Bayesiano de Inferencia (Resumen de Teoria)", *Notas de clase EST-46114*. (De momento, no presten atencion a la descripcion de los metodos de simulacion transdimensionales.)
est46114_s03_suplemento2.pdf
- Press (2005). "Applied multivariate analysis, using Bayesian and frequentist methods of inference." Dover Pub. (Capitulo 4.)

Table of Contents