

# La tosferina en menores de edad en la ciudad de Medellín

Mateo Murcia Valles      Christopher Andrés Obando  
Juan David García Zapata

Junio 2023

---

## Abstract

This study aims to use information from a database on the open data page of the city of Medellín, related to pertussis disease, to determine if hospitalization is required and important characteristics such as age in months, gender, security regime, and the presence of vomiting. For this purpose, several logistic models are proposed that use Bayesian statistics and the a priori information provided by this data in order to answer questions such as: what type of social security system can influence a person to be hospitalized. Additionally, attempts are made to predict possible cases that may occur based on gender, age in months, stage of the disease, and the presence of vomiting. Furthermore, we present the processes used, the results obtained, and the conclusions.

---

**Keywords:** *your keywords here*

**MSC (2020):** *your keywords here*

## 1 Introducción

La tosferina, o pertussis, es una enfermedad infecciosa causada por la bacteria *Bordetella pertussis*. Se caracteriza por ataques intensos de tos, con sonidos agudos al inhalar, que pueden durar varias semanas. Es altamente contagiosa y puede ser especialmente peligrosa en bebés y niños pequeños. La vacunación es fundamental para prevenir la tosferina y sus complicaciones. Con esta inducción se entra en la materia a la problemática que se abordará aplicando la estadística bayesiana. En este estudio se pretende realizar un análisis acerca de una base de datos extraída de la página de datos abiertos de la alcaldía de Medellín (MEDdata), la cual con la información de paciente que fueron diagnosticados con tosferina en el área metropolitana de Medellín, dicha base de datos, cuenta con 3968 registros y 25 variables desde la primera semana del 2008 hasta la última semana del 2021.

## 2 Descripción de los datos

Tal y como se mencionó anteriormente, la base de datos que se va a utilizar constaba inicialmente con 3968 datos y 25 variables, las cuales son:

- **ID:** Número consecutivo. Identificador único del registro llevado a cabo.
- **Semana:** Semanas del año de 1 a 53. Representa la semana del año en la que se llevó a cabo el registro.
- **edad:** Edad del paciente que sufre de tosferina.
- **uni med:** Unidad de medida: 0 = No aplica, 1 = Años, 2 = Meses, 3 = Días, 4 = Horas, 5 = Minutos SD = Sin información. Corresponde a la unidad medida en la que fue tomada la edad del paciente. Cabe notar que en su mayoría los datos fueron registrados en años y meses.
- **sexo:** M = Masculino, F = Femenino, SD = Sin información. Representa el sexo del paciente.
- **nombre barrio:** Texto asociado a la tabla de barrios definidos por la entidad territorial, vacíos se diligencian con "Sin información", Sin ubicación en zona urbana. Representa el nombre del barrio al cual pertenece el paciente.
- **comuna:** Texto asociado a la tabla de barrios definidos por la entidad territorial, vacíos se diligencian con "Sin información", Sin ubicación en zona urbana. Representa el nombre de la comuna a la cual pertenece el paciente.
- **Tip\_ss.:** Tipo de Régimen de seguridad social C = Contributivo, S = Subsidiado, P = Excepción, E = Especial, N = No asegurado, I = Indeterminado/Pendiente, SD = Sin información. Corresponden a los diferentes tipos de régimen de seguridad social que existen en Colombia.
- **cod\_ase:** Código de la aseguradora. Corresponde al código de la aseguradora a la cual está asociado el paciente.
- **fec\_con:** Fecha de Consulta. Corresponde a la fecha en la que el paciente realizó la consulta.
- **Ini\_sin:** Fecha de inicio de síntomas. Describe la fecha inicial en la que el paciente reportó síntomas.
- **Tip\_cas:** Tipo de caso: 1 = Sospechoso, 2 = Probable, 3 = Confirmado por laboratorio, 4 = Confirmado por clínica, 5 = Confirmado por nexo epidemiológico. Corresponde al estado del paciente frente a la valoración por parte de las IPS (Instituciones Prestadoras de Servicio en Salud) en Medellín, categoriza a los pacientes de tosferina. Cabe notar que en nuestro estudio se basó principalmente en los casos confirmados.
- **Pac\_hos:** Paciente Hospitalizado: 1 = Sí, 2 = No. Corresponde a una variable dicotómica que representa si el paciente fue hospitalizado o no. Representa en el estudio la variable respuesta que busca ser explicada a través de las variables disponibles y sobre la cual se realizarán los análisis y conclusiones al respecto.

- **Ira:** 1 = Sí, 2 = No, SD = Sin Información.
- **eta\_enf:** 1 = Catarral, 2 = Espasmódica, 3 = Convaleciente, SD = Sin Información.
- **tos:** 1 = Sí, 2 = No, SD = Sin Información.
- **dur\_tos:** Depende del campo anterior: si tos = 1 (Sí), SD = Sin Información.
- **tos\_par:** 1 = Sí, 2 = No, SD = Sin información.
- **estridor:** 1 = Sí, 2 = No, SD = Sin información.
- **apnea:** 1 = Sí, 2 = No, SD = Sin información.
- **cianosis:** 1 = Sí, 2 = No, SD = Sin información.
- **vomito:** 1 = Sí, 2 = No, SD = Sin información.
- **complicaci:** 1 = Sí, 2 = No, SD = Sin información.
- **tip\_com:** 1 = Convulsiones, 2 = Atelectasia, 3 = Neumotórax, 4 = Neumonía, 5 = Otro, SD = Sin información.
- **evento:** Texto asociado a los eventos notificados.
- **year\_:** Año en el que el paciente fue atendido.

## 2.1 Depuración

Principalmente en la depuración de los datos se tomaron solo los sujetos que se confirmaron que tenia la tosferina, después se procedió a unificar la unidad de medida de las fechas se dejaron todas en meses, luego se eliminaron variables que solo tenían en su mayoría NA's y variables tipo identificadoras, luego se eliminó filas también de NA'S. después se procedió a seleccionar las variables a utilizar como son: **edad**, **sexo\_**, **tipo\_ss\_** (tipo de seguridad social), **eta\_enf\_** (etapa de la enfermedad), **vomito**, **pac\_hos** (paciente hospitalizado).

Seguidamente, se procedió a filtrar la base de datos a solo menores de edad, en tipo de seguridad social no se tuvo en cuenta una categorica por su falta de registros, se organiza la variable **pac\_hos\_** para sea dicótoma, posteriormente las variables categóricas quedarón así:

- **Sexo:**
  - **M** = Masculino
  - **F** = Femenino
- **Pac\_hos:**
  - **0** = No
  - **1** = Si
- **Tipo\_ss:**
  - **C** = Contributivo
  - **E** = Especial
  - **N** = No asegurado
  - **P** = Excepción
  - **S** = Subsidiado
- **Eta\_enf:**
  - **1** = Catarral

- **2** = Espasmódica
- **3** = Convaleciente
- **Vomito:**
  - **1** = Si
  - **2** = No

En consecuencia, el conjunto final de registros consta de 3295 observaciones, que se distribuyen en 6 variables, de las cuales 5 son categóricas y 1 es numérica.

## 2.2 Análisis Descriptivo

A continuación se presentan tablas y gráficos descriptivos de estas variables:

Table 1: Descripción de las variables

Denominación de variables		
Variabes	Tipos	Niveles
Edad	Numérica	-
Paciente hospitalizado	Categórica	1 = Si 0 = No
Sexo	Categórica	M = Masculino F = Femenino
Tipo de régimen	Categórica	C = Contributivo E = Especial N = No asegurado P = Excepción S = Subsidiado
Vómito	Categórica	1 = Si 0 = No
Etapas de la enfermedad	Categórica	1 = Catarral 2 = Espasmódica 3 = Convaleciente

Table 2: Resumen numérico de variables categóricas

Tipo de régimen	Sexo	Paciente hospitalizado	Vómito	Etapas de la enfermedad
C = 1619	F = 1592	1 = 1649	1 = 1620	1 = 1969
E = 46	M = 1703	0 = 1646	2 = 1675	2 = 1202
N = 516	-	-	-	3 = 124
P = 55	-	-	-	-
S = 1059	-	-	-	-

Table 3: Tablas de contingencia: Sexo

	Femenino	Masculino
Hospitalizado	776	873
No hospitalizado	816	830

Table 4: Tablas de contingencia: Tipo de régimen

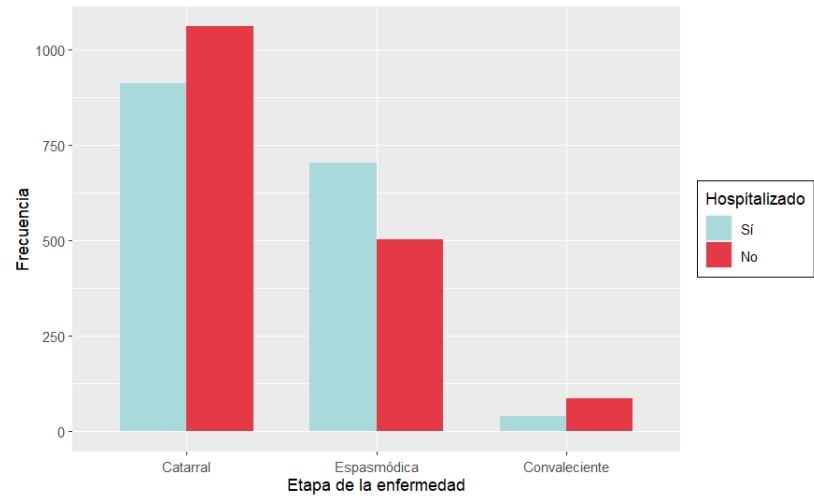
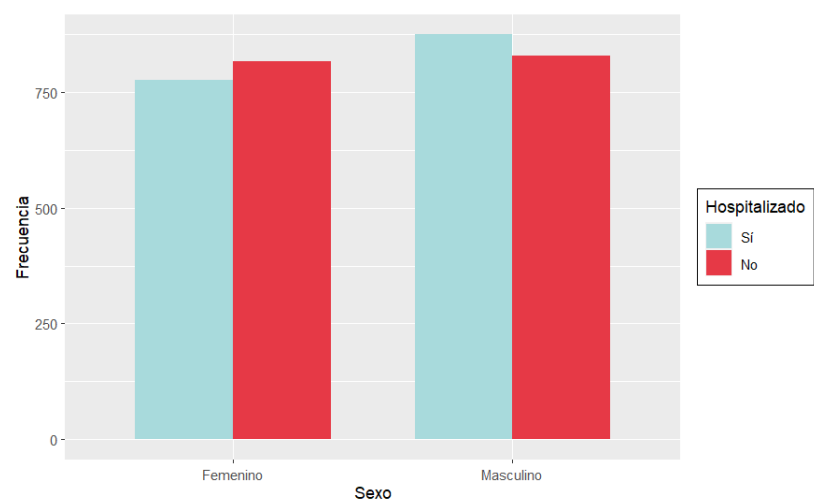
	Contributivo	Subsidiado	Especial	Excepción	No asegurado
Hospitalizado	753	670	22	20	184
No hospitalizado	866	389	24	35	332

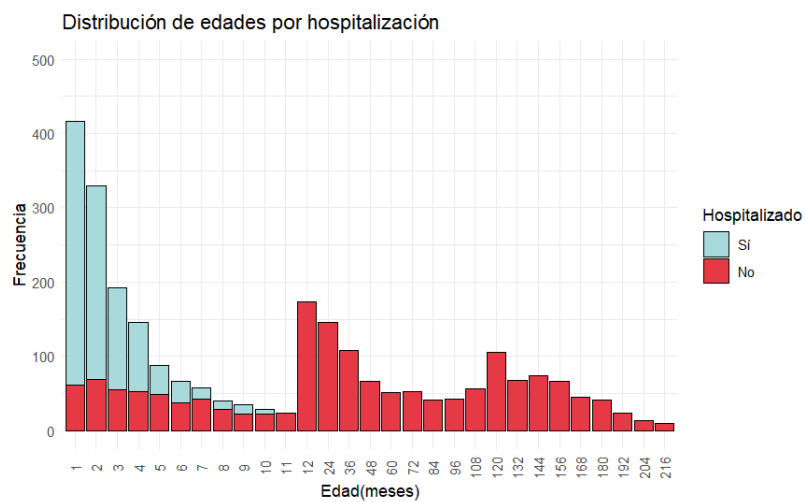
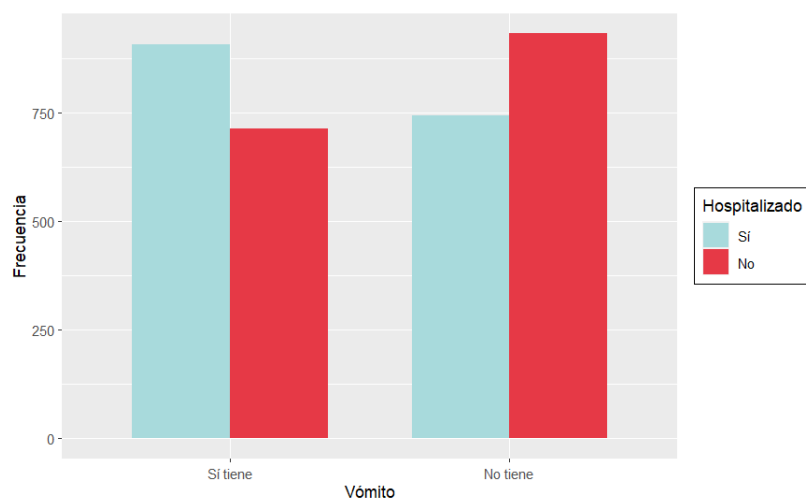
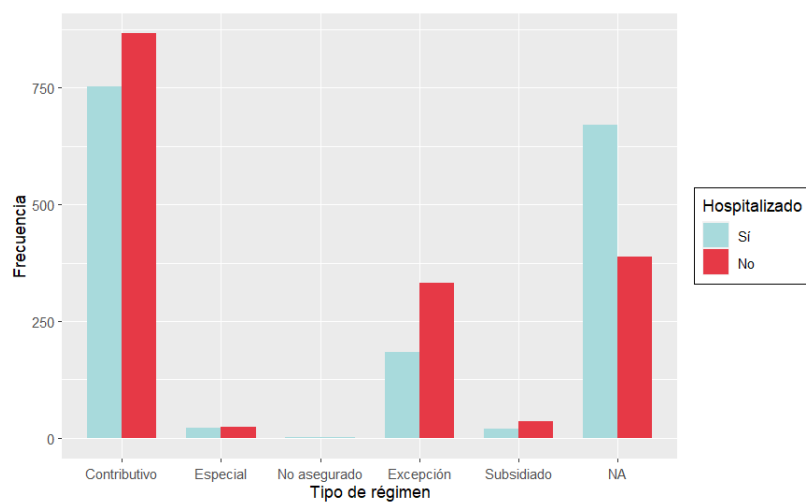
Table 5: Tablas de contingencia: Etapa de la enfermedad

	Catarral	Espasmódica	Convaleciente
Hospitalizado	909	701	39
No hospitalizado	1060	501	85

Table 6: Tablas de contingencia: Vómito

	Si tiene	No tiene
Hospitalizado	907	742
No hospitalizado	713	933





### 3 Metodología

Al realizar este estudio de la base de datos sobre la enfermedad de la tosferina, nuestra pregunta de interés es identificar las variables que pueden afectar la hospitalización de los pacientes en las respectivas IPS de las ciudades. Para responder a esta pregunta, se han propuesto tres modelos.

El primer modelo incluye las variables de edad, sexo, etapa de la enfermedad, tipo de régimen de seguridad y vómito. El segundo modelo incluye las variables de edad, tipo de régimen y etapa de la enfermedad. Por último, el tercer modelo solo contiene la variable de edad, tipo de régimen y etapa de la enfermedad.

Para responder a nuestra pregunta de interés, se utilizarán tres modelos de regresión logística. La regresión logística es uno de los métodos estadísticos más versátiles utilizados para modelar relaciones entre variables dependientes dicotómicas y varias variables independientes. Estas variables independientes pueden ser de cualquier naturaleza, lo que convierte a la regresión logística en un método estándar para el análisis de regresión cuando los datos son binarios.

En la regresión logística se tiene que:

$$Y_i \sim \text{Bernoulli}(\theta_i), \quad \beta \sim \text{Normal}(0, 100)$$

$$y_i \in \{0, 1\}, \quad i = 1, \dots, n$$

Cuyo predictor lineal es:

$$\text{Logit}(\theta_i) = \ln \left( \frac{\theta_i}{1 - \theta_i} \right) = \mathbf{X}\beta$$

$$\theta_i = \frac{e^{\mathbf{X}\beta}}{1 + e^{\mathbf{X}\beta}}$$

#### 3.1 Modelo 1

Para el primer modelo, la variable respuesta es:

- $Y$ : Paciente hospitalizado (Dicotómica)

Y las covariables son:

- $x_1$ : Edad (Continua)
- $x_2$ : Tipo de régimen de seguridad (Categórica)
- $x_3$ : Sexo (Categórica)
- $x_4$ : Etapa de la enfermedad (Categórica)
- $x_5$ : Vómito (Categórica)

Usando rstan, para este modelo se obtuvo:

Para este modelo las betas que dieron significativos son  $\beta[1], \beta[2], \beta[5], \beta[6]$  y  $\beta[8]$ , al no incluir el 0 en el intervalo del 95% esto se puede ver en el anexo.

#### 3.2 Modelo 2

Para este modelo no se tuvo en cuenta la variable sexo ya que en los gráficos descriptivo se comporta de forma muy similar en sus categorías, consideramos que esta variable no tendría tanto efecto al momento de responder nuestra pregunta de interés. En este caso se tiene que:

- $Y$ : Paciente hospitalizado (Dicotómica)

y las covariables son:

- $x_1$ : Edad (Continua)

Table 7: Resumen modelo 1

Betas	Media	SD	2.5%	97.5%	n_eff	Rhat
$\beta[1]$	0.56	0.26	0.05	1.06	4296	1
$\beta[2]$	-0.05	0.00	-0.06	-0.05	12776	1
$\beta[3]$	0.07	0.35	-0.60	0.76	11029	1
$\beta[4]$	0.04	0.13	-0.22	0.30	11272	1
$\beta[5]$	-0.66	0.31	-1.27	-0.06	11753	1
$\beta[6]$	0.88	0.10	0.68	1.07	10688	1
$\beta[7]$	0.18	0.25	-0.32	0.67	4392	1
$\beta[8]$	0.52	0.26	0.01	1.02	4477	1
$\beta[9]$	0.02	0.09	-0.16	0.18	11526	1
$\beta[10]$	-0.13	0.09	-0.30	0.04	12722	1

- $x_2$ : Tipo de régimen de seguridad (Categórica)
- $x_4$ : Etapa de la enfermedad (Categórica)
- $x_5$ : Vómito (Categórica)

Usando rstan, para este modelo se obtuvo:

Table 8: Resumen modelo 2

Betas	Medias	SD	2.5%	97.5%	n_eff	Rhat
$\beta[1]$	0.56	0.26	0.05	1.06	3652	1
$\beta[2]$	-0.05	0.00	-0.06	-0.05	15146	1
$\beta[3]$	0.06	0.35	-0.62	0.76	11020	1
$\beta[4]$	0.04	0.13	-0.22	0.30	10594	1
$\beta[5]$	-0.66	0.31	-1.28	-0.07	11941	1
$\beta[6]$	0.87	0.10	0.68	1.07	10618	1
$\beta[7]$	0.18	0.25	-0.32	0.68	3611	1
$\beta[8]$	0.52	0.26	0.01	1.04	3755	1
$\beta[9]$	-0.13	0.09	-0.30	0.04	10923	1

Para este modelo las betas que dieron significativos son  $\beta[1], \beta[2], \beta[5], \beta[6]$  y  $\beta[8]$ , al no incluir el 0 en el intervalo del 95% esto se puede ver en el anexo.

### 3.3 Modelo 3

Para este modelo no se obtuvo en cuenta las variables vómito y sexo. En este caso se tiene:

- $Y$ : Paciente hospitalizado (Dicotómica)

y las covariables son:

- $x_1$ : Edad (Continua)
- $x_2$ : Tipo de régimen de seguridad (Categórica)
- $x_4$ : Etapa de la enfermedad (Categórica)

Usando rstan, para este modelo se obtuvo:

Para este modelo las betas que dieron significativos son  $\beta[1], \beta[2], \beta[5]$  y  $\beta[6]$ , al no incluir el 0 en el intervalo del 95% esto se puede ver en el anexo.



Table 9: Resumen modelo 3

Betas	Media	SD	2.5%	97.5%	n_eff	Rhat
$\beta[1]$	0.50	0.25	0.00	1.00	4348	1
$\beta[2]$	-0.05	0.00	-0.06	-0.05	11553	1
$\beta[3]$	0.07	0.35	-0.59	0.76	7923	1
$\beta[4]$	0.04	0.13	-0.22	0.29	8100	1
$\beta[5]$	-0.64	0.30	-1.25	-0.05	8302	1
$\beta[6]$	0.87	0.10	0.68	1.08	6890	1
$\beta[7]$	0.18	0.25	-0.33	0.68	4421	1
$\beta[8]$	0.50	0.26	-0.01	1.01	4432	1

## 4 Resultados

Ya que se tiene 3 modelos se quiere escoger el modelo que mejor ajuste para resolver nuestra pregunta de interés. Para esto se utilizó 3 métodos de decisión: las diferencias del DIC, la curva ROC y LOO (Leave-one-out cross-validation).

### 4.1 DIC

DIC es un método bayesiano para la comparación de modelos es una versión bayesiana del criterio de información de AKAIKE (AIC).

El DIC está definido como:

$$DIC = -2\log(f(y|\hat{\theta}_{Bayes}) + 2P_{DIC})$$

Donde  $\hat{\theta}_{Bayes} = E[\theta|y]$  y  $P_{DIC}$  es el número efectivo de parámetros.

Usando el DIC, se selecciona el modelo con el DIC menor, por lo tanto, diferencias de más de 10 en el DIC permiten descartar el modelo con mayor DIC.

Se realizó las respectivas diferencias entre los DIC de cada modelo, en donde se obtuvo:

Table 10: Resultado de los DICs

MODELO	DIC
1	4428.323
2	4413.04
3	4379.276

Table 11: Diferencia entre DICs

$DIC_i - DIC_j$	1	2	3
1	0	15.283	49.047
2	-15.283	0	33.764
3	-49.047	-33.764	0

Ya que las diferencia entre los DIC de cada modelo dio mayor a 10. Se escogió el modelo con menor DIC, es decir el modelo 3, cuyo DIC es igual a 4379.276.

### 4.2 LOO

Es una técnica de validación cruzada que estima la capacidad predictiva de un modelo utilizando datos de validación generados dejando uno de los puntos de datos a la vez. El LOO proporciona una estimación de la pérdida esperada fuera de la muestra y se puede utilizar para comparar modelos. Un valor LOO más bajo indica un mejor rendimiento predictivo. Cuanto menor sea el valor, mayor será

la precisión de la estimación. De la tabla siguiente, tenemos que P\_LOO es la medida de precisión de la estimación de ELPD utilizando LOO.

Table 12: Resultado de LOO

MODELO	P_LOO
1	12
2	10.8
3	10

En este caso se escoge el modelo 3 al tener un P\_LOO más bajo.

### 4.3 Curva ROC

La curva ROC ilustra la sensibilidad y especificidad de cada uno de los posibles puntos de corte de un test diagnostico cuya escala de medición es continua, es una representación gráfica de la sensibilidad frente a la especificidad para un sistema de clasificación binario según se varía el umbral de discriminación. Otra interpretación de este grafico es la representación de la razón o proporción de verdaderos positivos (VPR = Razón de verdaderos positivos) frente a la razón o proporción de falsos positivos (FPR = razón de falsos positivos) también según se varía el umbral de discriminación (valor a partir el cual decidimos que caso es un positivo).

La curva ROC define como el mejor método posible de predicción aquel que se sitúa en la coordenada (0,1), del espacio ROC. REPRESENTANDO UN 100% de sensibilidad (ningún falso negativo) y un 100% también es especificidad (ningún falso positivo).

Se realizo la respectiva curva ROC para cada modelo, obteniendo los siguientes resultados:

Table 13: Resultado de curva ROC

MODELO	auc
1	0.8564
2	0.8562
3	0.8558

En este método, no hay una diferencia significativa entre los tres modelos, para escoger el mejor.

## 5 Modelo escogido

Con los métodos de comparación: DIC y LOO, el modelo 3 es el mejor y aunque en la curva ROC no es el mejor, tampoco hay una diferencia significativa con los demás, por eso es que nuestro modelo escogido es el 3.

A continuación presentamos la tabla del modelo escodigo, con el nombre de las variables que representa los betas.

Table 14: Estimación puntual modelo seleccionado

Betas	Variables	Media	2.5%	97.5%
$\beta[1]$	Interceto	0.50	0.00	1.00
$\beta[2]$	Edad	-0.05	-0.06	-0.05
$\beta[3]$	Seguridad social: Especial	0.07	-0.59	0.76
$\beta[4]$	Seguridad social: No asegurado	0.04	-0.22	0.29
$\beta[5]$	Seguridad social: Pendiente	-0.64	-1.25	-0.05
$\beta[6]$	Seguridad social: Subsidiado	0.87	0.68	1.08
$\beta[7]$	Etapas 1: Catarral	0.18	-0.33	0.68
$\beta[8]$	Etapas 2: Espasmódica	0.50	-0.01	1.01

En seguida, procedemos a realizar el cálculo de la exponencial de las estimaciones de los  $\beta_i$ , ya que el predictor lineal está dado de la siguiente forma:

$$\text{Logit}(\theta_i) = \log\left(\frac{\theta_i}{1 - \theta_i}\right) = \mathbf{X}\beta$$

$$\text{Logit}(\theta_i) = \log(odds) = \mathbf{X}\beta$$

$$odds = e^{\mathbf{X}\beta}$$

Table 15: Estimación Odds

$e^{\beta_i}$	Variables	Media	2.5%	97.5%
$e^{\beta_1}$	Intercepto	1.648	1	2.718
$e^{\beta_2}$	Edad	0.951	0.942	0.951
$e^{\beta_5}$	S.S: Pendiente	0.527	0.286	0.951
$e^{\beta_6}$	S.S: Subsidiado	2.387	1.974	2.945

## 6 Conclusiones

A continuación, presentamos conclusiones de los resultados obtenidos con los tres modelos, y también inferencias sobre el modelo escogido:

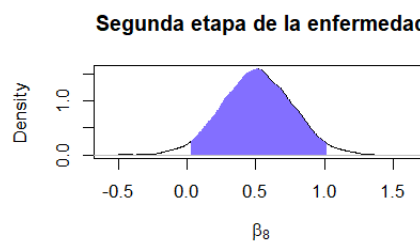
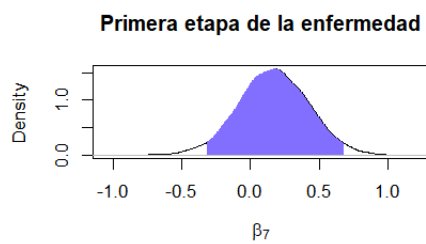
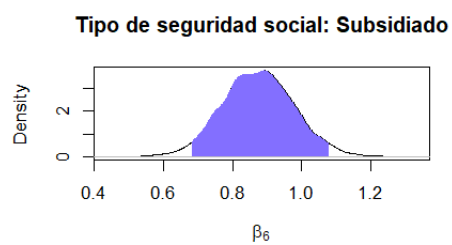
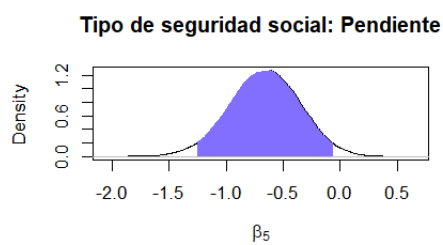
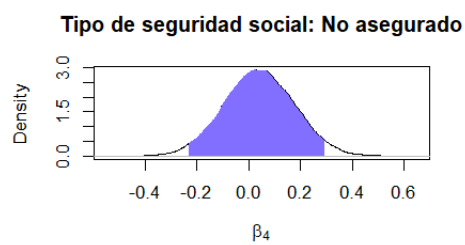
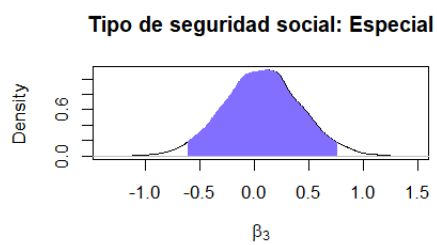
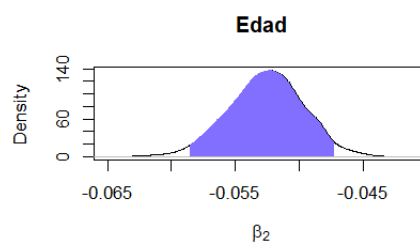
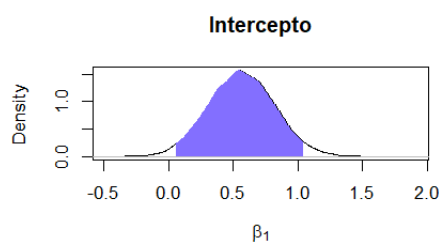
- El modelo 3, con menos covariables, es el que mejor se ajusta a la variable respuesta.
- Podemos observar que los tres modelos convergen, ya que tienen un Rhat igual a 1, y también podemos ver en los diagnósticos de convergencia (Figura), no hay cambios drásticos, entonces se puede decir que llegaron a la función estacionaria.
- Por el incremento de un mes en la edad de un niño, la probabilidad de ser hospitalizado se reduce en un 4.9%.
- Los niños pertenecientes a seguridad social pendiente, presentan una reducción del 47.3% en la probabilidad de ser hospitalizado, comparados con los niños pertenecientes a seguridad social contributivo.
- Los niños pertenecientes a seguridad social subsidiados, presentan un aumento del 138.7% en la probabilidad de ser hospitalizado, comparados con los niños pertenecientes a seguridad social contributivo.

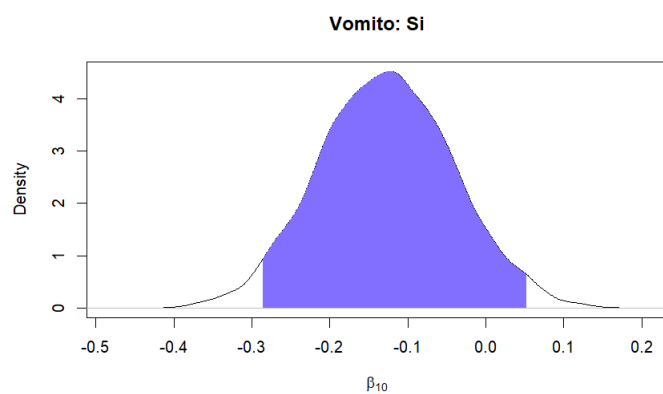
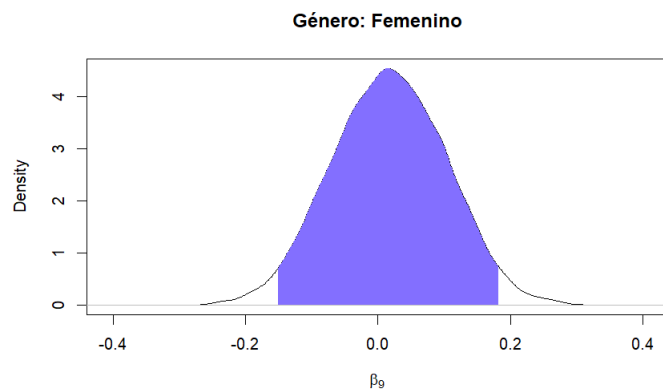
## 7 Discusión

## 8 Anexos

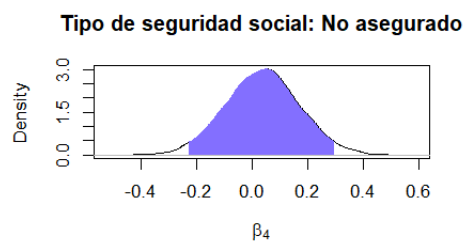
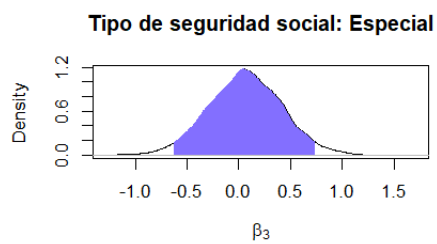
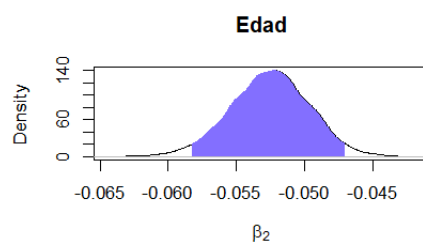
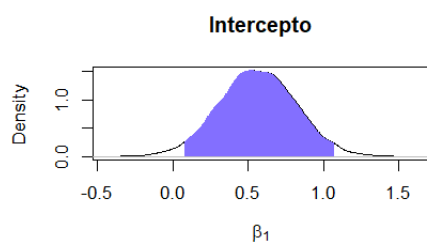
### 8.1 Intervalos de densidad

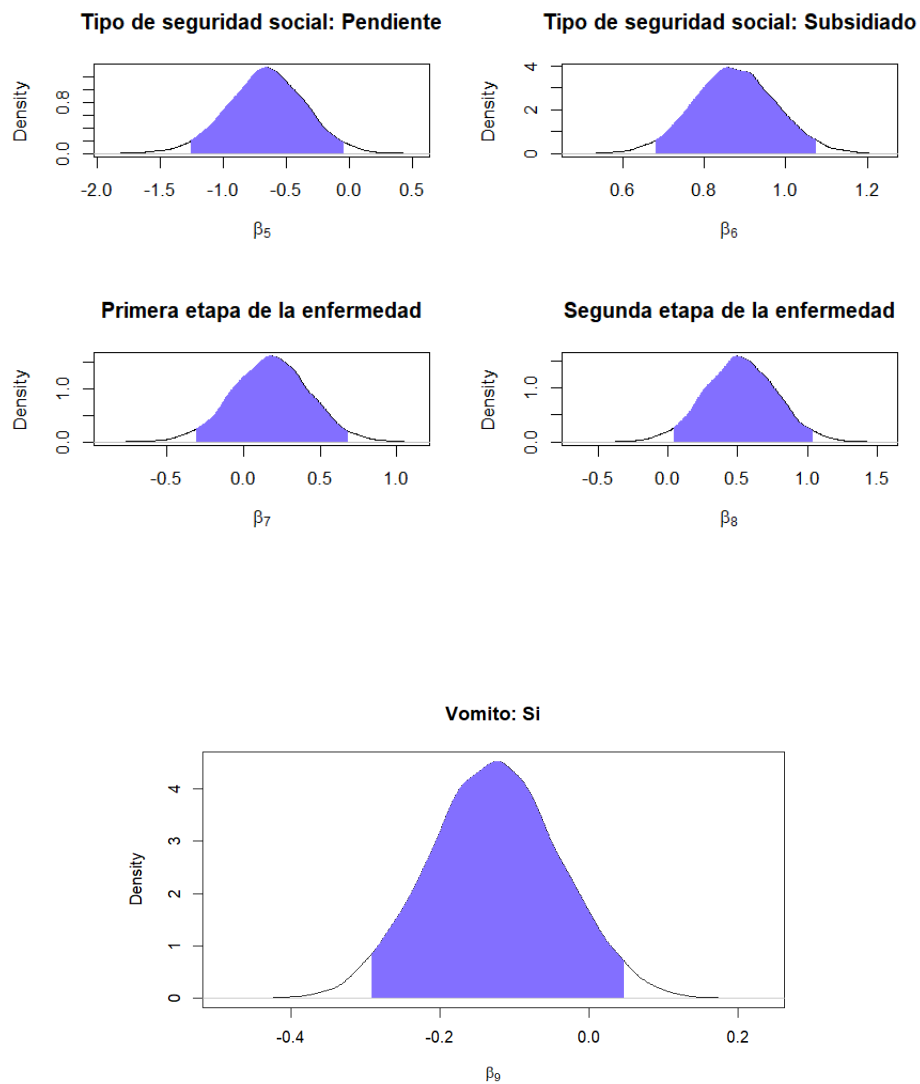
#### 8.1.1 Modelo 1



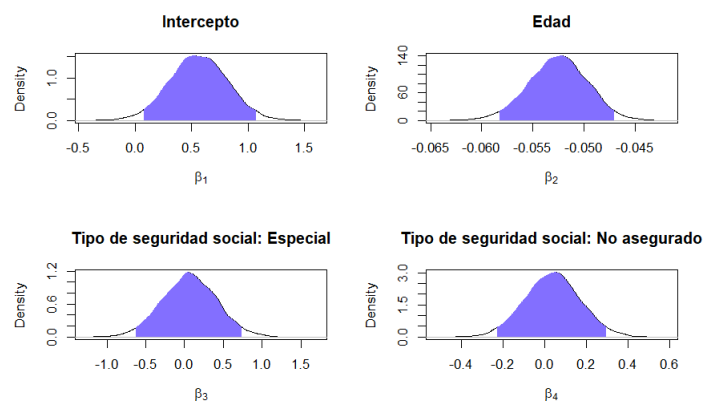


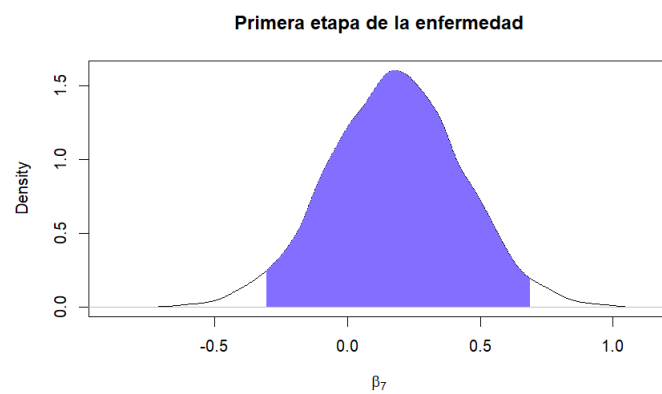
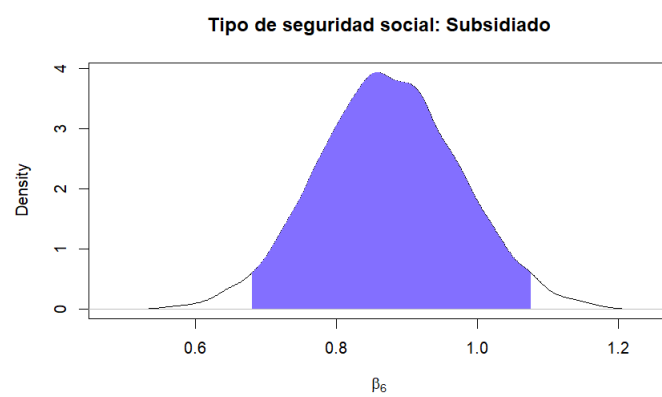
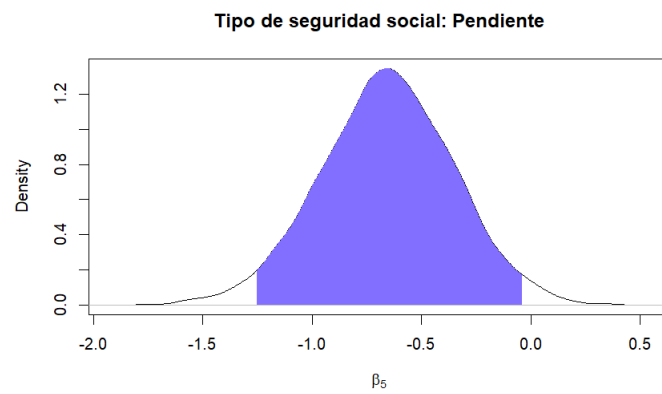
### 8.1.2 Modelo 2

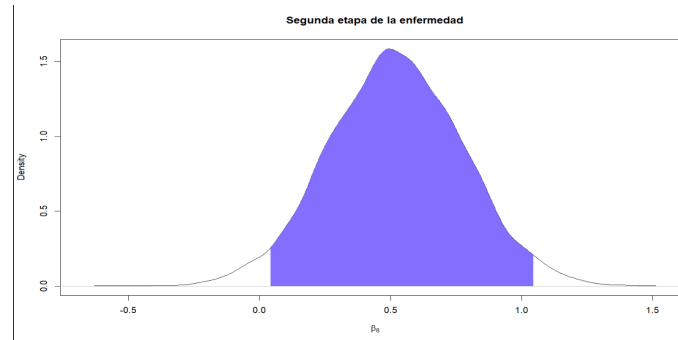




### 8.1.3 Modelo 3

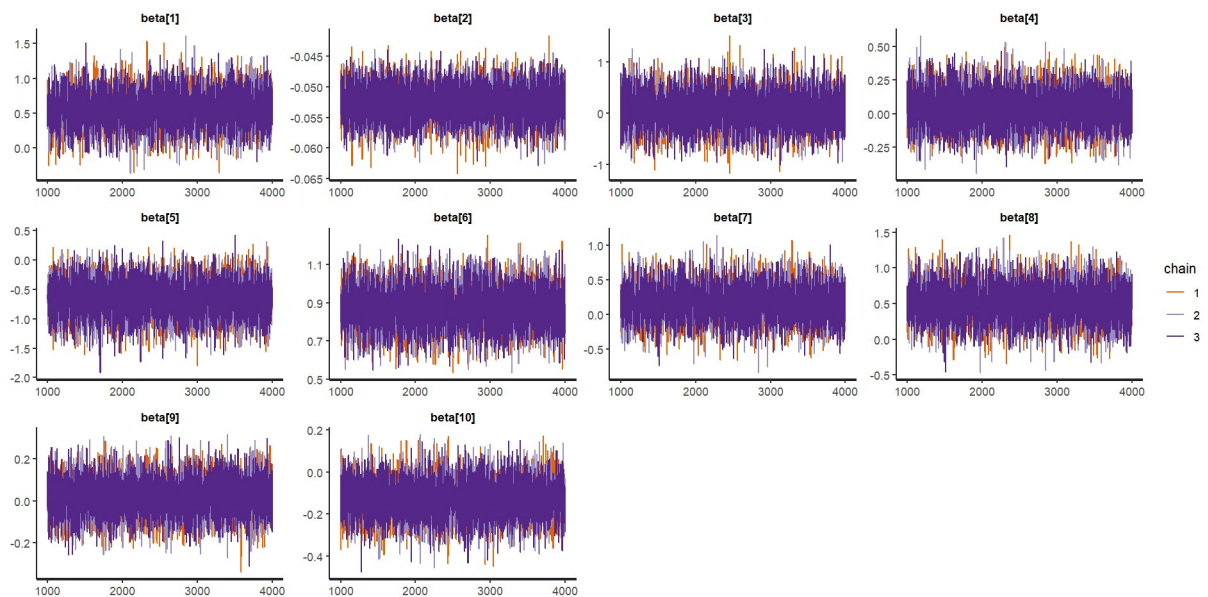






## 8.2 Traceplot

### 8.2.1 Modelo 1



### 8.2.2 Modelo 2

### 8.2.3 Modelo 3

## 8.3 Gráficas de autocorrelación

### 8.3.1 Modelo 1

### 8.3.2 Modelo 2

### 8.3.3 Modelo 3

## 8.4 Gráficas de Curva de ROC

### 8.4.1 Modelo 1

### 8.4.2 Modelo 2

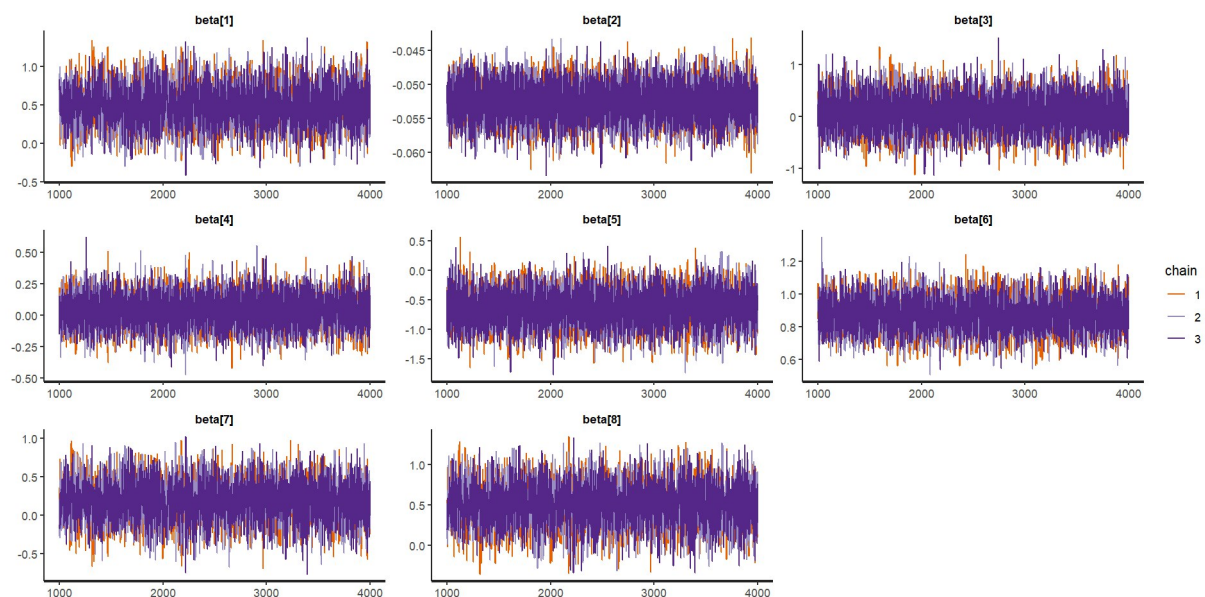
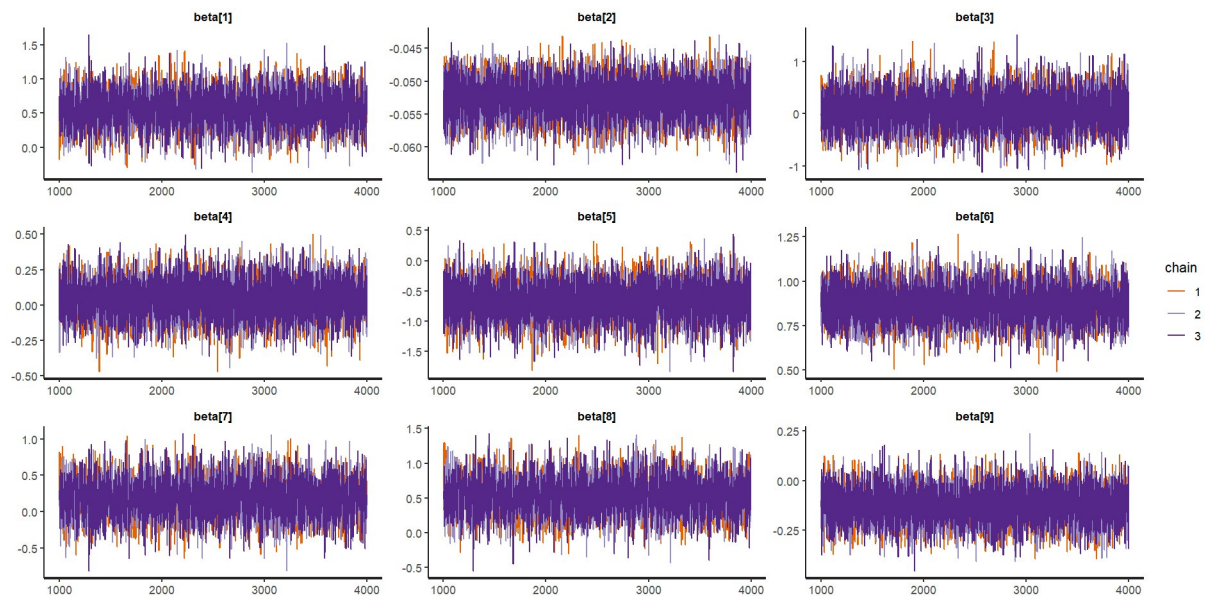
### 8.4.3 Modelo 3

## References

INFORMACIÓN DE MATEO MURCIA VALLES

*E-mail:* mmurciav@unal.edu.co





INFORMACIÓN DE CHRISTOPHER ANDRÉS OBANDO

*E-mail:* cobandor@unal.edu.co

INFORMACIÓN DE JUAN DAVID GARCÍA ZAPATA

*E-mail:* jgarciaza@unal.edu.co

