



UNIVERSIDAD NACIONAL DE COLOMBIA

Comparación de Intervalos de Confianza para el Coeficiente de Correlación

Liliana Vanessa Pacheco Galindo

Universidad Nacional de Colombia
Facultad de Ciencias, Escuela de Estadística
Medellín, Colombia
2013

Comparación de Intervalos de Confianza para el Coeficiente de Correlación

Liliana Vanessa Pacheco Galindo

Tesis o trabajo de grado presentada(o) como requisito parcial para optar al título de:
Magíster en Ciencias - Estadística

Director:
(Ph.D.) Juan Carlos Correa Morales

Línea de Investigación:
Análisis Multivariado
Universidad Nacional de Colombia
Facultad de Ciencias, Escuela de Estadística
Medellín, Colombia
2013

Dedicatoria

Dedico este trabajo de tesis principalmente a Dios, quien me ha acompañado hasta estas instancias de mi vida y ha hecho de mí quien ahora soy. A mis padres y a mis hermanos, que han sido siempre mi apoyo, guía y fuente de amor.

A mis compañeros de estudio y a mis profesores, quienes sin su ayuda nunca hubiera podido hacer esta tesis. Y en especial a aquel que me acompañó y fue mi apoyo incondicional, D.A.S.

Resumen

La construcción de intervalos de confianza para la estimación del coeficiente de correlación en la distribución normal bivariada, digamos ρ , es un problema importante en el trabajo estadístico aplicado. Uno de los propósitos principales de este trabajo es realizar una revisión de los diferentes procedimientos para su construcción, en la distribución normal bivariada. Mediante un estudio de simulación se analiza el comportamiento de los niveles de confianza reales y se comparan con los teóricos. Se estudia además el comportamiento de las longitudes de los intervalos de confianza logrados por nueve métodos considerados en la Estadística clásica y dos Intervalos de credibilidad construidos mediante el enfoque Bayesiano para así determinar cuál metodología provee los intervalos más cortos y de nivel real más cercano al nominal.

Además, se propone un indicador que resume de manera más efectiva la calidad del intervalo analizado. Dicho estudio de simulación se desarrolló empleando el software R (R Development Core Team 2010) para construir los intervalos de confianza, las distribuciones muestrales de diversos estadísticos utilizados y obtener las gráficas de resumen de resultados obtenidos.

Dentro del enfoque clásico hay ciertos procedimientos que generan mejores resultados para muestras pequeñas, mientras que en el enfoque Bayesiano las conclusiones no son homogéneas en cuanto a la selección de la “mejor” distribución a priori para ρ .

Palabras clave: Estimación por intervalo, Coeficiente de Correlación, Distribución Normal Bivariada, Muestreador de Gibbs, Intervalos de Credibilidad.

Abstract

The construction of confidence intervals for the estimate of the correlation coefficient in the bivariate normal distribution, say ρ , is an important problem in applied statistical work. One of the main purposes of this paper is to review the different procedures for their construction, in the case of the bivariate normal distribution.

Through a simulation study we analyse the behavior of real confidence levels and compare them with the theoretical ones. We also analyse the behavior of the confidence interval's lengths achieved by nine methods considered in classical statistics and two credibility intervals using the Bayesian methodology to determine which provides the shorter intervals and a coverage probability closer to the nominal one.

Furthermore, we propose an indicator that summarizes even more effectively the analyzed interval quality. This simulation study was developed using the R software (R Development Core Team 2010) to construct confidence intervals, sampling distributions of various statistics used in this paper and to get summary results graphs.

Within the classical approach there are certain procedures that generate better results for smaller samples, while the Bayesian approach conclusions are not homogeneous in terms of the selection of the “best” apriori distribution for ρ .

Keywords: Interval Estimation, Correlation Coefficient, Bivariate Normal Distribution, Gibbs Sampler, Credibility Intervals

Contenido

Resumen	vii
1. Introducción	2
2. Metodologías en la construcción de Intervalos de confianza caso Normal Bivariada	3
2.1. Distribución de Probabilidad Normal Bivariada y Multivariada	3
2.2. Coeficiente de correlación	4
2.3. Intervalos de confianza	10
2.3.1. Método I: Basado en la transformación Arco tangente	10
2.3.2. Método II: Intervalo de la Razón de Verosimilitud	11
2.3.3. Método III: Bootstrap	12
2.3.4. Método IV: Intervalo de Jeyaratnam	12
2.3.5. Método V: Test Generalizado para ρ	12
3. Resultados de la Simulación - Metodología Clásica I	14
3.1. Caso Normal Bivariada	14
4. Metodología Bayesiana para la construcción de Intervalos de Credibilidad	37
4.1. Inferencia Bayesiana	37
4.2. Intervalos Bayesianos o Intervalos de credibilidad	38
4.3. Inferencias para el coeficiente de correlación ρ	38
4.3.1. Selección de Distribuciones Apriori para ρ	39
4.3.2. Obtención de Distribuciones Aposteriori para ρ	42
5. Resultados de la Simulación - Metodología Bayesiana	47
6. Aplicaciones	54
6.1. Aplicación Base de datos Huevos	54
6.2. Aplicación Base de Datos Vinos	55
7. Conclusiones	57
8. Recomendaciones	58

9. Anexos	59
9.1. Tablas de Resultados Metodología Clásica	59
9.2. Inferencias para el coeficiente de correlación ρ	63
9.2.1. Verosimilitud Simplificada	63
9.2.2. Distribuciones condicionales	65
10. Códigos en R	67
10.1. Intervalos de Confianza parte clásica	67
10.2. Error cuadrático medio de los estimadores de ρ	72
10.3. Intervalo de confianza Bayesiano con la apriori 2	73
10.4. Índices de resumen	76
10.5. Construcción de Clústers	77
Bibliografía	78

1 Introducción

El coeficiente de correlación es una de las medidas estadísticas de más uso dentro del trabajo aplicado. Algunas de sus propiedades fueron estudiadas por Zheng and Matis (1994), quienes presentan y demuestran algunas de sus propiedades.

Debido a su amplia utilización, varias son sus interpretaciones. Falk and Well (1997) sustentan que el coeficiente de correlación de Pearson, ρ , es ampliamente usado en campos como la educación, psicología, y todas las ciencias sociales, y el concepto es empleado en diversas metodologías de tipo estadístico.

La estimación del coeficiente de correlación por medio de intervalos es importante, y para ello se disponen de diversos métodos. La metodología quizá más conocida es la propuesta originalmente por Fisher (1915) en la cual se realiza una transformación del coeficiente de correlación muestral, r , y asumiendo normalidad asintótica, se desarrolla un intervalo para el coeficiente de correlación poblacional ρ , (Krishnamoorthy and Xia, 2007). También se conocen transformaciones adicionales hechas por Hotelling (1953) a la propuesta inicial de Fisher. La estadística bayesiana presenta a su vez metodologías para la construcción de intervalos de credibilidad para parámetros distribucionales (Bernardo and Smith, 2000).

El problema para el analista es la carencia de reglas sobre cuál fórmula es preferible. Para esto se pretende realizar un estudio de simulación que permita analizar el comportamiento de los niveles de confianza reales y compararlos con los teóricos de los diversos intervalos disponibles. Así como también, hacer una comparación de las longitudes del intervalo obtenido por las diferentes metodologías y la implementación de un indicador que permita relacionar los dos criterios de evaluación anteriormente mencionados.

Algunas de las metodologías empleadas para la construcción de los intervalos de confianza pueden encontrarse en: (Fisher, 1921), (Hotelling, 1953), (Pawitan, 2001), (Efron, 1979) y (Krishnamoorthy and Xia, 2007). Y para la construcción de los intervalos de credibilidad se emplean las distribuciones a priori de (McCullagh, 1989) y el kernel de la distribución empleada por (Ghosh et al., 2010).

2 Metodologías en la construcción de Intervalos de confianza caso Normal Bivariada

2.1. Distribución de Probabilidad Normal Bivariada y Multivariada

Si se tiene una distribución normal univariada, con media μ y varianza σ^2 , esta tendrá la siguiente función de densidad (Johnson and Wichern, 2007):

$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2}\left[\frac{(x-\mu)}{\sigma}\right]^2}; \quad -\infty < x < \infty \quad (2-1)$$

Una función de densidad de una normal p -dimensional, para el vector aleatorio $X' = [X_1, X_2, \dots, X_p]$ tiene la forma:

$$f(x) = \frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} e^{-\frac{(x-\mu)'\Sigma^{-1}(x-\mu)}{2}} \quad (2-2)$$

$$-\infty < x_i < \infty; \quad i = 1, 2, \dots, p$$

donde el vector $\mu \in \Re^p$ representa el valor esperado del vector aleatorio \mathbf{X} y la matriz $\Sigma_{p \times p}$, simétrica y definida positiva, es la matriz de varianzas y covarianzas de \mathbf{X} .

La distribución Normal multivariada corresponde a la generalización de la densidad de una normal univariada a $p \geq 2$ dimensiones. Específicamente, se denota el caso de la distribución normal bivariada, es decir con $p = 2$, en términos de los parámetros $\mu_1 = E(X_1)$, $\mu_2 = E(X_2)$, $\sigma_{11} = \text{Var}(X_1)$, $\sigma_{22} = \text{Var}(X_2)$, $\sigma_{12} = \text{Cov}(X_1, X_2)$, así (Johnson and Wichern, 2007):

$$f(x_1, x_2) = \frac{1}{2\pi\sqrt{\sigma_{11}\sigma_{22}(1-\rho_{12}^2)}} \exp \left\{ -\frac{1}{2(1-\rho_{12}^2)} \left[\left(\frac{x_1 - \mu_1}{\sqrt{\sigma_{11}}} \right)^2 + \left(\frac{x_2 - \mu_2}{\sqrt{\sigma_{22}}} \right)^2 - 2\rho_{12} \left(\frac{x_1 - \mu_1}{\sqrt{\sigma_{11}}} \right) \left(\frac{x_2 - \mu_2}{\sqrt{\sigma_{22}}} \right) \right] \right\} \quad (2-3)$$

2.2. Coeficiente de correlación

La interpretación quizás más conocida del concepto de coeficiente de correlación es la siguiente: índice del grado de cercanía o relación lineal entre dos variables aleatorias. Para la distribución Normal Bivariada se tiene un estimador muestral para este coeficiente así:

Suponga que se tienen n pares de observaciones $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ de una distribución normal bivariable con vector de medias μ y matriz de varianzas y covarianzas Σ . El estimador máximo verosímil para ρ es (Zheng and Matis, 1994) está dado por:

$$R = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\left(\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (y_i - \bar{y})^2 \right)^{1/2}} \tag{2-4}$$

Este estimador es asintóticamente insesgado, y además, se puede demostrar que es un estimador consistente para ρ . La distribución de R está dada por, Graybill (1976):

$$f_R(r) = \frac{(n-2)(1-\rho^2)^{\frac{(n-1)}{2}}}{\pi} (1-r^2)^{\frac{(n-4)}{2}} \int_0^\infty (\cosh w - \rho r)^{-(n-1)} dw \tag{2-5}$$

donde $-1 < r < 1$ y $-1 < \rho < 1$. El único parámetro de la distribución es ρ , (Fisher, 1915). La igualdad en $-1 \leq r \leq 1$ se alcanza si y solo si los datos se distribuyen a lo largo de una perfecta línea recta en un diagrama de dispersión.

En las siguientes gráficas observamos la distribución de R dependiendo de diferentes tamaños de muestra.

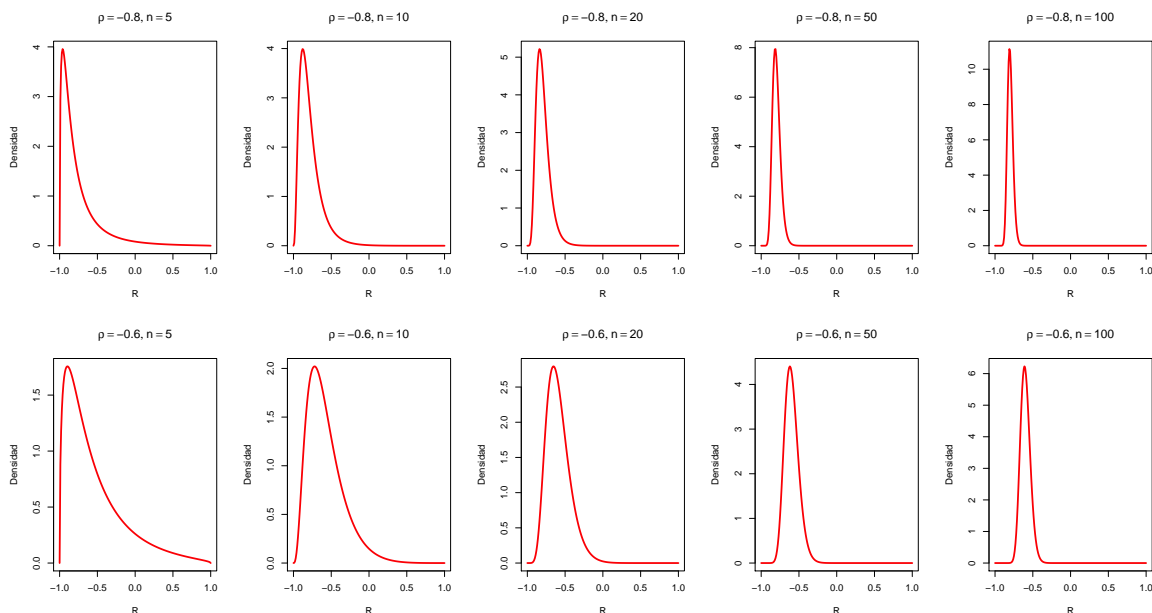


Figura 2-1: Distribución para R con diferentes tamaños de muestra y valores de $\rho < 0$.

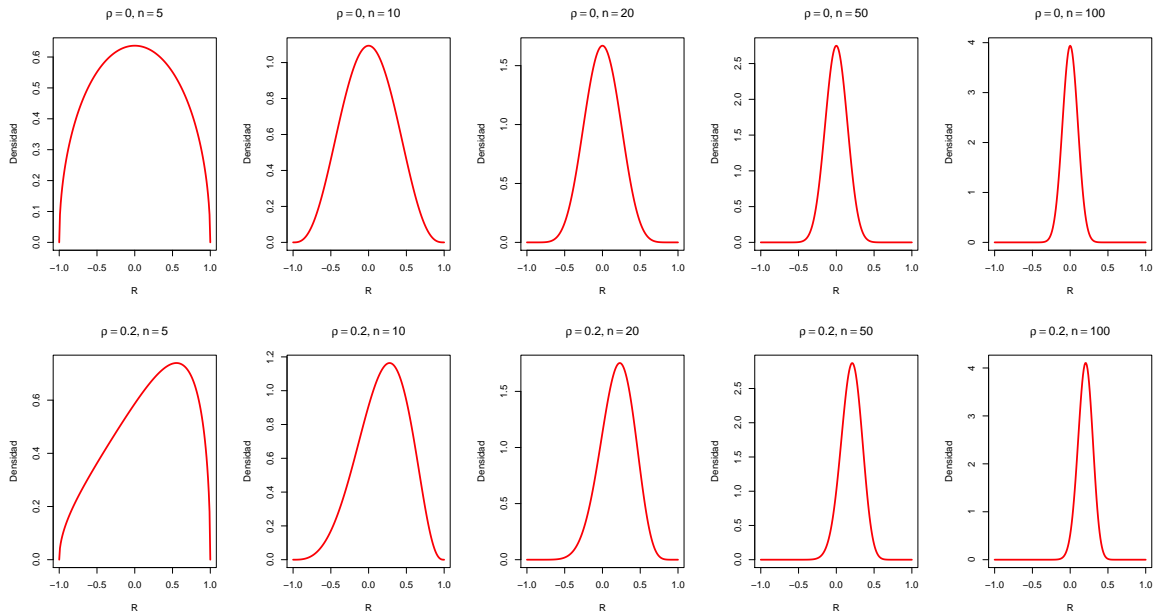


Figura 2-2: Distribución para R con diferentes tamaños de muestra y diferentes valores de $\rho \geq 0$.

El UMVUE para ρ es el siguiente, Graybill (1976):

$$\hat{\rho} = R \left(\frac{\Gamma(\frac{n-2}{2})}{\Gamma(\frac{1}{2})\Gamma(\frac{n-3}{2})} \right) \int_0^1 \frac{t^{-\frac{1}{2}}(1-t)^{\frac{(n-5)}{2}}}{\sqrt{1-t(1-R^2)}} dt \quad (2-6)$$

Comparación de los Estimadores para ρ En los siguientes resúmenes gráficos y tabulares se muestra una comparación sobre el desempeño de los estimadores presentados anteriormente, (2-4) y (2-6), teniendo en cuenta el criterio de Error Cuadrático medio, para tamaños de muestra pequeños ($n = 5, 10, 20, 30$):

Si se define el Error Cuadrático Medio para un estimador $\hat{\theta}$ de θ así:

$$ECM(\hat{\theta}) = E[(\hat{\theta} - \theta)^2] = V[\hat{\theta}] + [E[\hat{\theta}] - \theta]^2 \quad (2-7)$$

el error cuadrático medio de un estimador $\hat{\theta}$ del parámetro θ se puede descomponer en términos de su varianza y el cuadrado del sesgo.

En la tabla 1.1. se presentan los cálculos del E.C.M. para los dos estimadores, EMV y UMVUE, y diferentes combinaciones para ρ y n .

ρ	n	EMV	UMVUE	ρ	n	EMV	UMVUE
0.0	5	0.25642953	0.32251825	0.1	5	0.24752902	0.30568290
	10	0.10409727	0.12779608		10	0.10708297	0.11514579
	20	0.05510656	0.05348213		20	0.05012973	0.05278830
	30	0.03480801	0.03753325		30	0.03435757	0.03523505
0.2	5	0.23609830	0.30389655	0.3	5	0.22291210	0.28322734
	10	0.09814702	0.11358335		10	0.09554574	0.10725311
	20	0.05357289	0.04657601		20	0.04392274	0.04659670
	30	0.03234935	0.03214782		30	0.02676595	0.03007831
0.4	5	0.22451309	0.26995101	0.5	5	0.17602590	0.19658625
	10	0.08144690	0.09309404		10	0.06919239	0.08216749
	20	0.03903870	0.04038843		20	0.03196183	0.03185162
	30	0.02569571	0.02409671		30	0.02137345	0.02141864
0.6	5	0.16219215	0.15179799	0.7	5	0.122614325	0.121207262
	10	0.04946531	0.05837362		10	0.038438604	0.034287258
	20	0.02380145	0.02101844		20	0.016023055	0.017225477
	30	0.01584849	0.01587491		30	0.009752949	0.009047878
0.8	5	0.072465378	0.064198423	0.9	5	0.0366800982	0.026697498
	10	0.022266161	0.019435461		10	0.007787224	0.005516347
	20	0.007676734	0.008228831		20	0.002506932	0.002418807
	30	0.004804023	0.005413279		30	0.001337051	0.001221097

Tabla 2-1: ECM para ambos estimadores

en las figuras 1.5. a 1.9. se muestran los gráficos del E.C.M. contra el tamaño de muestra n , para diferentes valores de ρ .

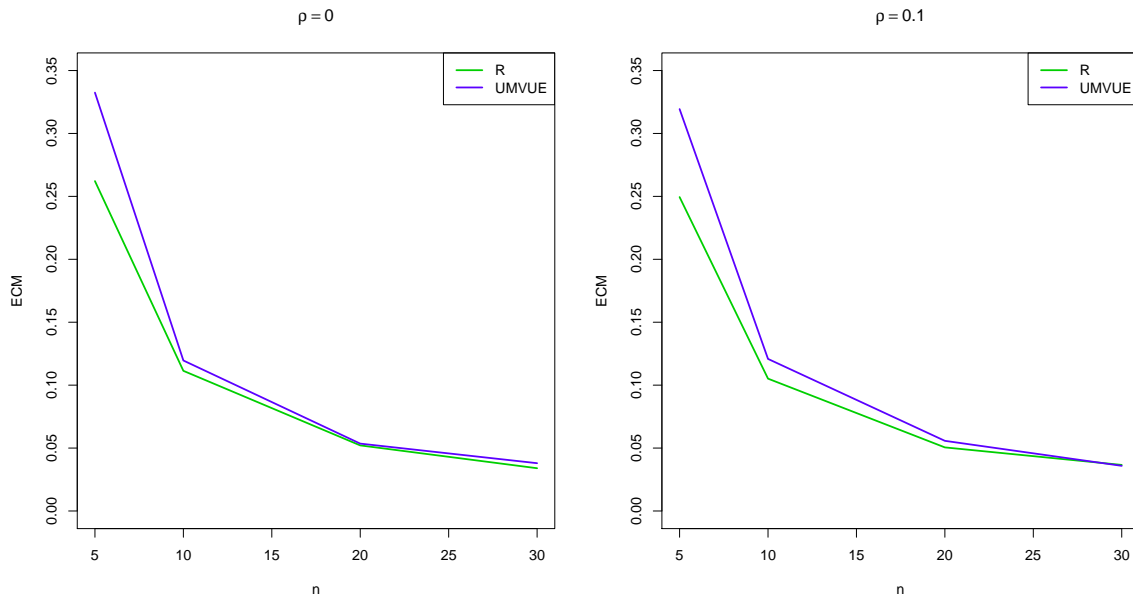


Figura 2-3: ECM para ambos estimadores con $\rho = 0.0$ y $\rho = 0.1$

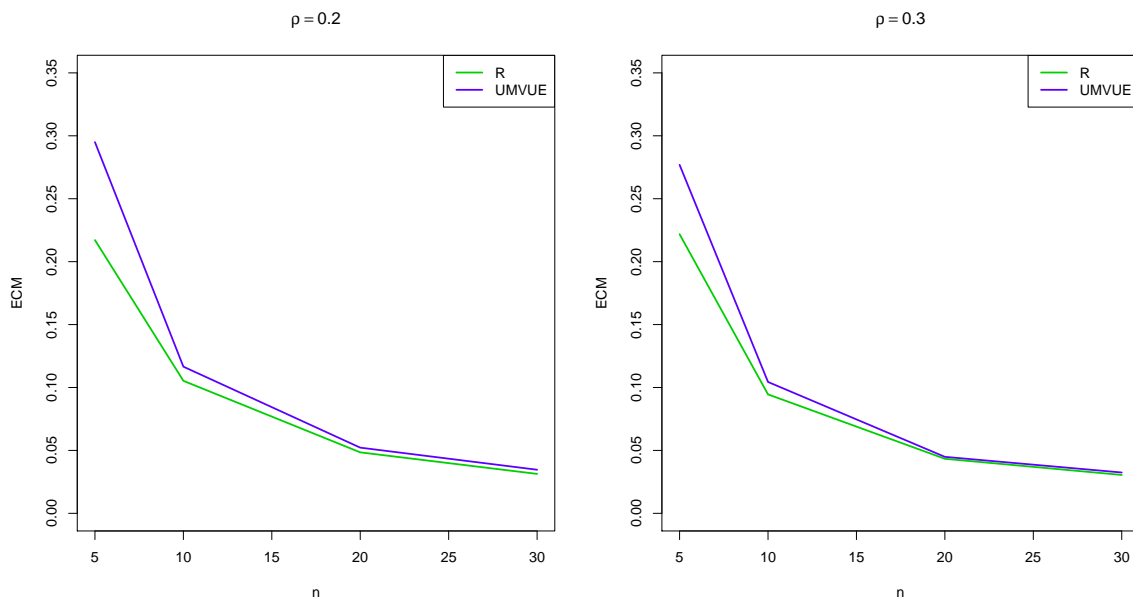


Figura 2-4: ECM para ambos estimadores con $\rho = 0.2$ y $\rho = 0.3$

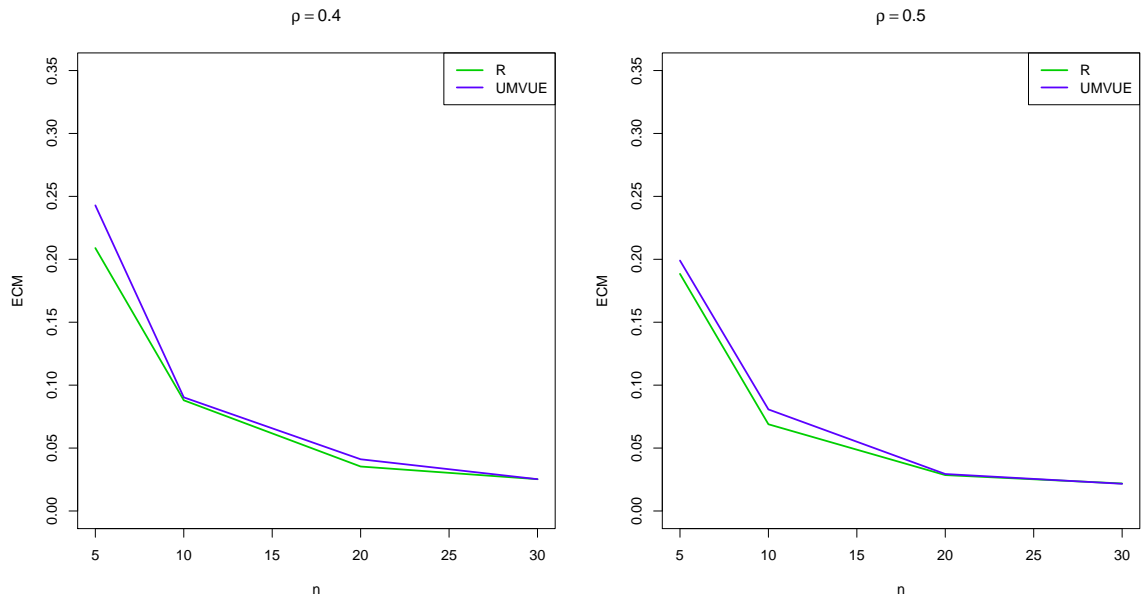


Figura 2-5: ECM para ambos estimadores con $\rho = 0.4$ y $\rho = 0.5$

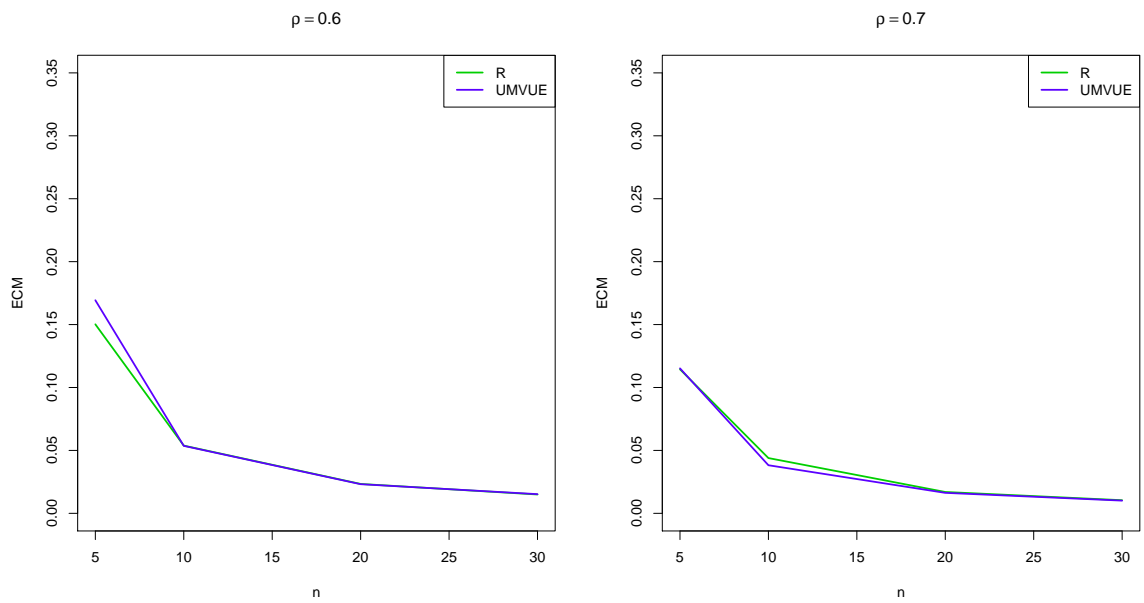


Figura 2-6: ECM para ambos estimadores con $\rho = 0.6$ y $\rho = 0.7$

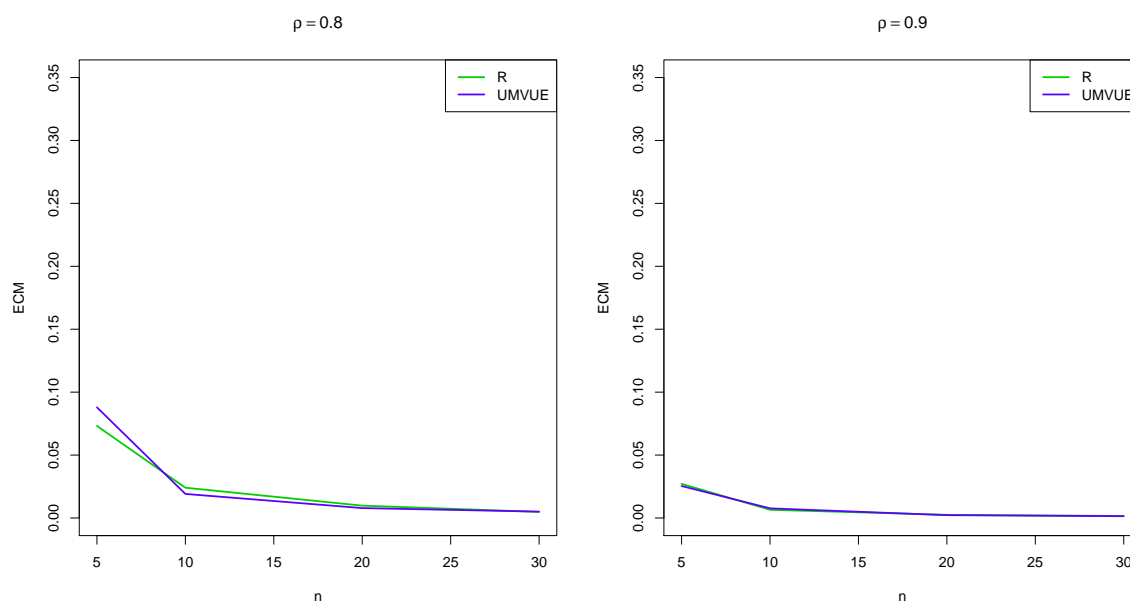


Figura 2-7: ECM para ambos estimadores con $\rho = 0.8$ y $\rho = 0.9$

De lo anterior se observa que a medida que los tamaños muestrales se hacen más grandes (dentro del rango que se estableció anteriormente), las diferencias entre los Errores Cuadráticos Medios desaparecen. Esto es mucho más evidente cuando se tiene un alto grado de relación lineal entre las variables (X, Y) de la Normal Bivariada que se consideró en la simulación. Indica esto que los dos estimadores para el coeficiente de correlación son parecidos cuando las muestras provenientes de la distribución tienen un tamaño cercano a 30.

En las figuras 1.3 a 1.5 se puede observar que para diferentes tamaños de muestra y siendo el valor verdadero de $\rho = (0, 0.1, 0.2, 0.3, 0.4, 0.5)$ el estimador máximo verosímil es, en la mayoría de las combinaciones, “mejor” en términos de menor error cuadrático medio, que el UMVUE salvo tamaños de muestra cercanos a 30, donde los resultados para ambos estimadores son muy similares. Pero si el valor verdadero de $\rho = (0.6, 0.7, 0.8, 0.9)$ la conclusión precedente ya no tiene igual validez. Las figuras 1.6. y 1.7. dan la evidencia que el UMVUE resulta ser en algunos casos “mejor” en términos de un menor E.C.M., mientras que en otros no hay diferencias grandes entre los dos estimadores analizados.

Entre las muchas características de R las más destacadas son las siguientes tres, Zheng and Matis (1994):

1. $|R| \leq 1$.
2. Si $|R| = 1$ entonces los pares $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ yacen en una línea recta.
3. Recíprocamente, si los $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ yacen en una línea recta, entonces $|R| = 1$.

Zheng y Matis demuestran por diferentes vías las anteriores tres propiedades.

2.3. Intervalos de confianza

2.3.1. Método I: Basado en la transformación Arco tangente

Este intervalo puede considerarse el intervalo clásico para este parámetro y fue propuesto por Fisher (1921). Debido a que la distribución del coeficiente de correlación no es centrada y/o simétrica, el cálculo de intervalos de confianza a partir de los cuantiles de la distribución no se hace sencillo. Por tanto, Fisher propone la transformación arcotangente hiperbólico:

$$r = \tanh(z) \Leftrightarrow z = \frac{1}{2} \log \frac{1+r}{1-r} \quad (2-8)$$

y demostró que z tiene una distribución aproximadamente Normal cuando el tamaño muestral es grande. Dicha distribución Normal se caracteriza por una media $\xi = \frac{1}{2} \log \left(\frac{1+\rho}{1-\rho} \right)$ y varianza $\frac{1}{n-3}$. El intervalo hallado a partir de la transformación Arcotangente hiperbólico tiene la siguiente forma:

$$\left(\tanh \left(\operatorname{arctanh}(r) - \frac{z_{\alpha/2}}{\sqrt{n-3}} \right), \tanh \left(\operatorname{arctanh}(r) + \frac{z_{\alpha/2}}{\sqrt{n-3}} \right) \right) \quad (2-9)$$

Modificaciones a la transformación Arco Tangente

Teniendo en cuenta el hecho de que la transformación propuesta por Fisher funciona adecuadamente siempre y cuando los tamaños muestrales sean grandes, se hizo necesario encontrar la manera de reducir el error al trabajar esta transformación en muestras pequeñas. Hotelling (1953) estudió esta situación y propuso 4 transformaciones z_i con $i = 1, \dots, 4$, para la transformación z original de Fisher, las cuales también, asintóticamente tienen una distribución Normal con media ξ_i y varianza $\frac{1}{n-1}$:

$$z_1 = z - \frac{7z+r}{8(n-1)}; \quad \xi_1 = \xi - \frac{7\xi+\rho}{8(n-1)} \quad (2-10)$$

$$z_2 = z - \frac{7z+r}{8(n-1)} - \frac{119z+57r+3r^2}{384(n-1)^2}; \quad \xi_2 = \xi - \frac{7\xi+\rho}{8(n-1)} - \frac{119\xi+57\rho+3\rho^2}{384(n-1)^2} \quad (2-11)$$

$$z_3 = z - \frac{3z+r}{4(n-1)}; \quad \xi_3 = \xi - \frac{3\xi+\rho}{4(n-1)} \quad (2-12)$$

$$z_4 = z - \frac{3z+r}{4(n-1)} - \frac{23z+33r-5r^2}{96(n-1)^2}; \quad \xi_4 = \xi - \frac{3\xi+\rho}{4(n-1)} - \frac{23\xi+33\rho-5\rho^2}{96(n-1)^2} \quad (2-13)$$

2.3.2. Método II: Intervalo de la Razón de Verosimilitud

Kalbfleish (1985) y Pawitan (2001) presentan la metodología para construir intervalos de verosimilitud. Un intervalo de la Razón de Verosimilitud para θ es definido como el conjunto de valores parametrales con valores altamente verosímiles:

$$\left\{ \theta, \frac{L(\theta)}{L(\hat{\theta})} > c \right\} \quad (2-14)$$

para un valor $c \in (0, 1)$.

Sabiendo que $2 \log \frac{L(\hat{\theta})}{L(\theta)} \xrightarrow{d} \chi_1^2$, la probabilidad de cubrimiento aproximada para este tipo de intervalos está dada por:

$$P \left(\frac{L(\theta)}{L(\hat{\theta})} > c \right) = P \left(2 \log \frac{L(\hat{\theta})}{L(\theta)} < -2 \log c \right) \quad (2-15)$$

$$\approx P \left(\chi_1^2 < -2 \log c \right). \quad (2-16)$$

Luego, para cualquier valor $0 < \alpha < 1$ el punto de corte c es:

$$c = \exp \left[-\frac{1}{2} \chi_{1, \alpha}^2 \right] \quad (2-17)$$

donde $\chi_{1, 1-\alpha}^2$ es el $100(1 - \alpha)$ percentil de una χ_1^2 . Por tanto:

$$P \left(\frac{L(\theta)}{L(\hat{\theta})} > c \right) = P \left(\chi_1^2 < \chi_{1, 1-\alpha}^2 \right) = 1 - \alpha. \quad (2-18)$$

Si $L(\rho)$ es la función de verosimilitud, se define la *función de verosimilitud relativa* como:

$$R(\rho) = \frac{L(\rho)}{L(r)} \quad (2-19)$$

El conjunto de valores de ρ para los cuales $R(\rho) > c$ es llamado *intervalo de $100c\%$ de verosimilitud* para ρ . Los intervalos del 14.7% y del 3.6% de verosimilitud corresponden a intervalos de confianza de niveles del 95% y 99% aproximadamente.

Lo que se debe hacer entonces es hallar las raíces que nos dan los límites del intervalo. Para el caso del parámetro ρ tenemos que un intervalo de confianza del 95% se halla encontrando el par de raíces tal que

$$R(\rho) = \frac{L(\rho)}{L(r)} \quad (2-20)$$

$$= \left(\frac{1 - \rho^2}{1 - r^2} \right)^{\frac{(n-1)}{2}} \frac{\int_0^\infty (\cosh w - \rho r)^{-(n-1)} dw}{\int_0^\infty (\cosh w - r^2)^{-(n-1)} dw} \geq K(k, \alpha) \quad (2-21)$$

donde $K(k, \alpha)$ es el valor crítico mínimo con el cual aseguramos una confianza deseada, ya sea del 95% o 99%, por ejemplo.

2.3.3. Método III: Bootstrap

La primera aplicación del método bootstrap fue en la determinación del intervalo de confianza del coeficiente de correlación en el artículo seminal de Efron (1979).

- A partir de la muestra $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ se estiman los parámetros de máxima verosimilitud del vector de medias y de la matriz de varianzas y covarianzas de la distribución normal bivariable.
- Se generan M muestras de tamaño n de una distribución normal bivariable con parámetros $\hat{\mu}$ y $\hat{\Sigma}$. Y para cada una de estas muestras se estima el parámetro ρ , por ejemplo, para la muestra j el valor del estimador para el coeficiente de correlación es r_j .
- Para los r_j , $j = 1, \dots, M$, se construye un histograma y se calculan los percentiles $0.025/(k-1)$ y $0.975/(k-1)$ los cuales se denotarán: $r_i^{\{0.025\}}$ y $r_i^{\{0.975\}}$.

2.3.4. Método IV: Intervalo de Jeyaratnam

Jeyaratnam propone un intervalo para el coeficiente de correlación de la distribución normal bivariada y este tiene la siguiente forma (Krishnamoorthy and Xia, 2007):

$$\left(\frac{r-w}{1-rw}, \frac{r+w}{1+rw} \right) \quad (2-22)$$

donde

$$w = \frac{\frac{t_{n-2, 1-\alpha/2}}{\sqrt{n-2}}}{\sqrt{1 + \frac{(t_{n-2, 1-\alpha/2})^2}{n-2}}} \quad (2-23)$$

y $t_{m,p}$ denota el p-ésimo cuantil de la distribución T-Student con m grados de libertad.

2.3.5. Método V: Test Generalizado para ρ

Krishnamoorthy and Xia (2007), proponen un algoritmo para construir un intervalo de confianza para ρ a partir de la distribución del Pivote Generalizado para el coeficiente de correlación:

$$G_{\rho_{ij}} = \frac{\sum_{k=1}^j b_{ik} b_{jk}}{\sqrt{\sum_{k=1}^i b_{ik}^2} \sqrt{\sum_{k=1}^j b_{jk}^2}} \quad (2-24)$$

para $i > j$. Que en el caso bivariado lo anterior se expresa de la siguiente forma:

$$G_{\rho_{21}} = \frac{b_{21}}{\sqrt{b_{21}^2 + b_{22}^2}} \quad (2-25)$$

Y simplificando la expresión anterior, se tiene:

$$G_{\rho_{21}} = \frac{r^*V_{22} - V_{21}}{\sqrt{(r^*V_{22} - V_{21})^2 + V_{11}^2}} \quad (2-26)$$

Donde $r^* = \frac{r}{\sqrt{1-r^2}}$ y V es una matriz triangular inferior, las V_{ij} 's son independientes con $V_{ii}^2 \sim \chi_{n-i}^2$ para $i = 1, \dots, p$ y $V_{ij} \sim N(0, 1)$ para $i < j$.

Entonces, según los autores, para un r dado la distribución de G_ρ no depende de algún parámetro que sea desconocido, y el intervalo se calcula empleando el siguiente algoritmo:

Algoritmo 1 Generar valores del Pivote $G_{\rho_{ij}}$

Requiere Un n y ρ fijo

Calcular: $r^* = \rho/\sqrt{1-\rho^2}$

Para $i = 1$ hasta m **Haga**

Generar: $Z_0 \sim N(0, 1)$.

Generar: $U1 \sim \chi_{n-1}^2$.

Generar: $U2 \sim \chi_{n-2}^2$.

Calcular: $Q_i = \frac{r^*\sqrt{U2} - Z_0}{\sqrt{(r^*\sqrt{U2} - Z_0)^2 + U1}}$

Fin para

Luego, los percentiles $\frac{\alpha}{2}$ y $(1 - \frac{\alpha}{2})$ de los valores calculados para el pivote $G_{\rho_{ij}}$ mediante el mencionado algoritmo conforman los límites del Intervalo de confianza al $100(1 - \alpha)\%$ para ρ .

3 Resultados de la Simulación - Metodología Clásica I

3.1. Caso Normal Bivariada

Para comparar los nueve métodos de construcción de intervalos de confianza en este caso se realizó una simulación en R en la cual se consideraron combinaciones de (ρ, n) con valores de $\rho = 0.0, 0.1, 0.2, \dots, 0.9$ y de $n = 5, 10, 20, 50, 100$. Para cada pareja se realizaron 1000 réplicas y se calcularon las fórmulas previas a un nivel de confianza del 95 % (Este es conocido como el nivel nominal). Para cada método y combinación se calculó la mediana de la longitud de los 1000 intervalos calculados y la proporción de intervalos que cubren el verdadero valor de ρ , esto es lo que se llama el nivel de confianza real.

Las siguientes tablas presentan los resultados (en esta parte se muestran solo dos a manera de ilustración de resultados brevemente, las restantes se encuentran en el anexo).

Se presenta además la siguiente serie de gráficas que permiten ilustrar de manera más clara la información obtenida y mostrada de manera tabular.

Nota: Para las gráficas 2.1 a 2.10, los puntos corresponden a cada tamaño muestral partiendo desde $n = 5$ hasta $n = 100$ en el sentido izquierda a derecha del gráfico.

$n = 5$										
ρ	Bootstrap	Arctanh	L.R	Jeyaratnam	Z1	Z2	Z3	Z4	P.G	
0.0	Longitud	1.68640	1.48640	1.67660	1.605600	1.62670	1.571800	1.588000	1.5498	
	Nivel	0.901	0.953	0.950	0.929	0.936	0.919	0.924	0.943	
0.1	Longitud	1.639230	1.69030	1.49050	1.68060	1.63110	1.577000	1.592800	1.5346	
	Nivel	0.913	0.947	0.933	0.946	0.938	0.926	0.931	0.944	
0.2	Longitud	1.622190	1.68488	1.48481	1.67501	1.62492	1.569920	1.586100	1.5360	
	Nivel	0.912	0.961	0.944	0.959	0.952	0.938	0.945	0.949	
0.3	Longitud	1.600380	1.66787	1.46717	1.65760	1.60561	1.548750	1.565460	1.5235	
	Nivel	0.892	0.939	0.914	0.934	0.922	0.904	0.908	0.942	
0.4	Longitud	1.584670	1.65975	1.45882	1.64929	1.59641	1.538680	1.555600	1.4994	
	Nivel	0.913	0.947	0.936	0.945	0.940	0.928	0.931	0.96	
0.5	Longitud	1.514100	1.60955	1.40818	1.59799	1.515710	1.477240	1.495570	1.43978	
	Nivel	0.894	0.955	0.929	0.954	0.929	0.937	0.918	0.942	
0.6	Longitud	1.296610	1.46645	1.27120	1.45224	1.353460	1.38212	1.329880	1.39341	
	Nivel	0.917	0.959	0.941	0.957	0.939	0.945	0.933	0.954	
0.7	Longitud	1.109360	1.33137	1.14926	1.31535	1.206140	1.23743	1.180640	1.24387	
	Nivel	0.876	0.935	0.914	0.933	0.917	0.920	0.901	0.913	0.959
0.8	Longitud	0.876584	1.15279	0.99496	1.13536	1.019500	1.05217	0.969860	0.993180	1.08383
	Nivel	0.897	0.951	0.927	0.949	0.936	0.941	0.925	0.931	0.952
0.9	Longitud	0.478751	0.75213	0.66126	0.73559	0.631446	0.65984	0.589663	0.609093	0.74498
	Nivel	0.896	0.955	0.917	0.951	0.933	0.939	0.924	0.929	0.946

Tabla 3-1: Longitud y nivel de confianza de los intervalos. Tamaño de muestra 5

$n = 10$									
ρ	Bootstrap	Arctanh	L.R	Jeyaratnam	Z1	Z2	Z3	Z4	P.G
0.0	Longitud	1.21010	1.13490	1.21460	1.18900	1.19280	1.17530	1.17820	1.1521
	Nivel	0.956	0.951	0.957	0.951	0.952	0.949	0.950	0.957
0.1	Longitud	1.21020	1.13500	1.21470	1.18910	1.19300	1.17540	1.17830	1.1456
	Nivel	0.947	0.941	0.949	0.943	0.944	0.939	0.940	0.947
0.2	Longitud	1.19560	1.12200	1.20020	1.17450	1.17840	1.16080	1.16370	1.1393
	Nivel	0.927	0.950	0.950	0.943	0.944	0.939	0.940	0.946
0.3	Longitud	1.16670	1.09680	1.17130	1.14600	1.14940	1.13180	1.13470	1.1105
	Nivel	0.931	0.925	0.934	0.926	0.926	0.921	0.922	0.957
0.4	Longitud	1.11610	1.05230	1.12070	1.09490	1.09880	1.08120	1.08410	1.0698
	Nivel	0.966	0.960	0.966	0.963	0.963	0.958	0.959	0.952
0.5	Longitud	1.03080	0.97700	1.03540	1.00970	1.01350	0.99610	0.99900	0.9997
	Nivel	0.935	0.946	0.956	0.945	0.945	0.944	0.944	0.959
0.6	Longitud	0.91380	0.87310	0.91830	0.89340	0.89710	0.88020	0.88300	0.8713
	Nivel	0.946	0.937	0.948	0.937	0.938	0.934	0.934	0.962
0.7	Longitud	0.70726	0.75815	0.76230	0.73920	0.74262	0.72710	0.72965	0.75789
	Nivel	0.928	0.946	0.947	0.941	0.942	0.939	0.940	0.944
0.8	Longitud	0.48596	0.54270	0.54612	0.52722	0.52999	0.51738	0.51944	0.58877
	Nivel	0.925	0.942	0.943	0.941	0.941	0.939	0.939	0.958
0.9	Longitud	0.27072	0.31574	0.31801	0.30556	0.30738	0.29914	0.30048	0.33409
	Nivel	0.942	0.960	0.960	0.957	0.957	0.951	0.953	0.945

Tabla 3-2: Longitud y nivel de confianza de los intervalos. Tamaño de muestra 10

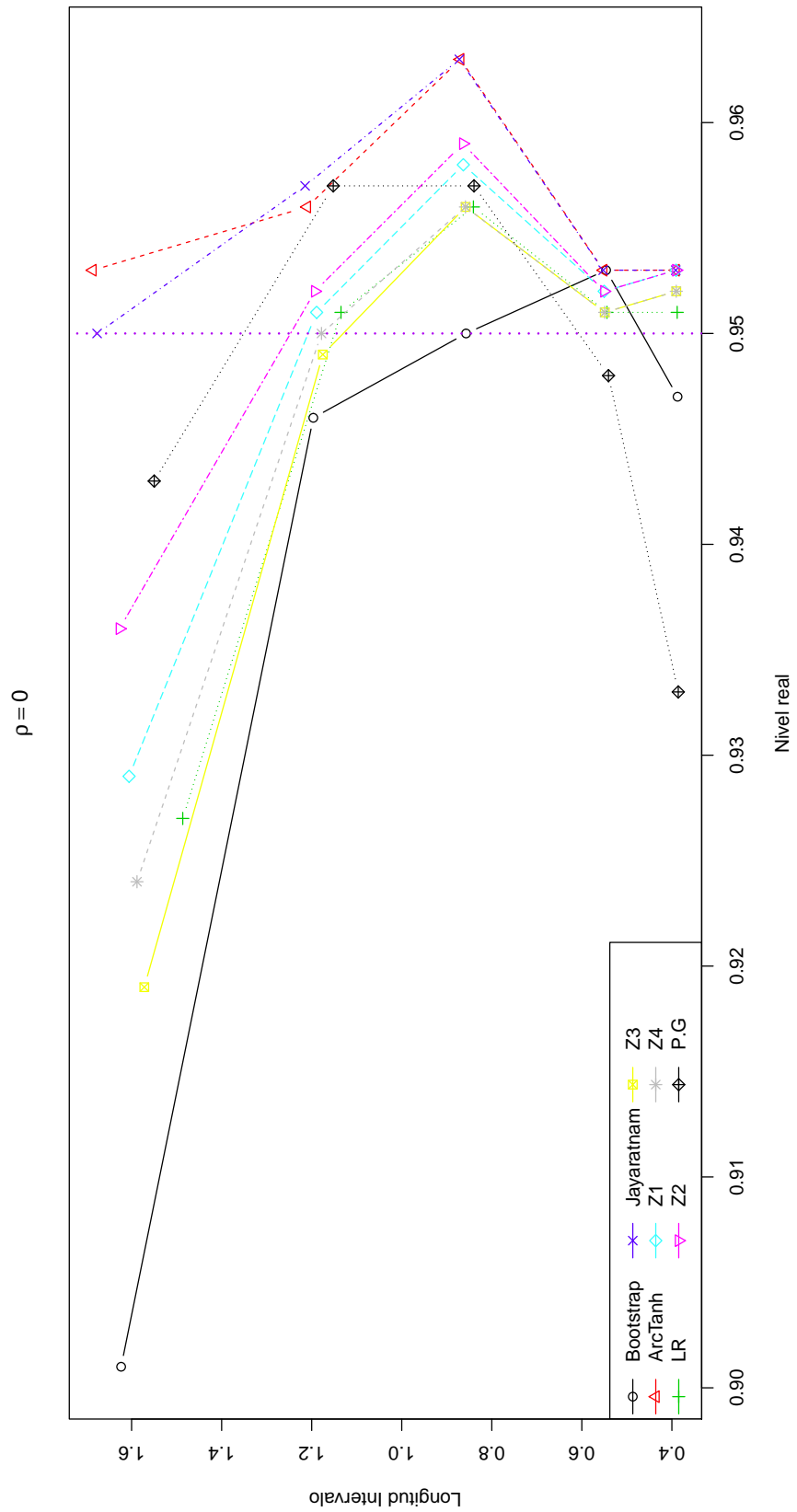


Figura 3-1: Amplitud y Nivel real alcanzado por cada intervalo con $\rho=0$

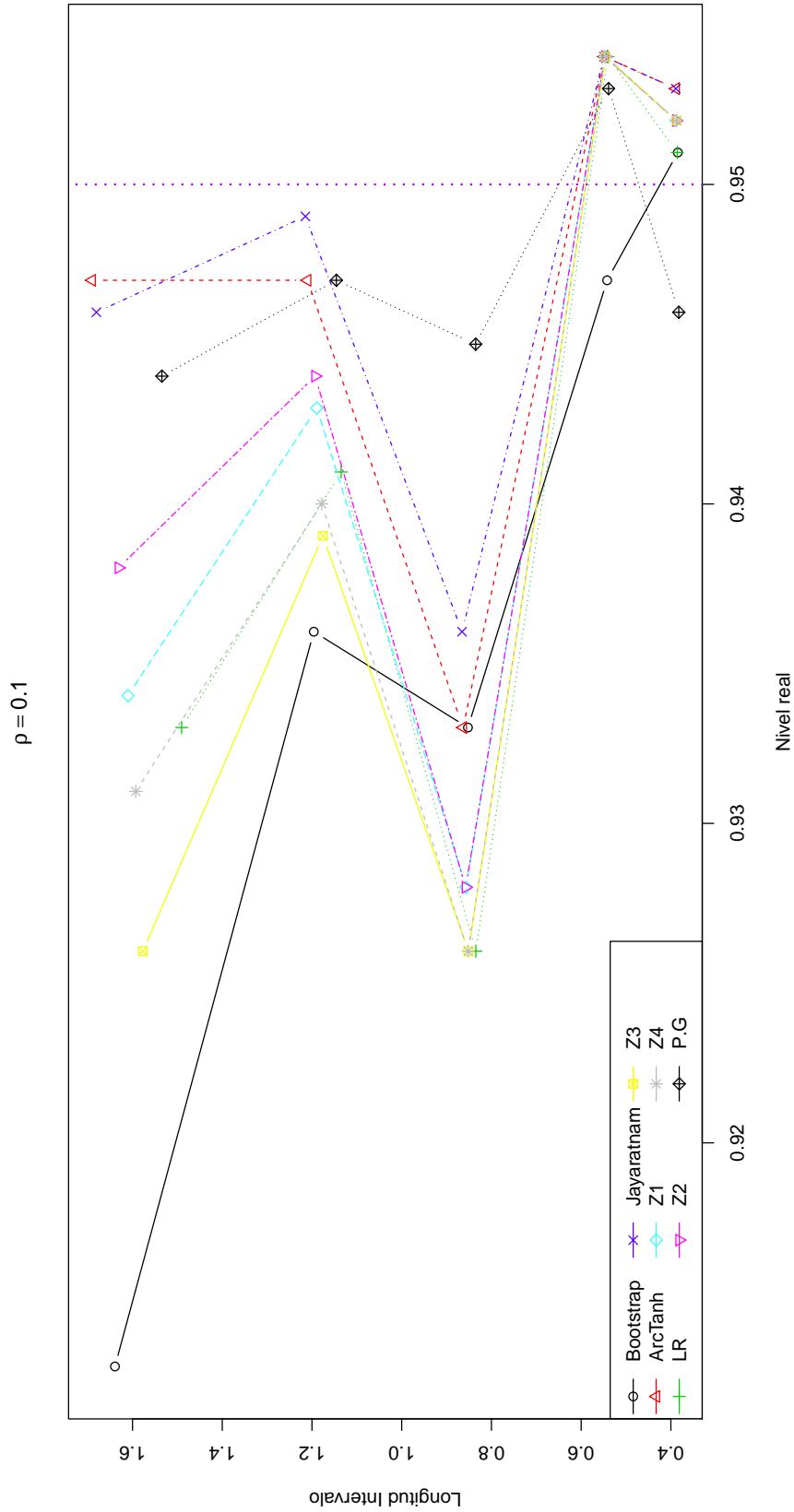


Figura 3-2: Amplitud y Nivel real alcanzado por cada intervalo con $\rho=0.1$

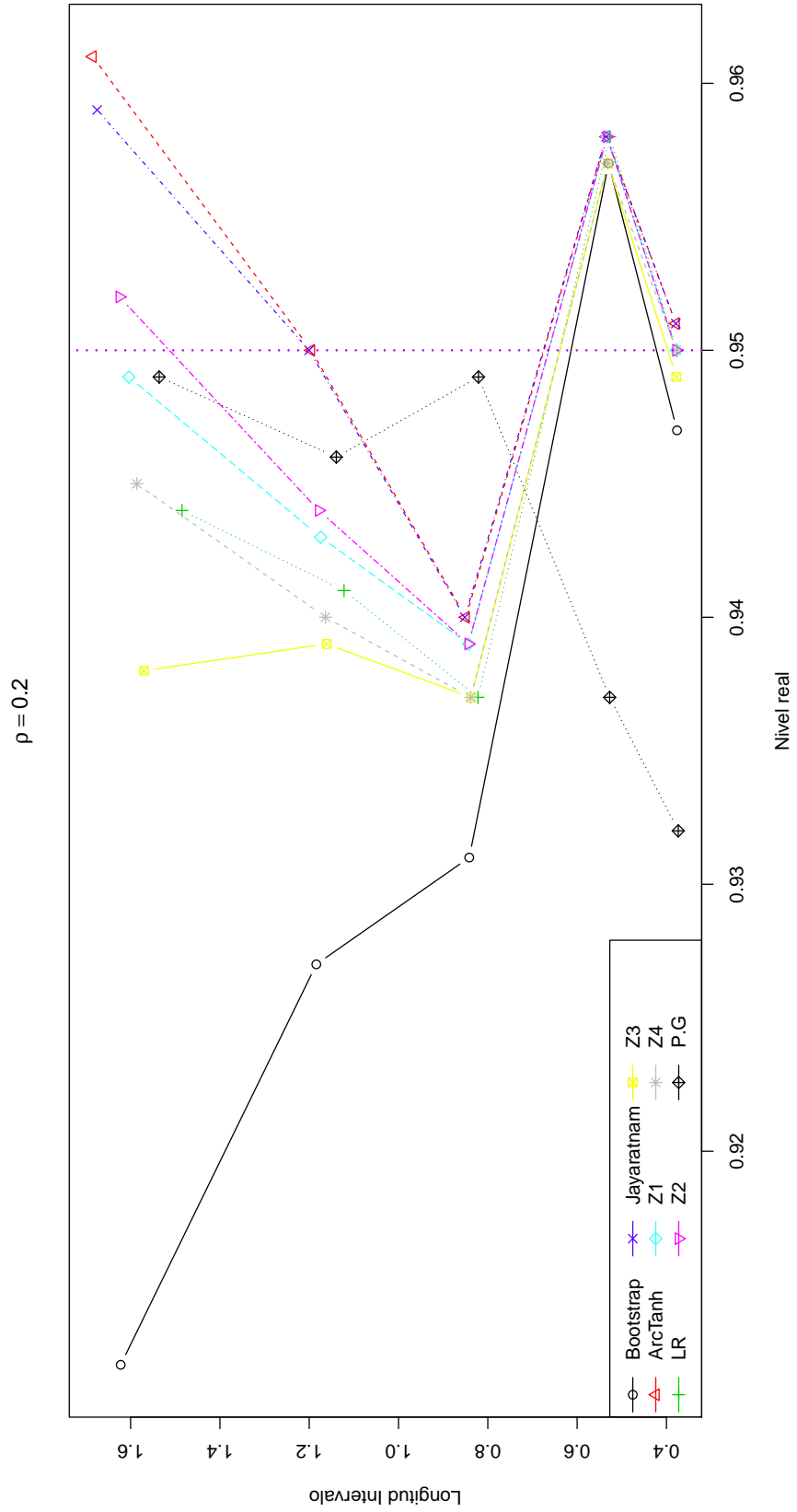


Figura 3-3: Amplitud y Nivel real alcanzado por cada intervalo con $\rho=0.2$

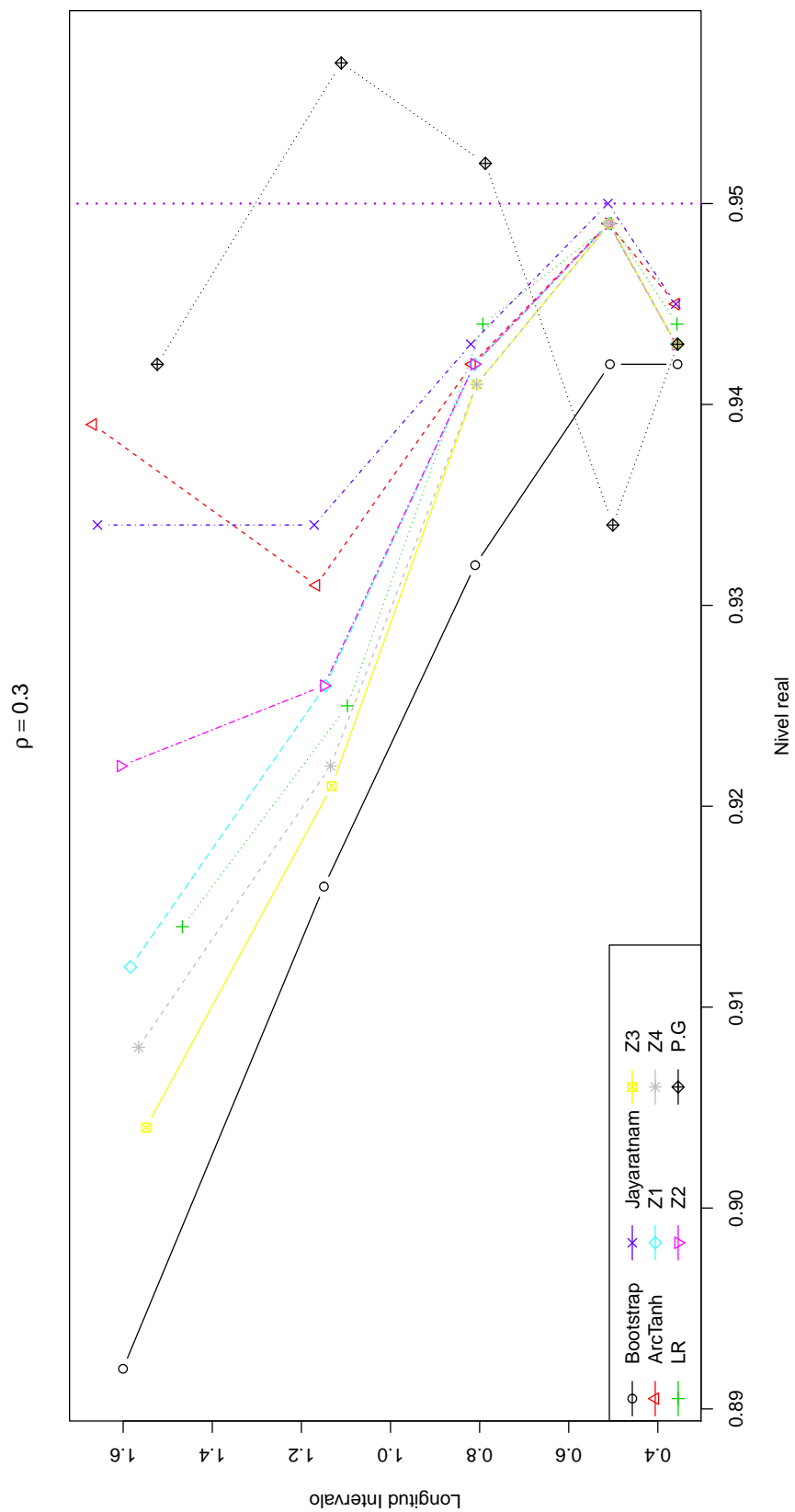


Figura 3-4: Amplitud y Nivel real alcanzado por cada intervalo con $\rho=0.3$

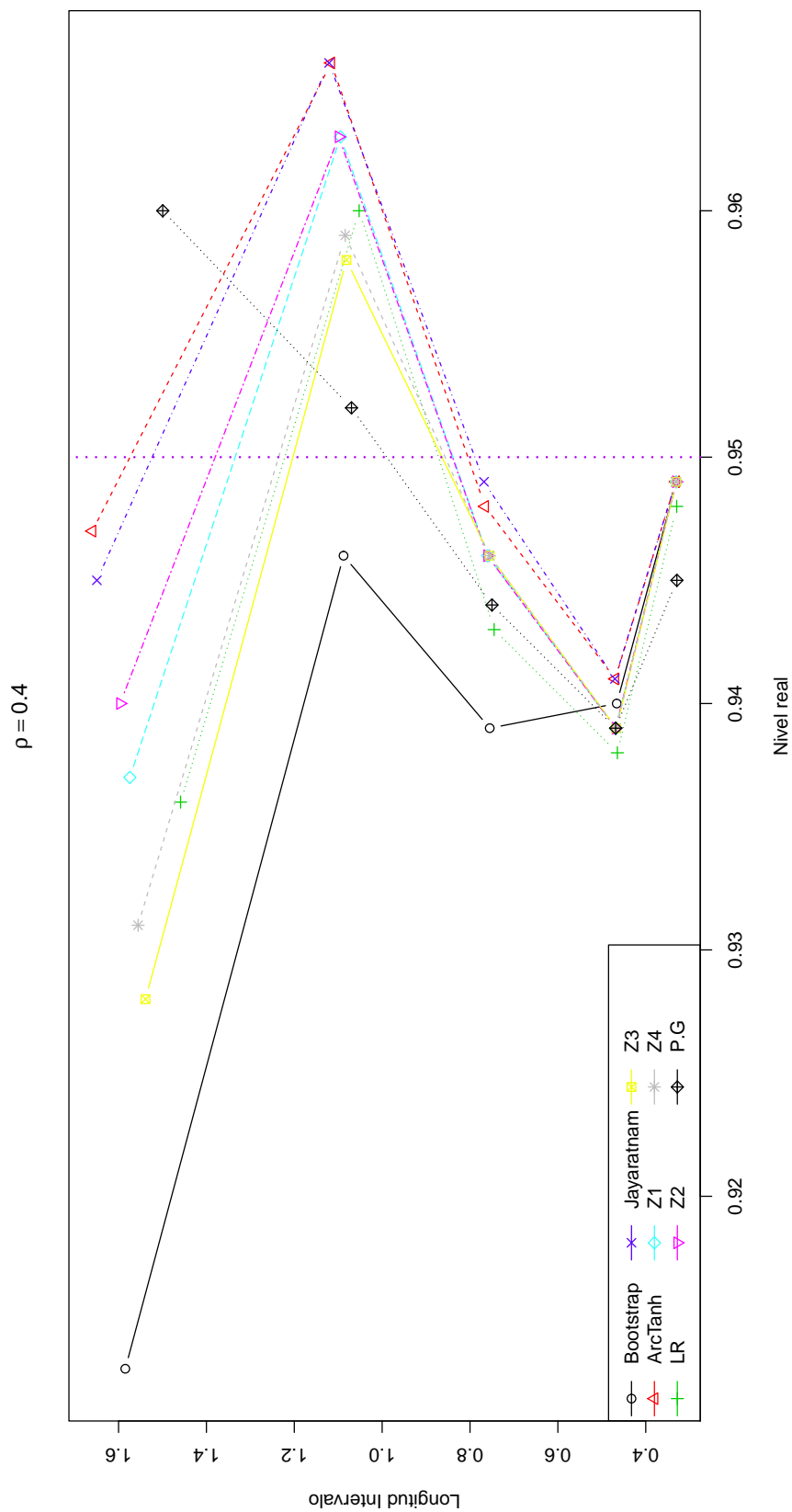


Figura 3-5: Amplitud y Nivel real alcanzado por cada intervalo con $\rho=0.4$

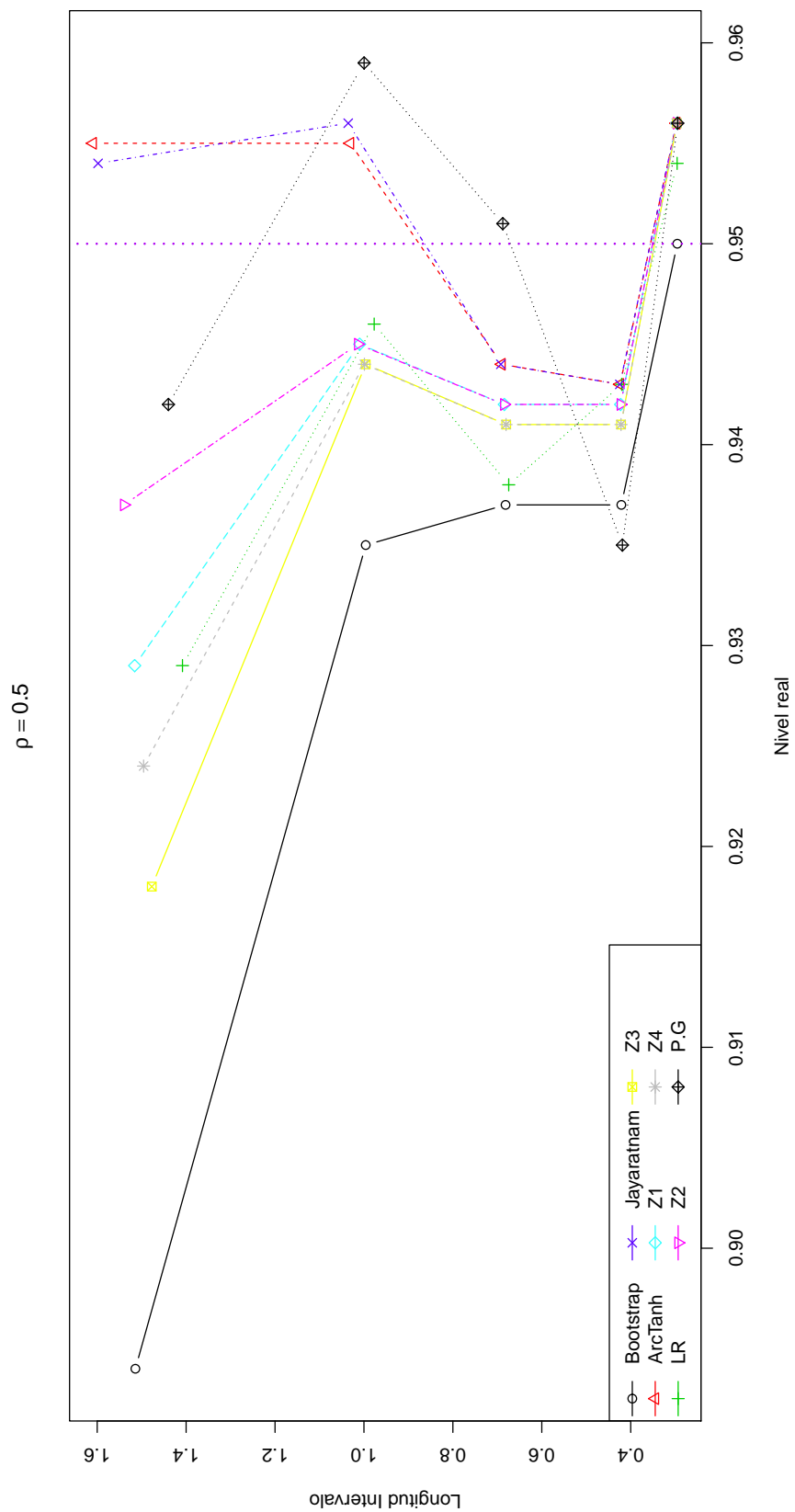


Figura 3-6: Amplitud y Nivel real alcanzado por cada intervalo con $\rho=0.5$

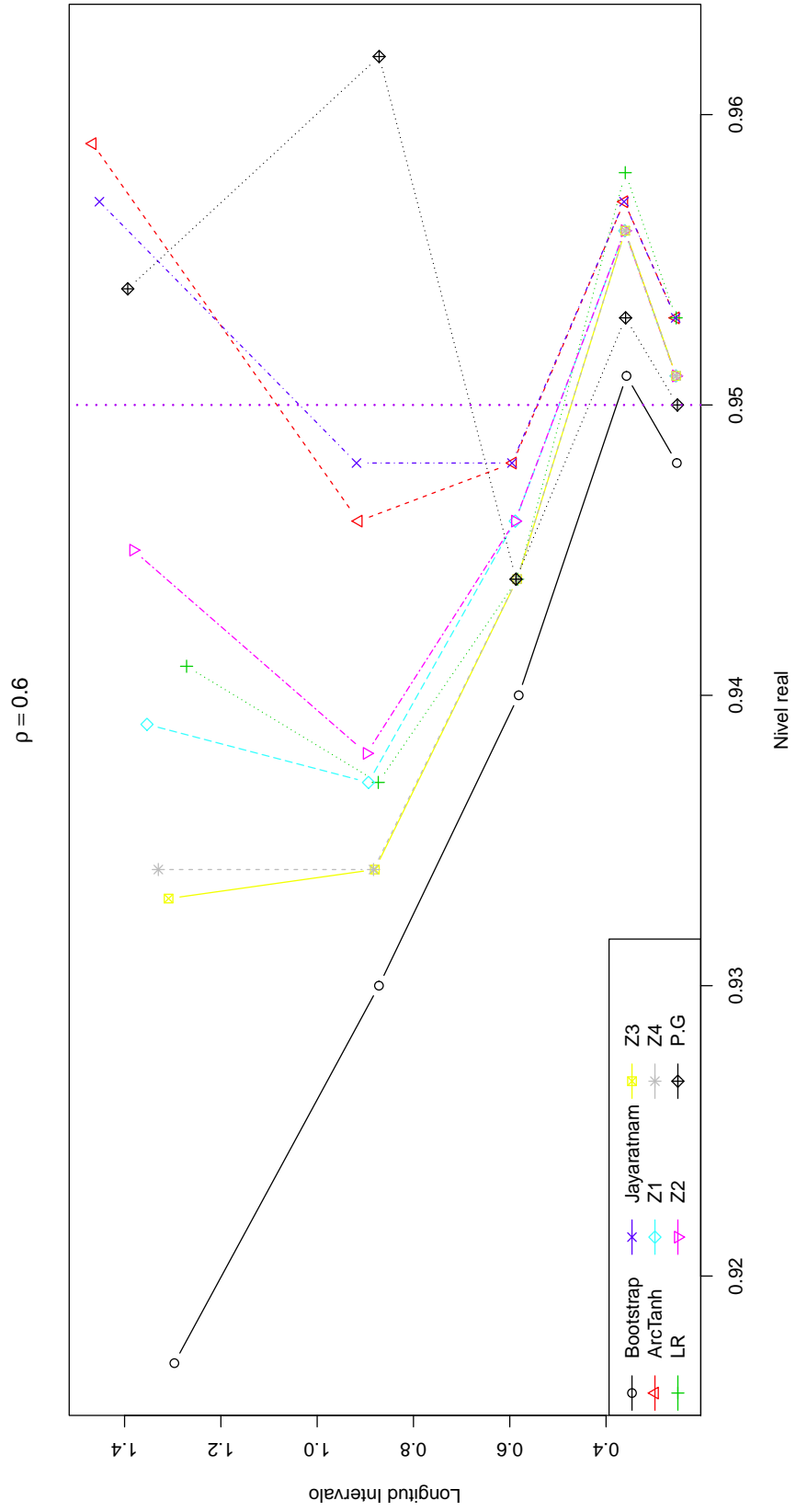


Figura 3-7: Amplitud y Nivel real alcanzado por cada intervalo con $\rho=0.6$

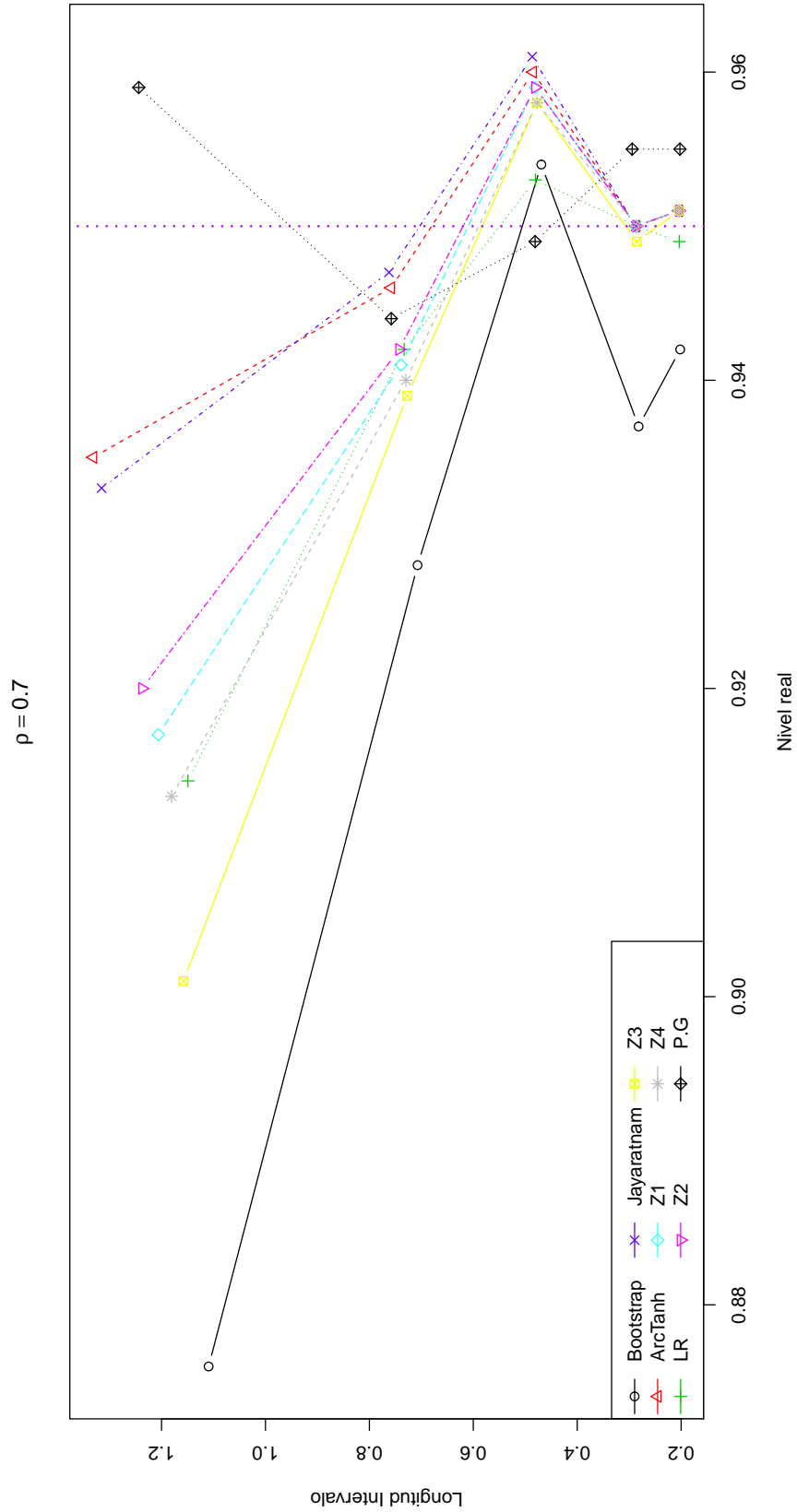


Figura 3-8: Amplitud y Nivel real alcanzado por cada intervalo con $\rho=0.7$

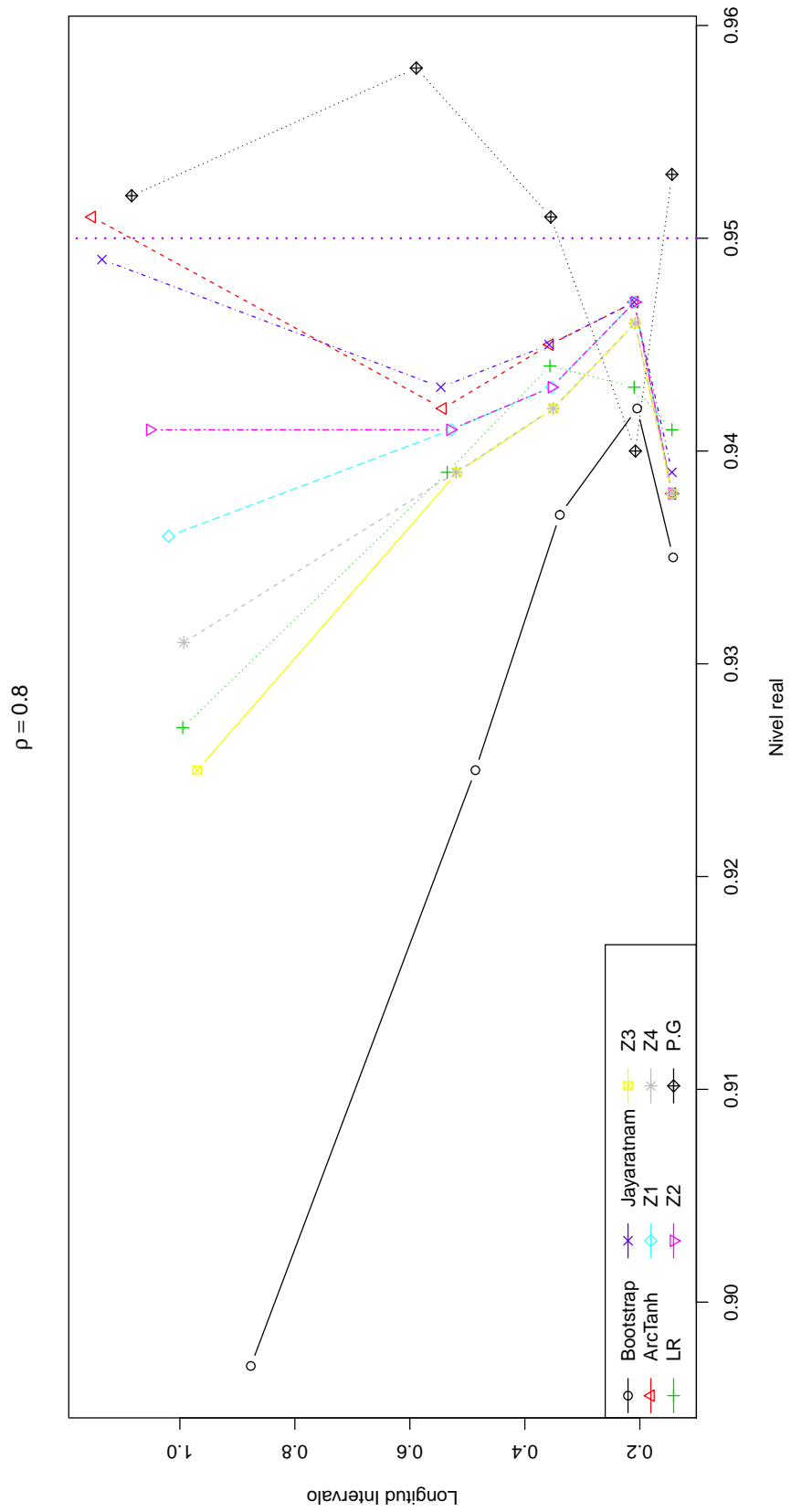


Figura 3-9: Amplitud y Nivel real alcanzado por cada intervalo con $\rho=0.8$

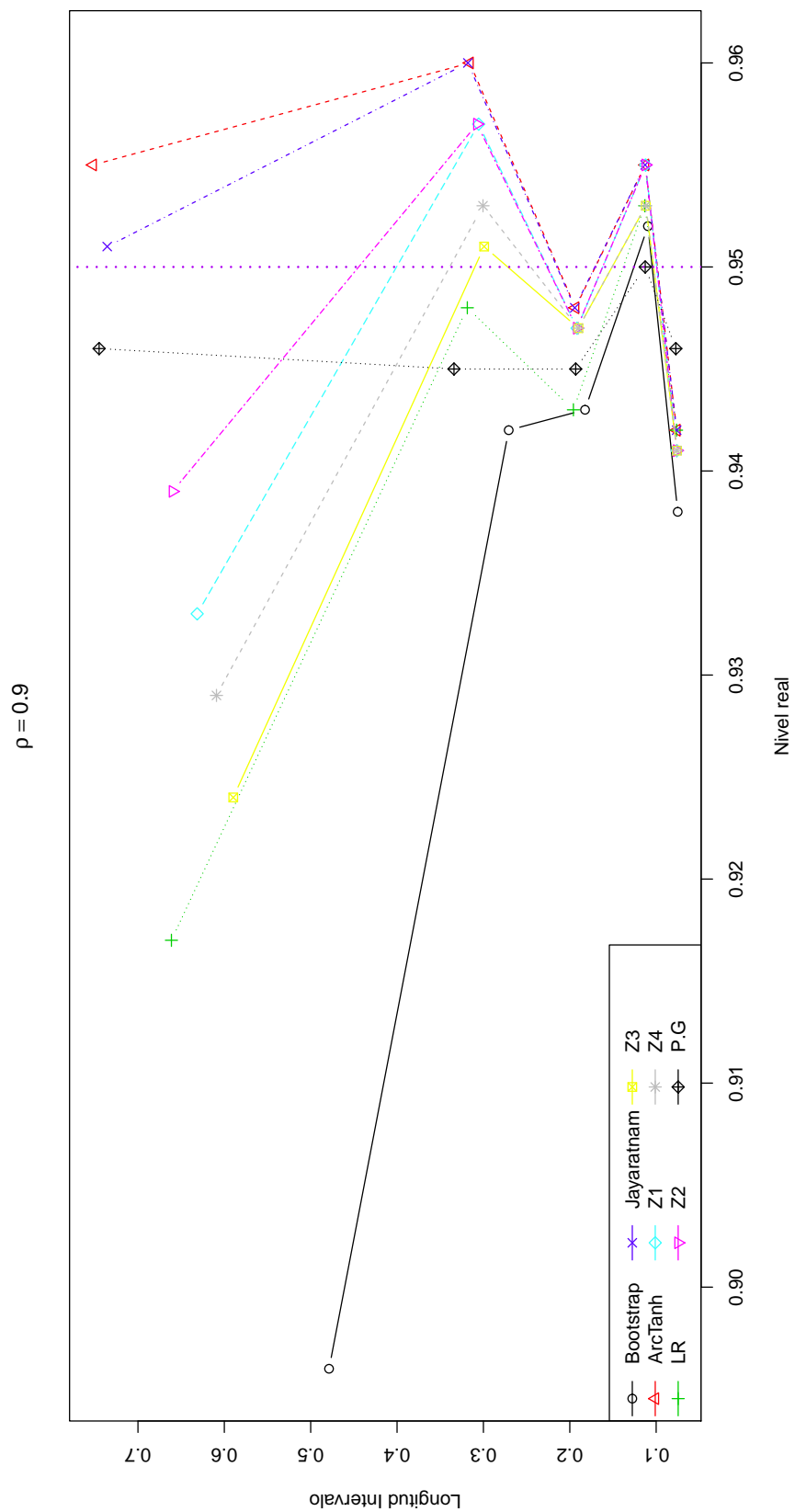


Figura 3-10: Amplitud y Nivel real alcanzado por cada intervalo con $\rho=0.9$

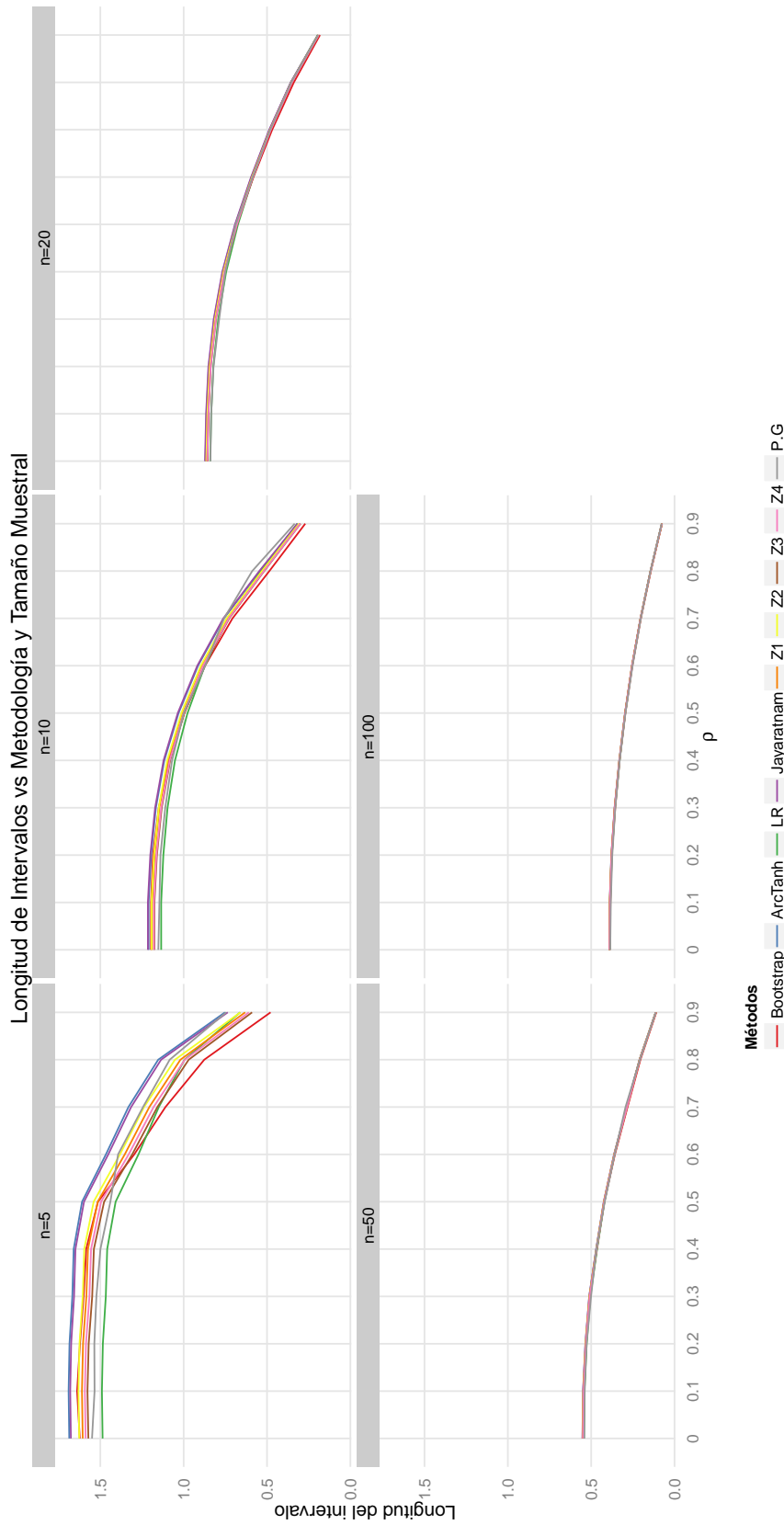


Figura 3-11: Amplitud por cada intervalo

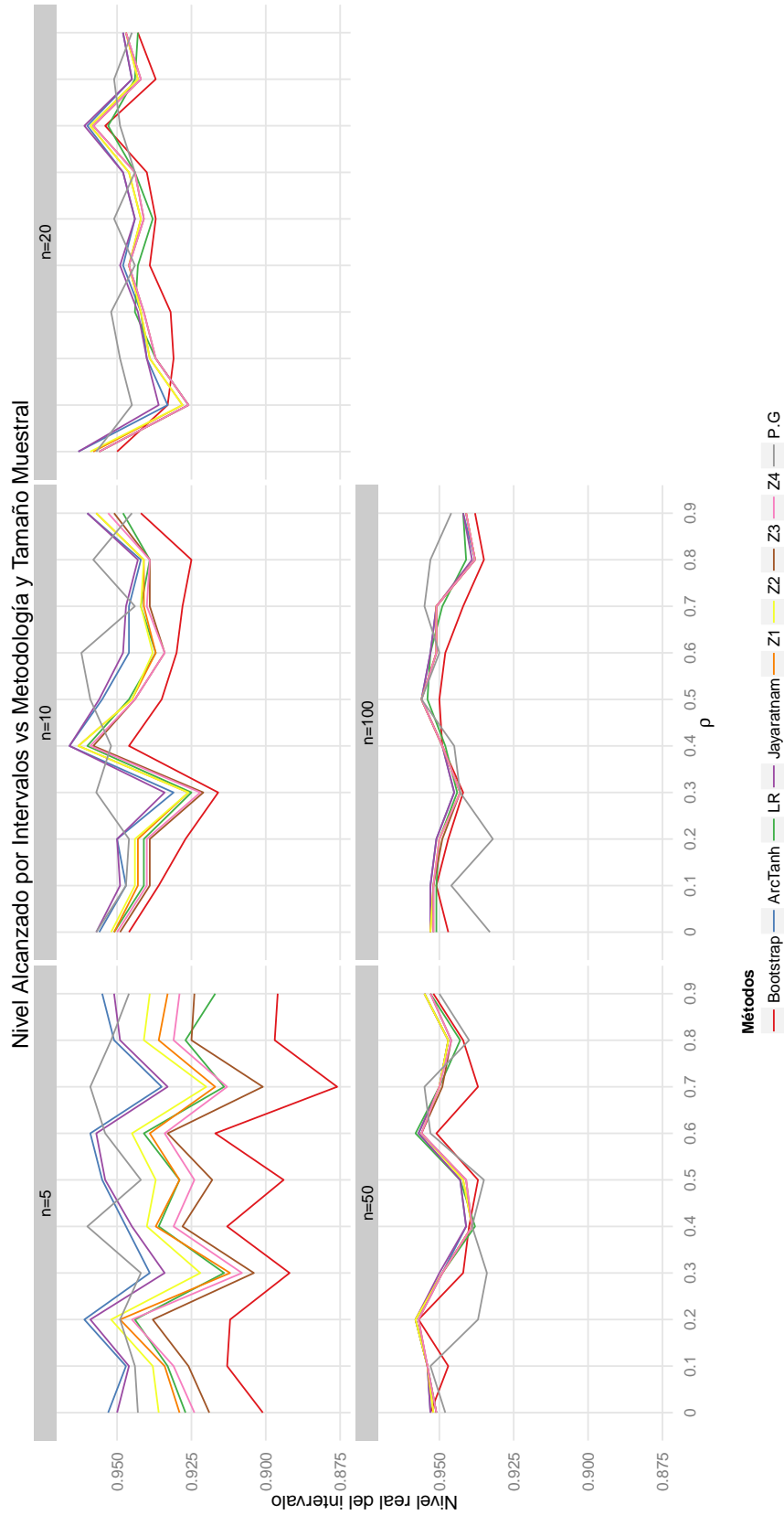


Figura 3-12: Nivel real alcanzado por cada intervalo

Observamos que en la gráfica para el caso $\rho = 0$ en el primer punto de cada método, es decir, cuando la muestra es de tamaño 5, se tienen las longitudes de intervalo más amplias, siendo el método LR el más corto en este caso. También se observa en que el método Bootstrap es el de menor probabilidad de cobertura, cercana al 90 %, y que de hecho se encuentra muy alejado de las coberturas de los demás métodos en el menor tamaño de muestra simulado en el estudio. Observamos que en los tamaños muestrales grandes, es decir, 50 y 100, los métodos se comportan de forma muy similar en cuanto a longitud de intervalo y nivel real; siendo este último superior al 95 % en todos los métodos excepto Bootstrap y Pivote Generalizado.

En términos generales, de las gráficas y las tablas observamos que las longitudes más amplias en los intervalos se encuentran en los tamaños de muestra más pequeños, siendo $n = 5$ el caso en donde los nueve tipos de intervalo alcanzan las mayores longitudes para el respectivo valor de ρ . Además, vemos que a medida que se amplía el valor de ρ en un tamaño de muestra particular, las longitudes son cada vez menores lo que sugiere que los intervalos de confianza alcanzan menores longitudes cuando el valor verdadero de ρ se va acercando a 1.

Con respecto al nivel real alcanzado por cada intervalo se observa que cada uno de los métodos cuando el tamaño muestral es bastante pequeño por ejemplo, $n = 5$, tienen una probabilidad de cobertura diferente a la deseada 95 %, y no es homogénea en cada valor de ρ , es decir, que en cada método se observa que algunas veces el nivel real supera al nominal y a veces es al contrario conforme se varía el valor de ρ . Los más cercanos a 95 % de nivel real, cuando los tamaños de muestra son pequeños (5, 10, 20) son ArcTanh, Jeyaratnam y P.G; y el que se comporta peor es Bootstrap.

A medida que aumenta el tamaño de muestra ($n \geq 50$) se nota una tendencia en todos los intervalos a estar cerca del nivel nominal deseado, 95 %. Se recomiendan para n muestrales grandes el método de Pivote Generalizado, L.R y Bootstrap por tener mejor comportamiento frente a los demás en todos los valores de ρ empleados en este análisis.

Índice de resúmenes A manera de resumen se presenta un índice que señala la calidad de las metodologías anteriormente mencionadas. Este índice busca favorecer a aquellos métodos que presenten longitudes de intervalo pequeños y niveles reales de confianza cercanos o mayores al 95 %:

$$I = (2 - LI) \frac{NR}{NN} \quad (3-1)$$

donde:

- LI: Mediana de la Longitud del intervalo.
- NR: Promedio del nivel real del intervalo.

- NN: Nivel nominal de los intervalos, que en este caso es 0.95.

Y el rango de este índice corresponde a (0, 2.1052) dado que si el Nivel Real, NR, se acerca al 100% y/o la longitud máxima es lo más pequeña posible, entonces I será cercano a 2.1052; o si la Longitud o el Nivel Real del intervalo es cercana a cero entonces el $I = 0$. Luego, a valores mayores del índice propuesto, mejor el intervalo obtenido.

Las siguientes tablas muestran el cálculo del índice de resumen a partir de los resultados ya obtenidos de la simulación:

$n = 5$									
ρ	Bootstrap	ArcTanh	LR	Jayaratnam	Z1	Z2	Z3	Z4	P.G
0	0.3572702	0.3145903	0.5011655	0.3234000	0.3856817	0.3677987	0.4142272	0.4007242	0.4468827
0.1	0.3467190	0.3087220	0.5003826	0.3180552	0.3832349	0.3642402	0.4123137	0.3990560	0.4624606
0.2	0.3626976	0.3187688	0.5119362	0.3280689	0.3957829	0.3758696	0.4246474	0.4117216	0.4635116
0.3	0.3752221	0.3282843	0.5126385	0.3366333	0.3996192	0.3827659	0.4294000	0.4153288	0.4724874
0.4	0.3991540	0.3391755	0.5332047	0.3488642	0.4199831	0.3993417	0.4506368	0.4355120	0.5058695
0.5	0.4572575	0.3925050	0.5787377	0.4037027	0.4735846	0.4537743	0.5051512	0.4906245	0.5555024
0.6	0.6789565	0.5386047	0.7218956	0.5517961	0.6390537	0.6146280	0.6790079	0.6588338	0.6091441
0.7	0.8212638	0.6580727	0.8185014	0.6723984	0.7662838	0.7384888	0.7987412	0.7874481	0.7632933
0.8	1.0607412	0.8481018	0.9807075	0.8637299	0.9660505	0.9388506	1.0030311	0.9866836	0.9180988
0.9	1.4347780	1.2544377	1.2922364	1.2657410	1.3440641	1.3246424	1.3717383	1.3601606	1.2497357

Tabla 3-3: Índice para un Tamaño muestral de 5

$n = 10$									
ρ	Bootstrap	ArcTanh	LR	Jayaratnam	Z1	Z2	Z3	Z4	P.G
0	0.8002164	0.7948888	0.8660106	0.7911872	0.8118537	0.8088994	0.8238319	0.8218000	0.8541477
0.1	0.7921516	0.7873059	0.8568053	0.7844734	0.8049249	0.8019032	0.8150520	0.8130505	0.8517019
0.2	0.7961466	0.8044000	0.8696821	0.7998000	0.8194174	0.8164109	0.8294829	0.8274968	0.8570760
0.3	0.8201575	0.8166340	0.8794316	0.8147429	0.8324253	0.8291112	0.8416971	0.8397964	0.8960542
0.4	0.9082596	0.8987867	0.9576758	0.8941093	0.9174856	0.9135322	0.9265373	0.9245769	0.9321583
0.5	0.9881474	0.9743011	1.0186926	0.9706922	0.9850879	0.9813079	0.9975596	0.9946779	1.0097765
0.6	1.1045463	1.0816265	1.1114793	1.0794227	1.0914571	1.0889686	1.1009402	1.0981874	1.1429573
0.7	1.2628029	1.2366212	1.2562215	1.2337915	1.2488556	1.2467915	1.2581612	1.2569779	1.2342651
0.8	1.4741968	1.4450280	1.4480961	1.4431672	1.4588273	1.4560836	1.4654528	1.4634167	1.4231140
0.9	1.7147176	1.7019891	1.6778003	1.6996952	1.7069253	1.7050919	1.7026504	1.7048869	1.6571421

Tabla 3-4: Índice para un Tamaño muestral de 10

$n = 20$									
ρ	Bootstrap	ArcTanh	LR	Jayaratnam	Z1	Z2	Z3	Z4	P.G
0	1.142900	1.145869	1.166521	1.143436	1.146373	1.146863	1.149213	1.148709	1.169051
0.1	1.127064	1.116555	1.135763	1.117781	1.116824	1.116140	1.119485	1.118998	1.158769
0.2	1.135232	1.138686	1.161781	1.136312	1.143702	1.143010	1.146296	1.145803	1.177859
0.3	1.167355	1.172740	1.199675	1.171603	1.178789	1.178194	1.182490	1.181995	1.215554
0.4	1.230485	1.231502	1.245653	1.230503	1.234879	1.234181	1.239658	1.239160	1.242105
0.5	1.300457	1.301726	1.308856	1.299640	1.304521	1.303926	1.307495	1.307099	1.313681
0.6	1.403667	1.403439	1.407355	1.401543	1.405457	1.404959	1.406560	1.406163	1.404076
0.7	1.537346	1.530947	1.524399	1.530924	1.533794	1.533290	1.535624	1.535321	1.517401
0.8	1.638142	1.634850	1.633259	1.633567	1.634705	1.634358	1.635649	1.635381	1.646832
0.9	1.804257	1.802697	1.790906	1.801958	1.802739	1.802530	1.804294	1.804135	1.797340

Tabla 3-5: Índice para un Tamaño muestral de 20

$n = 50$									
ρ	Bootstrap	ArcTanh	LR	Jayaratnam	Z1	Z2	Z3	Z4	P.G
0	1.458491	1.452673	1.457433	1.451971	1.452652	1.452552	1.452427	1.452427	1.456128
0.1	1.453296	1.457812	1.465545	1.457109	1.459319	1.459218	1.460624	1.460624	1.465714
0.2	1.481335	1.478043	1.485203	1.477337	1.479455	1.479354	1.479220	1.479220	1.452745
0.3	1.479932	1.487033	1.493426	1.488000	1.488432	1.488432	1.489730	1.489730	1.473754
0.4	1.517853	1.515901	1.516005	1.515307	1.513965	1.513965	1.515151	1.515052	1.513668
0.5	1.557294	1.564387	1.568060	1.563891	1.563918	1.563819	1.563348	1.563249	1.556332
0.6	1.644129	1.649767	1.653710	1.649364	1.649150	1.649050	1.650056	1.650056	1.645781
0.7	1.694589	1.713100	1.713800	1.712700	1.713900	1.713900	1.712895	1.714600	1.714778
0.8	1.780192	1.784726	1.776969	1.784447	1.785334	1.785304	1.784056	1.783957	1.773444
0.9	1.894450	1.897334	1.892728	1.897143	1.897655	1.897635	1.893992	1.893982	1.887330

Tabla 3-6: Índice para un Tamaño muestral de 50

$n = 100$									
ρ	Bootstrap	ArcTanh	LR	Jayaratnam	Z1	Z2	Z3	Z4	P.G
0	1.607508	1.613880	1.613497	1.613680	1.614382	1.614382	1.613189	1.613189	1.585314
0.1	1.617100	1.616388	1.615899	1.616087	1.615193	1.615193	1.615694	1.615694	1.610789
0.2	1.618573	1.622907	1.624000	1.622706	1.621700	1.621700	1.620492	1.622200	1.595682
0.3	1.630354	1.631866	1.632524	1.631667	1.628908	1.628908	1.629305	1.629305	1.632482
0.4	1.668542	1.666844	1.666983	1.666644	1.667343	1.667243	1.667743	1.667743	1.662106
0.5	1.704800	1.713756	1.711576	1.713555	1.714158	1.714158	1.714561	1.714561	1.716473
0.6	1.743322	1.750310	1.751112	1.750109	1.746937	1.746937	1.747237	1.747237	1.748200
0.7	1.783057	1.797790	1.794409	1.797690	1.798091	1.798091	1.798291	1.798291	1.807162
0.8	1.828624	1.832526	1.838318	1.834381	1.832724	1.832714	1.832901	1.832901	1.861680
0.9	1.900655	1.907619	1.907411	1.907570	1.905703	1.905703	1.905802	1.905802	1.914684

Tabla 3-7: Índice para un Tamaño muestral de 100

De las anteriores tablas se observa que, según el criterio establecido para concluir con el Índice de resumen propuesto, en tamaños de muestra pequeños y para correlaciones menores que 0.7 (dentro de las empleadas en este estudio), el mejor método para la construcción de intervalos de confianza para el coeficiente de correlación agrupando las características deseadas (longitud del intervalo corta y mayor porcentaje de cobertura) es el de la Razón de verosimilitud. Le sigue el método Bootstrap en calidad. El método menos eficiente es el de la Transformación de Fisher. Esto es apenas lógico de esperar ya que establece que z tiene una distribución aproximadamente Normal cuando el tamaño muestral es grande.

Este comportamiento cambia cuando se comienza a aumentar el tamaño de muestra ya que, cuando este es igual a 20, en correlaciones pequeñas (de 0 a 0.4) el mejor método para construir los intervalos es el que parte del Pivote Generalizado. Se observa que en el resto de las correlaciones algunos de los métodos también muy eficientes son el de la Razón de verosimilitud, Bootstrap y la Transformación de Fisher modificada Z3. En este caso resulta complicado afirmar cuál método es el menos efectivo ya que las diferencias en el índice de resumen para el resto de los métodos son mínimas y por tanto no muy significativas.

El método de la Razón de Verosimilitud en tamaños de muestra iguales a 50 muestra los mejores resultados para estos intervalos de confianza y en menor medida el del Pivote Generalizado y la Transformación de Fisher modificada Z1. Cabe aclarar que, las diferencias que hicieron que cada uno de estos métodos sobresaliese en tamaños de muestra grandes, es decir, iguales a 100, fueron mínimas. Lo anterior indica que con tamaños de muestra grandes, todos los métodos ofrecen índices de resumen casi iguales y por tanto, el uso de un método u otro no arrojaría resultados significativamente diferentes.

Conforme al comportamiento observado en todos los tamaños muestrales se sugiere emplear en primera instancia el método de Razón de Verosimilitud, y en segunda, el Pivote Generalizado ya que mostraron buenos comportamientos en muchos de los estadios de la simulación.

Clusters Para el análisis de los resultados obtenidos en las tablas 2.6 a 2.10 las cuales contienen la información del Índice de resumen propuesto en 3-1 para las nueve diferentes metodologías utilizadas, valores verdaderos de $\rho = 0.0, 0.1, \dots, 0.8, 0.9$, y tamaños muestrales $n = (5, 10, 20, 50, 100)$, se utilizó un método de clústering jerárquico.

En dicho método se calcula una matriz de distancias euclidianas a partir de una matriz de información en donde se tiene como individuos a cada uno de los métodos de construcción del intervalo (filas de la matriz), y se tienen también observaciones en cada una de las 10 variables (resultados del índice en cada valor de ρ simulado). Luego, si en particular se toman dos métodos de construcción del intervalo de confianza para el coeficiente de correlación ρ , $A = (A_1, A_2, \dots, A_k)$ y $B = (B_1, B_2, \dots, B_k)$ se calcula la distancia entre los dos de la siguiente forma:

$$d_E(A, B) = \sqrt{(A_1 - B_1)^2 + (A_2 - B_2)^2 + \dots + (A_k - B_k)^2} \quad (3-2)$$

donde k es igual a 10 (10 valores de ρ empleados en este estudio de simulación) y cada A_i o B_i , con $i = 1, 2, \dots, 10$, representa el valor del índice en el respectivo método y valor de ρ .

La matriz de distancias obtenida en la parte anterior refleja las similitudes entre los nueve métodos considerados en este estudio y el método de clústering jerárquico posteriormente agrupa paso a paso los métodos mediante un agrupamiento Completo.

El agrupamiento de los nueve métodos de construcción de los intervalos de confianza se hizo dependiendo de los tamaños muestrales debido a que la ecuación del Índice de resumen depende de n y de los valores verdaderos de ρ . Cabe recordar que dicho Índice resume ambos criterios de decisión sobre el desempeño de cada metodología, nivel real y longitud de intervalo. Por otra parte, según los resultados de la simulación, las metodologías tienen un comportamiento diferente dependiendo del tamaño muestral.

Por último, se presenta la agrupación de los métodos mediante un análisis de clusters con el fin de determinar qué métodos pueden ser empleados en cada tamaño muestral y obtener resultados aproximadamente iguales en términos de amplitud de intervalo y nivel de cobertura real:

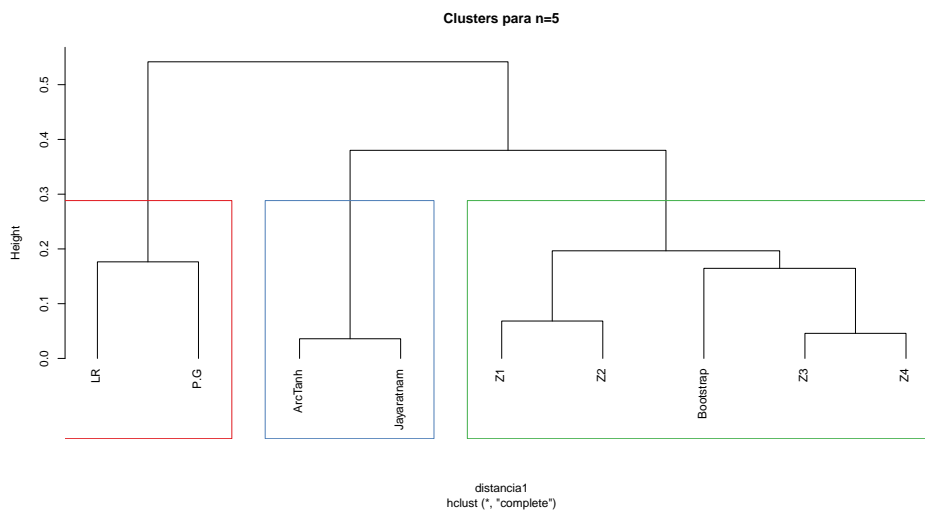


Figura 3-13: Dendogramas para el Análisis de clusters n=5

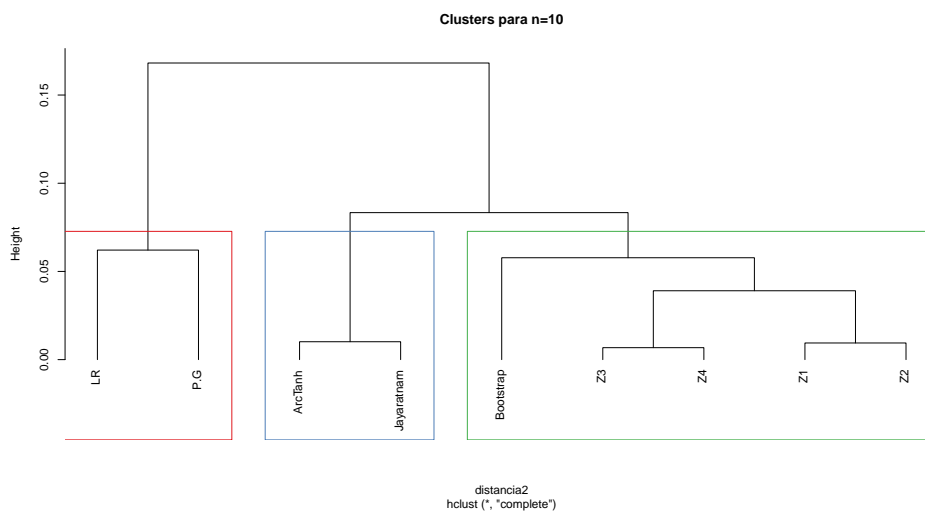


Figura 3-14: Dendogramas para el Análisis de clusters n=10

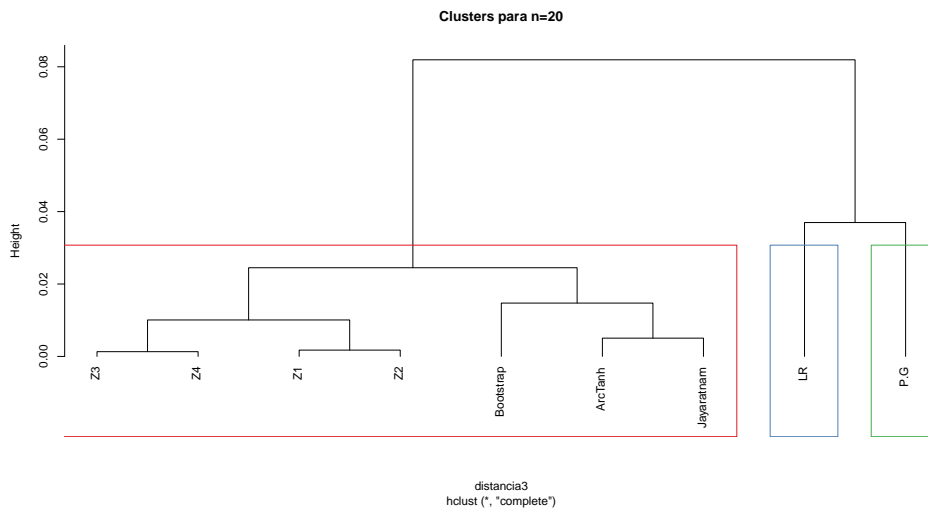


Figura 3-15: Dendogramas para el Análisis de clusters n=20

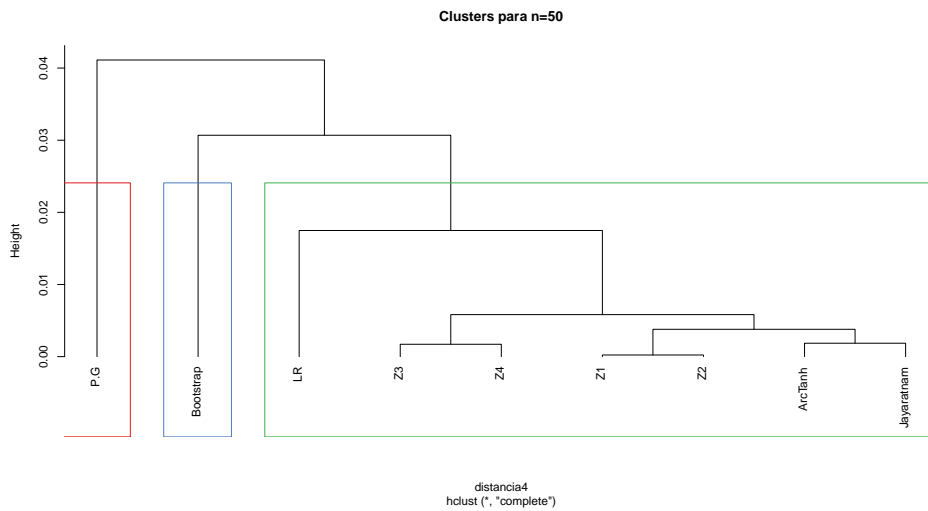


Figura 3-16: Dendogramas para el Análisis de clusters n=50

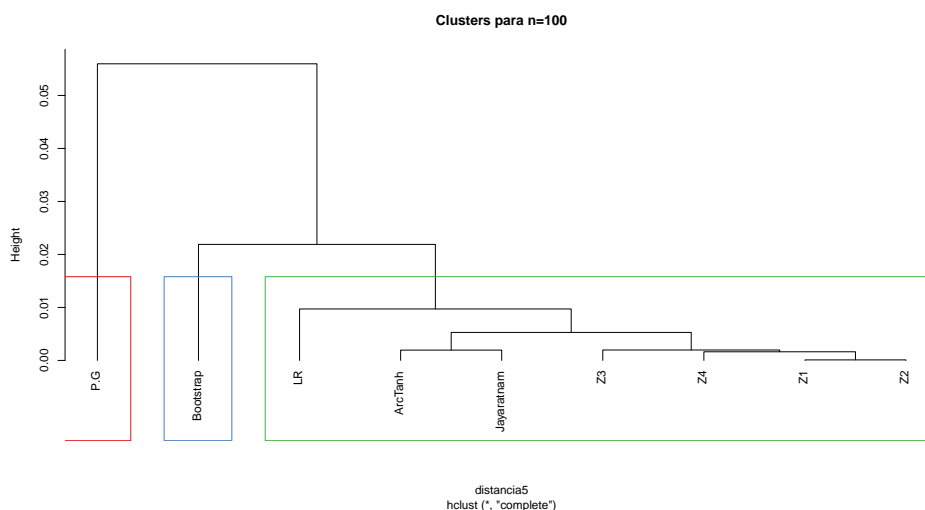


Figura 3-17: Dendogramas para el Análisis de clusters $n=100$

Los anteriores dendogramas y los resultados arrojados por el Índice de resumen que aportan la información para la construcción de los mismos, nos permiten afirmar que, en el caso de muestras pequeñas, es decir cuando $n=5$, el mejor método es el de la Razón de Verosimilitud y que el método que se comporta de manera similar a este es el del Pivote Generalizado. El dendograma nos muestra que estos dos métodos no se parecen en su calidad a la de los demás, y la tabla nos confirma que para estos dos es en donde el índice de resumen muestra los mejores resultados, es decir, los más altos en la gran mayoría de los valores de ρ simulados.

El otro método que está muy cercano a obtener valores similares al primero mencionado es el de Bootstrap debido al grado de separación que se observa entre el cluster L.R y P.G y Bootstrap. Los dendogramas, cuando el tamaño de muestra aumenta, muestran un agrupamiento entre los métodos de Razón de verosimilitud, Pivote Generalizado y Bootstrap, los cuales nos permiten afirmar que el resto de los métodos arrojan resultados considerablemente distintos a los anteriormente mencionados.

Se observa de las tablas que a medida que se aumenta el tamaño de la muestra, los métodos de Razón de Verosimilitud y Pivote Generalizado continúan proporcionando resultados (Longitud de intervalo y Porcentaje de cobertura) muy parecidos entre sí. Y en términos generales, las diferencias entre todos los métodos se van haciendo mínimas a medida que n aumenta, observándose en particular resultados muy parejos para los métodos de Arco tangente hiperbólico, Jeyaratnam, Z1, Z2, Z3 y Z4.

4 Metodología Bayesiana para la construcción de Intervalos de Credibilidad

4.1. Inferencia Bayesiana

El análisis estadístico bayesiano se fundamenta sobre el teorema de Bayes, el cual muestra la manera como se puede implementar una creencia apriori sobre algún modelo, y actualizarla a partir de la información disponible, es decir, los datos, para así obtener una creencia aposteriori sobre dicho modelo. Por ejemplo, si θ es algún objeto de interés, pero está sujeto a la incertidumbre — un parámetro, una hipótesis, un modelo, un punto muestral — entonces el Teorema de Bayes indica como ajustar racionalmente las creencias apriori sobre θ , $p(\theta)$, con respecto a la información o datos disponibles, y , para obtener una creencia aposteriori de $p(\theta|y)$, (Jackman, 2009).

Definición 1. (Probabilidad condicional) Sean A y B eventos con $P(B) > 0$. Entonces la probabilidad condicional de A dado B es:

$$P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{P(A, B)}{P(B)}$$

Proposición 1. (Teorema de Bayes) Para cualquier partición finita $\{E_j, j \in J\}$ de Ω , donde Ω es el espacio muestral asociado a un experimento aleatorio, y $P(G) > 0$

$$P(E_j|G) = \frac{P(G|E_j)P(E_j)}{\sum_{j \in J} P(G|E_j)P(E_j)}$$

Proposición 2. (Teorema de Bayes, caso continuo) Sean Y y θ variables aleatorias con función de distribución dadas por $p(y|\theta)$ y $p(\theta)$. Se define la distribución posterior de θ dado y como:

$$p(\theta|y) = \frac{p(y|\theta)p(\theta)}{\int p(y|\theta)p(\theta) d\theta}$$

4.2. Intervalos Bayesianos o Intervalos de credibilidad

Mediante la metodología bayesiana, es posible construir Intervalos de Confianza, que propiamente se llamarían Intervalos de Credibilidad, para parámetros desconocidos de una función de probabilidad correctamente definida para la v.a X , $f(x|\theta)$, donde $\theta \subseteq \Theta$. Específicamente, se desea conocer la región C en donde sea más probable o plausible encontrar al parámetro de interés θ . Para ello se define lo siguiente, (Jackman, 2009; Bernardo and Smith, 2000):

Definición 2. (*Región de Credibilidad*). Una región $C \subseteq \Omega$, donde Ω es el espacio muestral asociado a un experimento aleatorio, tal que $\int_C p(\theta)d\theta = 1 - \alpha$, $0 \leq \alpha \leq 1$, es una región de $100(1 - \alpha)\%$ de credibilidad para θ .

Para los casos en donde sea un solo parámetro (i.e. $\Omega \subseteq \mathfrak{R}$), si C no es una unión de intervalos disjuntos, entonces C es un intervalo de credibilidad.

Si $p(\theta)$ es una densidad (apriori/aposteriori), entonces C es una región de credibilidad (apriori/aposteriori).

Debido a que pueden resultar múltiples regiones de credibilidad para un parámetro desconocido θ se suele limitar o restringir el espacio solamente a aquellas regiones que cumplan con condiciones particulares tales como: menor volumen o longitud, (Jackman, 2009).

Definición 3. (*Región de la Densidad Posterior Más Alta*) Una región $C \subseteq \Omega$ es una Región de la Densidad Posterior Más Alta al $100(1 - \alpha)\%$ para θ bajo $p(\theta)$ si:

1. $P(\theta \in C) = 1 - \alpha$
2. $P(\theta_1) \geq P(\theta_2), \forall \theta_1 \in C \wedge \forall \theta_2 \notin C$

4.3. Inferencias para el coeficiente de correlación ρ

Si se tiene una muestra aleatoria $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ de una Normal Bivariada con vector de medias μ y matriz de varianzas y covarianzas Σ :

$$f(x, y|\mu_1, \mu_2, \sigma_1^2, \sigma_2^2, \rho) = \frac{1}{2\pi\sigma_1\sigma_2\sqrt{1-\rho^2}} \exp \left\{ -\frac{1}{2(1-\rho^2)} \left[\left(\frac{x-\mu_1}{\sigma_1} \right)^2 - 2\rho \left(\frac{x-\mu_1}{\sigma_1} \right) \left(\frac{y-\mu_2}{\sigma_2} \right) + \left(\frac{y-\mu_2}{\sigma_2} \right)^2 \right] \right\} \quad (4-1)$$

entonces su función de verosimilitud es la siguiente:

$$L(\mu_1, \mu_2, \sigma_1^2, \sigma_2^2, \rho) = \left(\frac{1}{2\pi} \right)^n \left(\frac{1}{\sigma_1} \right)^n \left(\frac{1}{\sigma_2} \right)^n \left(\frac{1}{1-\rho^2} \right)^{\frac{n}{2}} \exp \left\{ -\frac{1}{2(1-\rho^2)} \sum_{i=1}^n \left[\left(\frac{x_i-\mu_1}{\sigma_1} \right)^2 - 2\rho \left(\frac{x_i-\mu_1}{\sigma_1} \right) \left(\frac{y_i-\mu_2}{\sigma_2} \right) + \left(\frac{y_i-\mu_2}{\sigma_2} \right)^2 \right] \right\} \quad (4-2)$$

Pero, manipulando las expresiones que se encuentran en el exponente (ver Apéndice A), se obtiene la siguiente verosimilitud simplificada:

$$L(\mu_1, \mu_2, \sigma_1^2, \sigma_2^2, \rho) = \left(\frac{1}{2\pi}\right)^n \left(\frac{1}{\sigma_1}\right)^n \left(\frac{1}{\sigma_2}\right)^n \left(\frac{1}{1-\rho^2}\right)^{\frac{n}{2}} \exp\left\{-\frac{1}{2(1-\rho^2)}A\right\} \quad (4-3)$$

donde A es:

$$A = \frac{1}{\sigma_1^2} [(n-1)S_1^2 + n(\bar{x} - \mu_1)^2] - \frac{2\rho}{\sigma_1\sigma_2} [(n-1)S_{12} + n(\bar{x} - \mu_1)(\bar{y} - \mu_2)] + \frac{1}{\sigma_2^2} [(n-1)S_2^2 + n(\bar{y} - \mu_2)^2] \quad (4-4)$$

Simplificando lo anterior tenemos finalmente que:

$$L(\mu_1, \mu_2, \sigma_1^2, \sigma_2^2, \rho) \propto \left(\frac{1}{\sigma_1}\right)^n \left(\frac{1}{\sigma_2}\right)^n \left(\frac{1}{1-\rho^2}\right)^{\frac{n}{2}} \exp\left\{-\frac{1}{2(1-\rho^2)}A\right\} \quad (4-5)$$

4.3.1. Selección de Distribuciones Apriori para ρ

Debido a que es necesario plantear una distribución apriori para el parámetro del cual se quiere realizar la inferencia, en este caso de ρ , se emplea una distribución de probabilidad que esté definida en el rango de dicho parámetro, es decir, entre -1 y 1.

Distribución de McCullagh (1)

Como distribución apriori informativa para el coeficiente de correlación se decidió trabajar en primera instancia con la distribución univariada propuesta por McCullagh (1989) cuyo rango precisamente es entre -1 y 1. Una descripción de dicha apriori se presenta a continuación: Sea X una variable aleatoria definida en el intervalo $(-1, 1)$ cuya p.d.f es de la siguiente forma (McCullagh, 1989):

$$f_x(x; \theta, v) = \frac{(1-x^2)^{v-\frac{1}{2}}}{(1-2\theta x + \theta^2)^v B(v + \frac{1}{2}, \frac{1}{2})} \quad (4-6)$$

Si X sigue la distribución dada en 4-6, entonces se dice que X sigue una distribución de McCullagh.

Dicha densidad está relacionada con la densidad de la variable X' :

$$f_{X'}(X'; \theta, v) = \frac{(1-x'^2)^{v-\frac{1}{2}}(1-\theta^2)}{(1-2\theta x' + \theta^2)^{v+1} B(v + \frac{1}{2}, \frac{1}{2})} \quad (4-7)$$

con $-1 < x' < 1$ y donde $B(v + \frac{1}{2}, \frac{1}{2})$ es la función Beta, evaluada en $v + \frac{1}{2}$ y $\frac{1}{2}$, la cual se define de la siguiente manera:

$$B(x, y) = \frac{\Gamma(x)\Gamma(y)}{\Gamma(x+y)} \quad (4-8)$$

Y a su vez, $\Gamma(x)$ es la función Gamma evaluada en x .

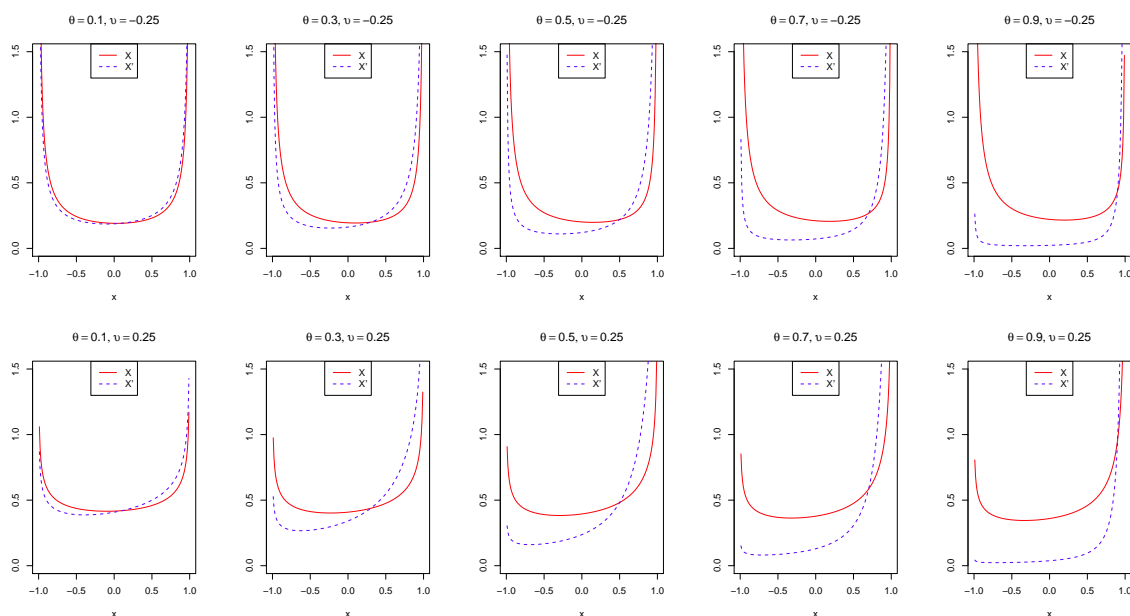
Ambas variables aleatorias se relacionan de la siguiente manera:

$$X' - \theta = \frac{(X - \theta)(\theta^2 - 1)}{1 - 2\theta X + \theta^2} \quad (4-9)$$

para todo $-1 < \theta < 1$ y $v > -\frac{1}{2}$ en ambas densidades.

En las siguientes gráficas se observan diferentes distribuciones apriori para diferentes conjuntos de valores de θ y v . Se consideran los escenarios siguientes:

1. Valores para el parámetro $\theta = 0.1, 0.3, 0.5, 0.7, 0.9$.
2. Valores para el parámetro $v = -0.25, 0.25, 0.5, 1, 2, 4$.
3. Observación: La línea en rojo corresponde a la pdf para X y la línea azul a la pdf para X' .
4. Observación: Como se evidencia gráficamente, cuando $v \rightarrow \infty$ ambas densidades tienden a adquirir una forma de campana para θ fijo. McCullagh en (McCullagh, 1989), menciona que ambas variables, X y X' , convergen en distribución a una variable aleatoria Normal.



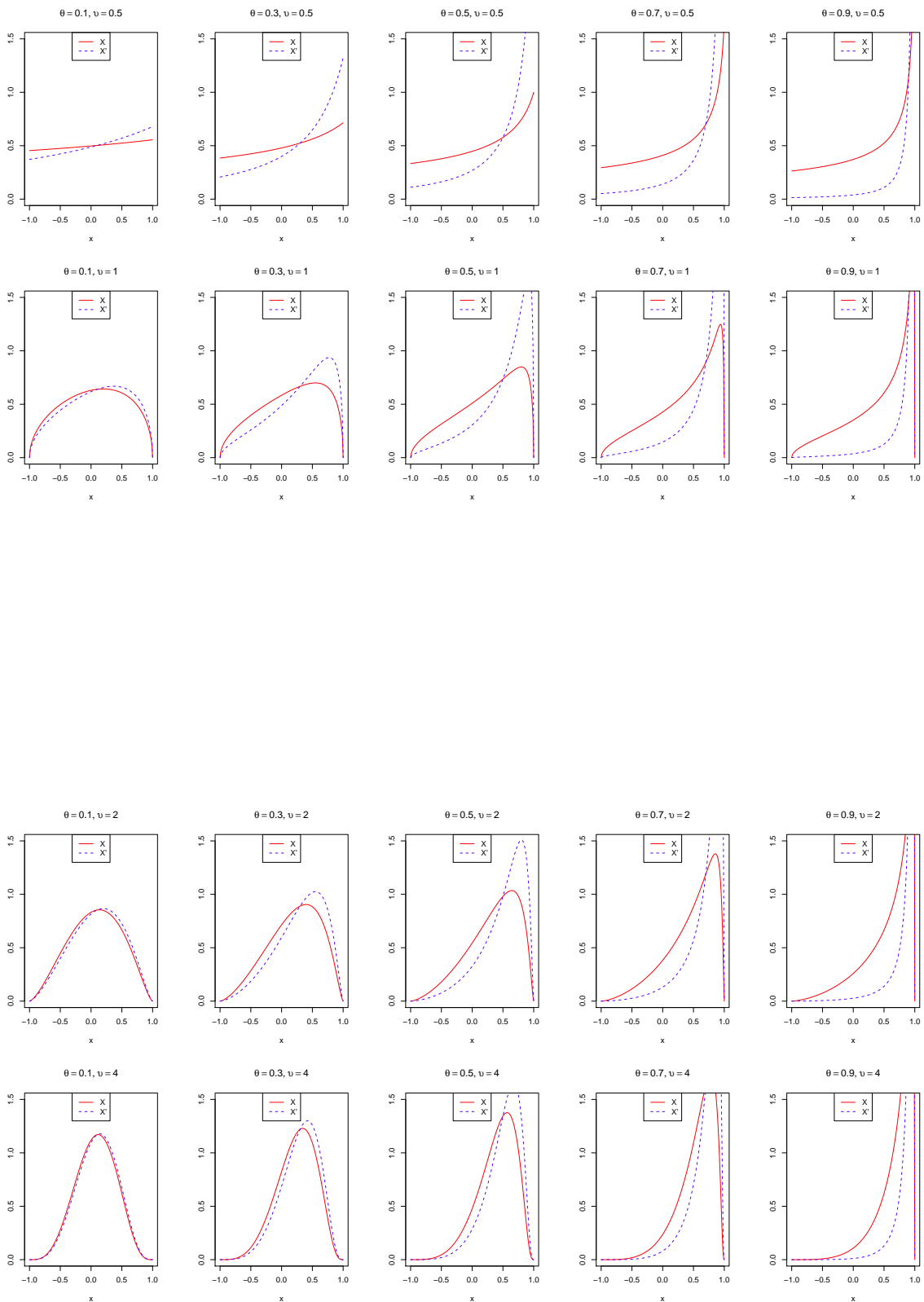


Figura 4-1: Distribución a priori para ρ con diferentes valores parametrales

Por último, McCullagh expresa los respectivos valores esperados y varianzas para X y X' :

$$E(X) = \frac{v\theta}{v+1} \quad (4-10)$$

$$V(X) = \frac{\left[1 - \frac{v(v-1)}{(v+1)(v+2)}\theta^2\right]}{2(v+1)} \quad (4-11)$$

$$E(X') = \theta \quad (4-12)$$

$$V(X') = \frac{1 - \theta^2}{2(v+1)} \quad (4-13)$$

Distribución a partir de: “Quantile matching priors” (2)

Como distribución apriori (2) no informativa para el parámetro de interés se especificó la siguiente distribución, (Ghosh et al., 2010):

$$p(\rho) \propto \frac{1}{\sigma_1\sigma_2(1-\rho^2)} \quad (4-14)$$

4.3.2. Obtención de Distribuciones Aposteriori para ρ

El Teorema de Bayes para parámetros continuos es comúnmente expresado así, (Jackman, 2009).:

$$p(\theta|y) \propto p(y|\theta)p(\theta) \quad (4-15)$$

donde la constante de proporcionalidad es:

$$\left[\int p(y|\theta)p(\theta) d\theta \right]^{-1} \quad (4-16)$$

Esto indica que la distribución aposteriori es proporcional al producto de la verosimilitud por la apriori. Por tanto, las distribuciones aposterioris para ρ se calcularán bajo este principio:

Distribución de McCullagh (1)

$$p(\mu_1, \mu_2, \sigma_1^2, \sigma_2^2, \rho | \text{Datos}) \propto L(\mu_1, \mu_2, \sigma_1^2, \sigma_2^2, \rho) \times p(\rho) \quad (4-17)$$

Si se tiene que la distribución apriori es de la forma:

$$p(\rho|\theta, v) = \frac{(1-\rho^2)^{v-\frac{1}{2}}}{(1-2\theta\rho+\theta^2)^v B(v+\frac{1}{2}, \frac{1}{2})} \quad (4-18)$$

donde

$$-1 < \theta < 1 \text{ y } v > -\frac{1}{2}$$

Entonces:

$$p(\mu_1, \mu_2, \sigma_1^2, \sigma_2^2, \rho | \text{Datos}) \propto \left(\frac{1}{\sigma_1}\right)^n \left(\frac{1}{\sigma_2}\right)^n \left(\frac{1}{1-\rho^2}\right)^{\frac{n}{2}} \exp\left\{-\frac{1}{2(1-\rho^2)}A\right\} \frac{(1-\rho^2)^{v-\frac{1}{2}}}{(1-2\theta\rho+\theta^2)^v B(v+\frac{1}{2}, \frac{1}{2})} \quad (4-19)$$

Que es igual a:

$$p(\mu_1, \mu_2, \sigma_1^2, \sigma_2^2, \rho | \text{Datos}) \propto \left(\frac{1}{\sigma_1}\right)^n \left(\frac{1}{\sigma_2}\right)^n (1-\rho^2)^{v-\frac{1}{2}-\frac{n}{2}} \exp\left\{-\frac{1}{2(1-\rho^2)}A\right\} \frac{1}{(1-2\theta\rho+\theta^2)^v B(v+\frac{1}{2}, \frac{1}{2})} \quad (4-20)$$

Es decir:

$$p(\mu_1, \mu_2, \sigma_1^2, \sigma_2^2, \rho | \text{Datos}) \propto \left(\frac{1}{\sigma_1}\right)^n \left(\frac{1}{\sigma_2}\right)^n (1-\rho^2)^{v-\frac{1+n}{2}} \exp\left\{-\frac{1}{2(1-\rho^2)}A\right\} \frac{1}{(1-2\theta\rho+\theta^2)^v B(v+\frac{1}{2}, \frac{1}{2})} \quad (4-21)$$

Los intervalos bayesianos se obtienen generando una gran cantidad de muestras para la distribución conjunta a posteriori de ρ a partir de cada una de las dos distribuciones a priori descritas anteriormente (McCullagh y Quantile Matching Priors), empleando la metodología conocida como muestreador de Gibbs, la cual hace parte de los métodos MCMC. Y a partir de dichas muestras, obtener para cada una los cuantiles 0.025 y 0.975, los cuales constituirán los límites para los intervalos de credibilidad.

El proceso en general es el siguiente:

1. Obtener una expresión para la distribución posterior conjunta $p(\mu_1, \mu_2, \sigma_1, \sigma_2, \rho | \text{Datos})$. Igualmente obtener las expresiones para las distribuciones marginales posteriores.
2. Determinar un valor inicial para los parámetros $\mu_1, \mu_2, \sigma_1, \sigma_2, \rho$.
3. Generar μ_1 de la distribución condicional para μ_1 teniendo en cuenta los valores iniciales.
4. Reemplazar μ_1 obtenido del paso anterior en la expresión para la distribución condicional de μ_2 y generar un valor para μ_2 .
5. Repetir esta secuencia en los demás parámetros usando las expresiones de las distribuciones condicionales respectivas.

6. Finalmente se obtiene una muestra con valores de cada uno de los parámetros de la distribución a posteriori conjunta.

Esto se repite las veces que se consideren necesarias para obtener una muestra de tamaño adecuado. Se realiza un “quemado” de muestras iniciales y finalmente se calculan los cuantiles de la muestra obtenida para ρ .

Para hacer uso del Muestreador de Gibbs en el estudio de simulación es necesario obtener las distribuciones condicionales de cada parámetro. Luego,

Para μ_1 : Considerando solo los términos de la a posteriori que contengan a μ_1 , es decir, la parte de la exponencial en A; y manipulando esa expresión (ver anexo) se llega a lo siguiente:

$$p(\mu_1 | \mu_2, \sigma_1^2, \sigma_2^2, \rho, \text{Datos}) \propto \exp \left\{ -\frac{1}{2(1-\rho^2)} \frac{n}{\sigma_1^2} \left[\mu_1 - \left(\bar{x} - \rho \frac{\sigma_1}{\sigma_2} (\bar{y} - \mu_2) \right) \right]^2 \right\} \quad (4-22)$$

El cual es el Kernel de una normal:

$$(\mu_1 | \mu_2, \sigma_1^2, \sigma_2^2, \rho, \text{Datos}) \sim N \left(\bar{x} - \rho \frac{\sigma_1}{\sigma_2} (\bar{y} - \mu_2), \frac{\sigma_1^2(1-\rho^2)}{n} \right) \quad (4-23)$$

Para μ_2 : De igual manera, siguiendo el mismo procedimiento que se hizo en el paso anterior se llega a que

$$(\mu_2 | \mu_1, \sigma_1^2, \sigma_2^2, \rho, \text{Datos}) \sim N \left(\bar{y} - \rho \frac{\sigma_2}{\sigma_1} (\bar{x} - \mu_1), \frac{\sigma_2^2(1-\rho^2)}{n} \right) \quad (4-24)$$

Para σ_1^2 : Observando en especial la parte de la exponencial que se obtuvo en la distribución a posteriori, específicamente en A; utilizando solo aquello que dependa de σ_1^2 y manipulando la expresión (ver anexo), se tiene la siguiente distribución condicional para σ_1^2 :

$$p(\sigma_1^2 | \mu_2, \mu_1, \sigma_2^2, \rho, \text{Datos}) \propto \frac{1}{\sigma_1^n} \exp \left\{ -\frac{1}{2(1-\rho^2)} \left[\frac{1}{\sigma_1^2} B - \frac{1}{\sigma_1} C \right] \right\} \quad (4-25)$$

donde:

$$B = (n-1)S_1^2 + n(\bar{x} - \mu_1)^2 \quad (4-26)$$

y

$$C = \frac{2\rho}{\sigma_2} [(n-1)S_{12} + n(\bar{x} - \mu_1)(\bar{y} - \mu_2)] \quad (4-27)$$

Para σ_2^2 : De igual manera se llega a que la distribución condicional para σ_2^2 es:

$$p(\sigma_2^2 | \mu_2, \mu_1, \sigma_1^2, \rho, \text{Datos}) \propto \frac{1}{\sigma_2^n} \exp \left\{ -\frac{1}{2(1-\rho^2)} \left[\frac{1}{\sigma_2^2} D - \frac{1}{\sigma_2} E \right] \right\} \quad (4-28)$$

donde:

$$D = (n-1)S_2^2 + n(\bar{y} - \mu_2)^2 \quad (4-29)$$

y

$$E = \frac{2\rho}{\sigma_1} [(n-1)S_{12} + n(\bar{x} - \mu_1)(\bar{y} - \mu_2)] \quad (4-30)$$

Para ρ :

$$p(\rho | \mu_1, \mu_2, \sigma_1^2, \sigma_2^2, \text{Datos}) \propto (1-\rho^2)^{v-\frac{1+n}{2}} \exp \left\{ -\frac{1}{2(1-\rho^2)} A \right\} \frac{1}{(1-2\theta\rho + \theta^2)^v B(v + \frac{1}{2}, \frac{1}{2})} \quad (4-31)$$

Distribución a partir de: “Quantile matching priors” (2)

$$p(\mu_1, \mu_2, \sigma_1^2, \sigma_2^2, \rho | \text{Datos}) \propto L(\mu_1, \mu_2, \sigma_1^2, \sigma_2^2, \rho) \times p(\rho) \quad (4-32)$$

Por tanto,

$$p(\mu_1, \mu_2, \sigma_1^2, \sigma_2^2, \rho | \text{Datos}) \propto \left(\frac{1}{\sigma_1} \right)^n \left(\frac{1}{\sigma_2} \right)^n \left(\frac{1}{1-\rho^2} \right)^{\frac{n}{2}} \exp \left\{ -\frac{1}{2(1-\rho^2)} A \right\} \frac{1}{\sigma_1 \sigma_2 (1-\rho^2)} \quad (4-33)$$

Las siguientes son las distribuciones condicionales de cada parámetro, necesarias para emplear el Muestreador de Gibbs en el estudio de simulación:

Para μ_1 : Ver anexo:

$$p(\mu_1 | \mu_2, \sigma_1^2, \sigma_2^2, \rho, \text{Datos}) \propto \exp \left\{ -\frac{1}{2(1-\rho^2)} \frac{n}{\sigma_1^2} \left[\mu_1 - \left(\bar{x} - \rho \frac{\sigma_1}{\sigma_2} (\bar{y} - \mu_2) \right) \right]^2 \right\} \quad (4-34)$$

El cual es el Kernel de una normal:

$$(\mu_1 | \mu_2, \sigma_1^2, \sigma_2^2, \rho, \text{Datos}) \sim N \left(\bar{x} - \rho \frac{\sigma_1}{\sigma_2} (\bar{y} - \mu_2), \frac{\sigma_1^2 (1-\rho^2)}{n} \right) \quad (4-35)$$

Para μ_2 : De igual manera:

$$(\mu_2 | \mu_1, \sigma_1^2, \sigma_2^2, \rho, \text{Datos}) \sim N \left(\bar{y} - \rho \frac{\sigma_2}{\sigma_1} (\bar{x} - \mu_1), \frac{\sigma_2^2(1 - \rho^2)}{n} \right) \quad (4-36)$$

Para σ_1^2 : Ver anexo:

$$p(\sigma_1^2 | \mu_2, \mu_1, \sigma_2^2, \rho, \text{Datos}) \propto \frac{1}{\sigma_1^{n+1}} \exp \left\{ -\frac{1}{2(1 - \rho^2)} \left[\frac{1}{\sigma_1^2} B - \frac{1}{\sigma_1} C \right] \right\} \quad (4-37)$$

donde:

$$B = (n - 1)S_1^2 + n(\bar{x} - \mu_1)^2 \quad (4-38)$$

y

$$C = \frac{2\rho}{\sigma_2} [(n - 1)S_{12} + n(\bar{x} - \mu_1)(\bar{y} - \mu_2)] \quad (4-39)$$

Para σ_2^2 : De igual manera:

$$p(\sigma_2^2 | \mu_2, \mu_1, \sigma_1^2, \rho, \text{Datos}) \propto \frac{1}{\sigma_2^{n+1}} \exp \left\{ -\frac{1}{2(1 - \rho^2)} \left[\frac{1}{\sigma_2^2} D - \frac{1}{\sigma_2} E \right] \right\} \quad (4-40)$$

donde:

$$D = (n - 1)S_2^2 + n(\bar{y} - \mu_2)^2 \quad (4-41)$$

y

$$E = \frac{2\rho}{\sigma_1} [(n - 1)S_{12} + n(\bar{x} - \mu_1)(\bar{y} - \mu_2)] \quad (4-42)$$

Para ρ :

$$p(\rho | \mu_1, \mu_2, \sigma_1^2, \sigma_2^2, \text{Datos}) \propto \frac{1}{(1 - \rho^2)^{\frac{n}{2}+1}} \exp \left\{ -\frac{1}{2(1 - \rho^2)} A \right\} \quad (4-43)$$

5 Resultados de la Simulación - Metodología Bayesiana

De igual manera, para comparar los intervalos de credibilidad derivados a partir de las dos distribuciones a priori y a posteriori consideradas se realizó una simulación en R en la cual se consideraron combinaciones de (ρ, n) con valores de $\rho = 0.0, 0.1, 0.3, 0.5, 0.7, 0.9$ y de $n = 5, 10, 20, 50, 100$. Para cada pareja se realizaron cadenas de Markov, las cuales contienen cada una 7000 valores muestrales de ρ . Para cada a priori y combinación se calculó la mediana de la longitud de los 1000 intervalos calculados y la proporción de intervalos que cubren el verdadero valor de ρ , esto es lo que se llama el nivel de confianza real. Las siguientes tablas presentan los resultados.

Nota: Para las gráficas 4.1 a 4.6, los puntos corresponden a cada tamaño muestral partiendo desde $n = 5$ hasta $n = 100$.

ρ		$n = 5$		$n = 10$		$n = 20$	
		Apriori 1	Apriori 2	Apriori 1	Apriori 2	Apriori 1	Apriori 2
0.0	Longitud	1.232323	1.494950	1.010101	1.151515	0.7878788	0.8287878
	Nivel	0.74	0.88	0.87	0.88	0.96	0.79
0.1	Longitud	1.212121	1.525252	0.989899	1.171717	0.7676768	0.8484848
	Nivel	0.75	0.87	0.86	0.87	0.86	0.89
0.3	Longitud	1.151515	1.515151	0.929293	1.151515	0.7272727	0.8080808
	Nivel	0.77	0.9	0.85	0.86	0.84	0.84
0.5	Longitud	1.090909	1.484849	0.8686869	1	0.6265152	0.6464646
	Nivel	0.78	0.93	0.92	0.94	0.92	0.81
0.7	Longitud	1	1.272727	0.7272728	0.7676768	0.4848484	0.5050505
	Nivel	0.86	0.88	0.82	0.96	0.86	0.91
0.9	Longitud	0.868687	0.7272727	0.4545455	0.3636364	0.2222223	0.2020202
	Nivel	0.59	0.72	0.91	0.97	0.87	0.95

Tabla 5-1: Longitud y nivel de confianza de los intervalos $n = 5$, $n = 10$ y $n = 20$

ρ		$n = 50$		$n = 100$	
		Apriori 1	Apriori 2	Apriori 1	Apriori 2
0.0	Longitud	0.5252525	0.5454546	0.3838384	0.3838384
	Nivel	0.94	0.91	0.97	0.9
0.1	Longitud	0.5252525	0.5454545	0.3838384	0.3838384
	Nivel	0.87	0.87	0.95	0.92
0.3	Longitud	0.4848485	0.5050505	0.3434343	0.3636363
	Nivel	0.9	0.97	0.95	0.95
0.5	Longitud	0.4040404	0.4242424	0.3030303	0.3030303
	Nivel	0.92	0.89	0.93	0.99
0.7	Longitud	0.2929293	0.2929293	0.2020202	0.2020202
	Nivel	0.79	0.94	0.9	0.89
0.9	Longitud	0.1212121	0.1010102	0.08080815	0.0808081
	Nivel	0.76	0.94	0.82	0.84

Tabla 5-2: Longitud y nivel de confianza de los intervalos $n = 50$ y $n = 100$

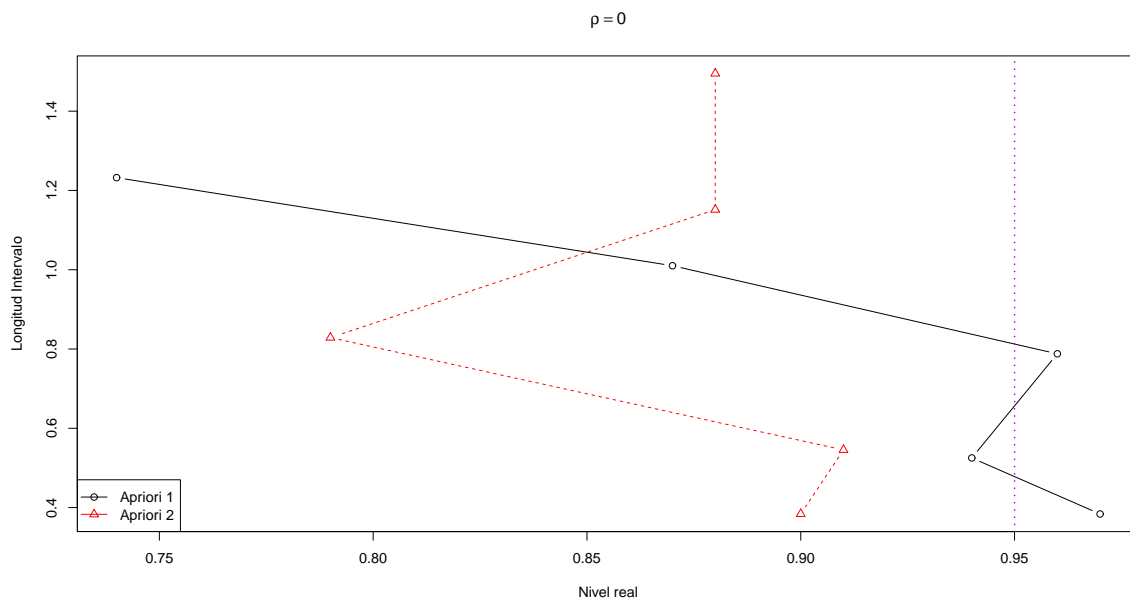


Figura 5-1: Amplitud y Nivel real alcanzado por cada intervalo con $\rho=0$

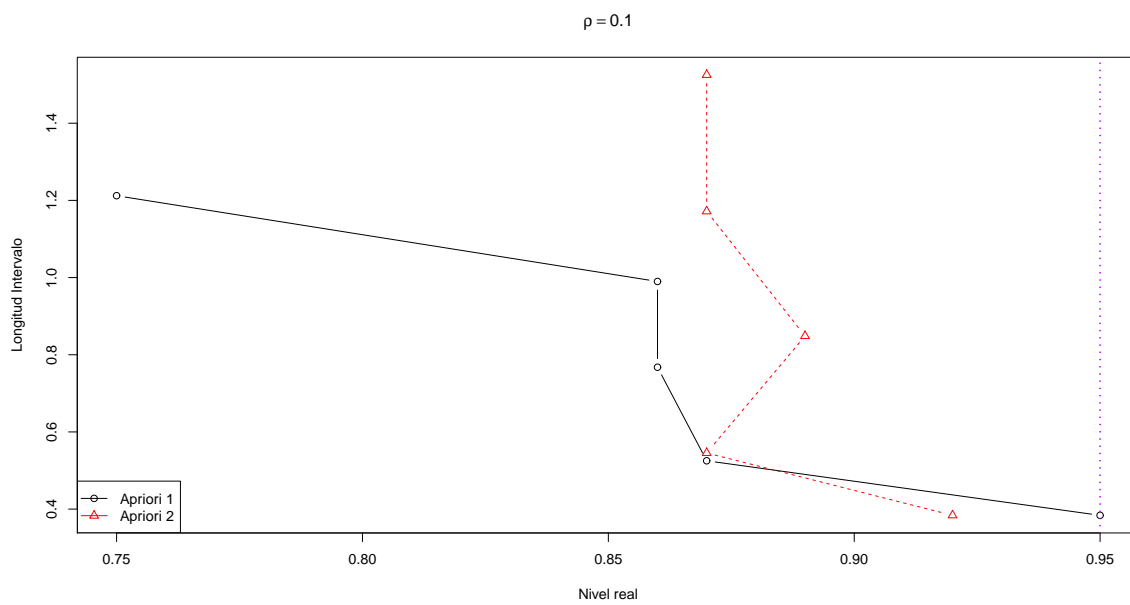


Figura 5-2: Amplitud y Nivel real alcanzado por cada intervalo con $\rho=0.1$

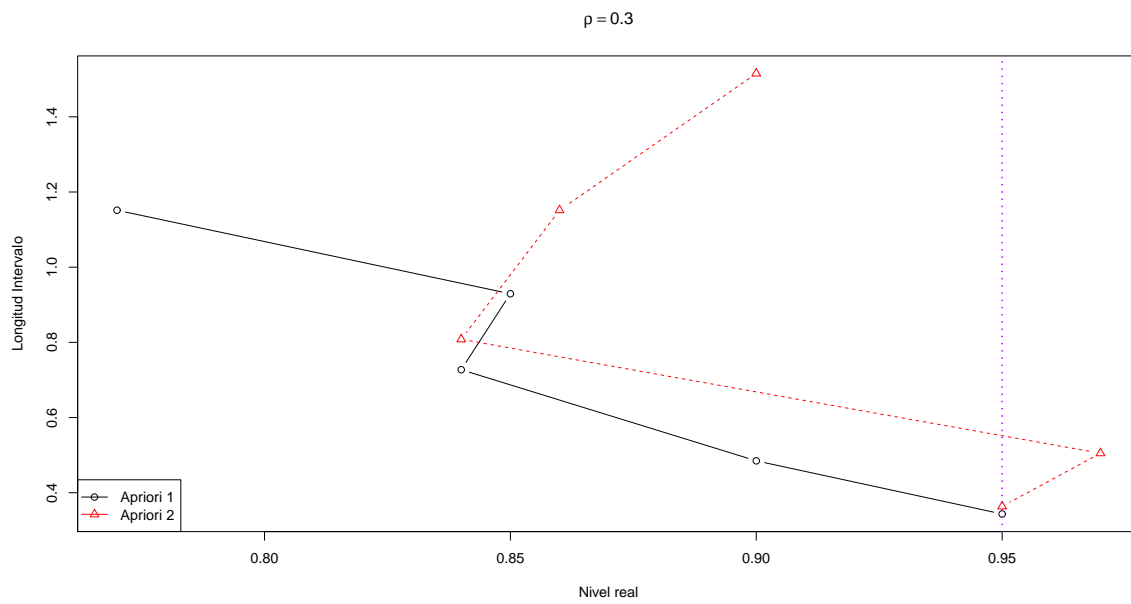


Figura 5-3: Amplitud y Nivel real alcanzado por cada intervalo con $\rho=0.3$

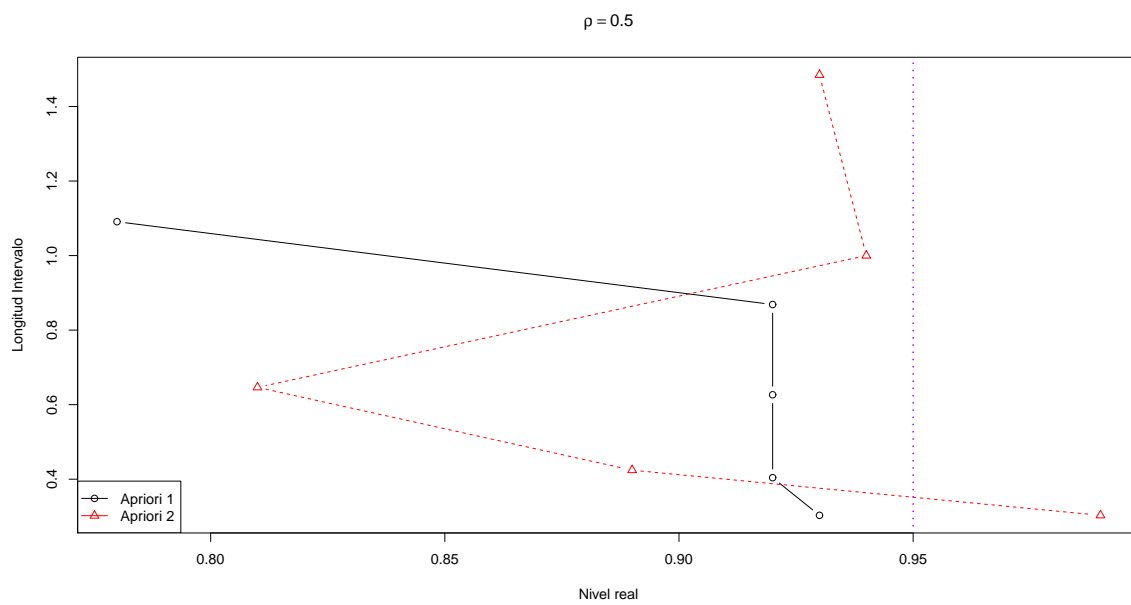


Figura 5-4: Amplitud y Nivel real alcanzado por cada intervalo con $\rho=0.5$

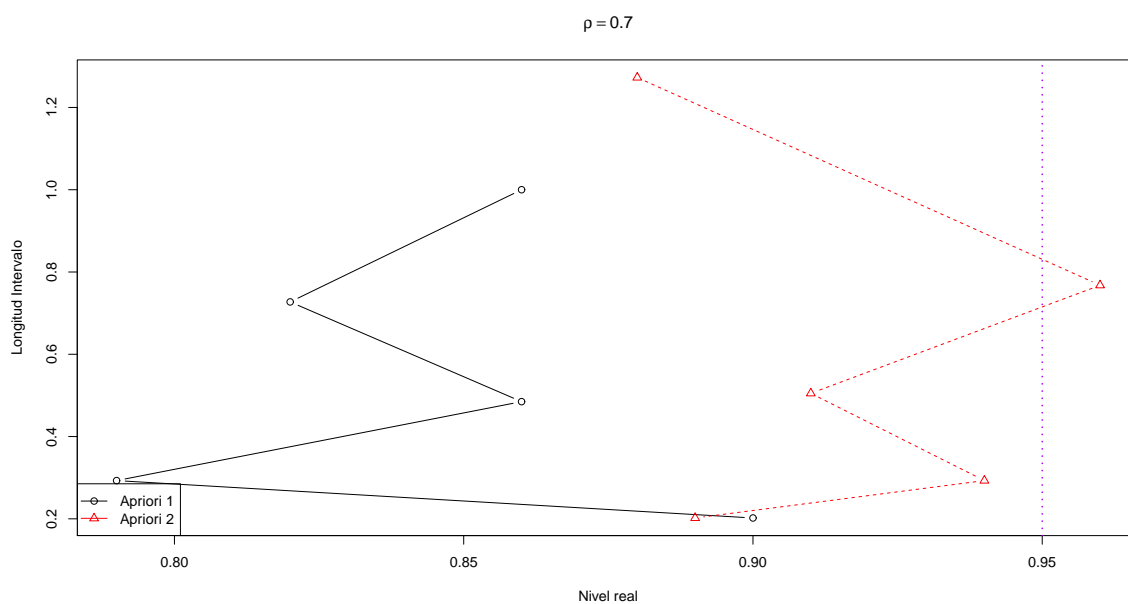


Figura 5-5: Amplitud y Nivel real alcanzado por cada intervalo con $\rho=0.7$

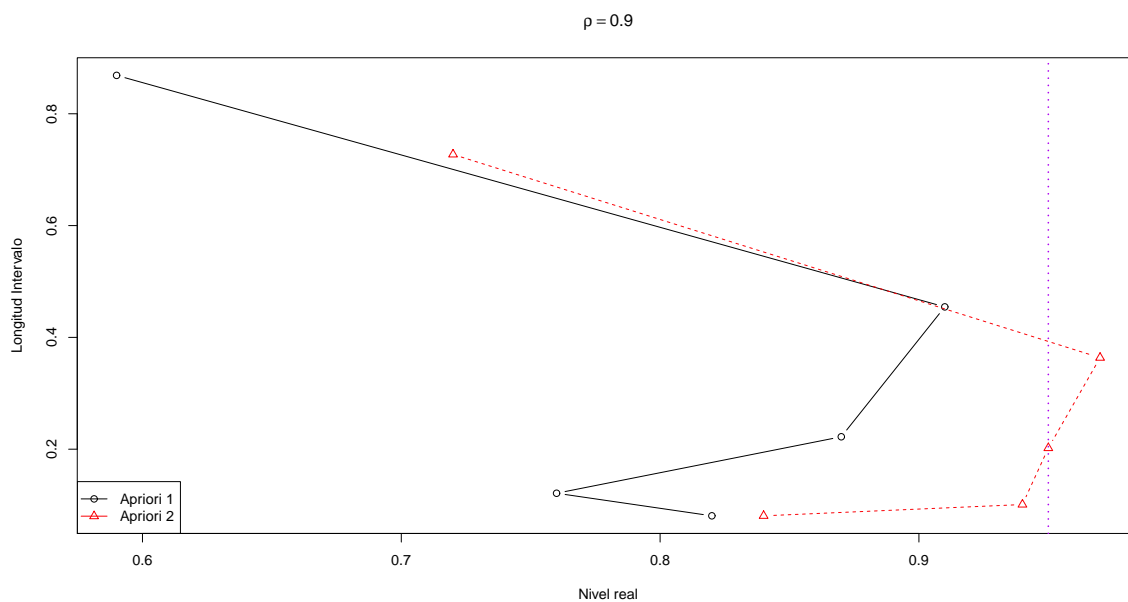


Figura 5-6: Amplitud y Nivel real alcanzado por cada intervalo con $\rho=0.9$

De las anteriores tablas y gráficas se puede concluir que no se puede considerar longitud o nivel real aisladamente para seleccionar la mejor distribución apriori en la construcción de un intervalo de credibilidad debido a que en muestras pequeñas:

- En la mayoría de los casos la apriori 2 es mejor en cuanto a nivel real del intervalo. Aunque cabe resaltar que los porcentajes de cobertura reales distan mucho del nivel teórico, en este caso el 95 %.
- El comportamiento de la longitud del intervalo para la distribución apriori 1, hace que este sea el mejor debido a que son en su mayoría los más cortos.

Pero a medida que el tamaño de muestra aumenta:

- La longitud de los intervalos construidos por las dos distribuciones apriori disminuye considerablemente, siendo para ambas aprioris los resultados muy similares entre sí.
- El nivel real en dichos intervalos se va acercando más al nivel nominal (95 %).

Índice de resúmenes De manera similar, se realiza un resumen de los resultados anteriores a partir del Índice propuesto en la sección anterior (parte clásica). Recordemos que este índice busca favorecer a aquellos métodos que presenten longitudes de intervalo pequeños y niveles reales de confianza cercanos o mayores al 95 %:

$$I = (2 - LI) \frac{NR}{NN}$$

donde:

- LI: Mediana de la Longitud del intervalo.
- NR: Promedio del nivel real del intervalo.
- NN: Nivel nominal de los intervalos, que en este caso es 0.95.

Luego, a valores mayores del índice propuesto, mejor el intervalo obtenido.

ρ	$n = 5$		$n = 10$		$n = 20$	
	Apriori 1	Apriori 2	Apriori 1	Apriori 2	Apriori 1	Apriori 2
0	0.60	0.47	0.91	0.79	1.22	0.97
0.1	0.62	0.43	0.91	0.76	1.12	1.08
0.3	0.69	0.46	0.96	0.77	1.13	1.05
0.5	0.75	0.50	1.10	0.99	1.33	1.15
0.7	0.91	0.67	1.10	1.25	1.37	1.43
0.9	0.86	0.96	1.48	1.67	1.63	1.80

Tabla 5-3: Índice Tamaño de muestra 5, 10 y 20

ρ	$n = 50$		$n = 100$	
	Apriori 1	Apriori 2	Apriori 1	Apriori 2
0	1.46	1.39	1.65	1.53
0.1	1.35	1.33	1.62	1.57
0.3	1.44	1.53	1.66	1.64
0.5	1.55	1.48	1.66	1.77
0.7	1.42	1.69	1.70	1.68
0.9	1.50	1.88	1.66	1.70

Tabla 5-4: Índice Tamaño de muestra 50 y 100

En términos generales se observa que el mejor método o apriori para construir Intervalos de Credibilidad para el Coeficiente de correlación en la distribución normal bivariada en el caso de muestras pequeñas es la distribución apriori de McCullagh, notándose esto principalmente para las correlaciones aproximadamente inferiores a 0.7 .

Cuando el tamaño de las muestras empleadas aumenta ($n = 50$ y $n = 100$) el comportamiento de las dos distribuciones apriori es similar, siendo levemente mejor la apriori de McCullagh en correlaciones menores que 0.7 y mejor la apriori a partir de “Quantile Matching priors” en correlaciones mayores o iguales a 0.7; según los casos de simulación considerados en este estudio.

Decir que los intervalos para ρ calculados bajo la metodología bayesiana son categóricamente mejores o peores que los obtenidos mediante la metodología clásica sería erróneo debido a que, por ejemplo, el Índice de resúmenes para $n = 50$ y $\rho = 0$, muestra que los intervalos calculados a partir la apriori de “Quantile Matching priors” tienen un índice menor a los de las metodologías clásicas, mientras que para $n = 50$ y $\rho = 0.9$ el comportamiento de los intervalos obtenidos mediante el uso de la apriori mencionada es similar al de las metodologías clásicas. Por último si se considera un tamaño de muestra igual a 100 y un coeficiente de correlación de 0.9, los intervalos mediante las metodologías bayesianas son considerablemente de menor calidad en comparación con los intervalos a partir de la metodología clásica.

La respuesta a la pregunta de cuál es la metodología, clásica o bayesiana, se debe usar para la construcción de un intervalo para el coeficiente de correlación de una distribución Normal Bivariada es, como lo muestran los resultados de las simulaciones, sensible a los tamaños de las muestras y al coeficiente de correlación verdadero.

6 Aplicaciones

6.1. Aplicación Base de datos Huevos

Estos datos fueron tomados en un laboratorio de la Facultad de Minas por un grupo de estudiantes de posgrado en estadística. Se tomaron 27 huevos y se calentaron a temperatura de 86°; su temperatura fue medida con un termostato. Las variables que aparecen en la muestra son:

- X_1 : Largo del huevo.
- X_2 : Ancho del huevo.
- X_3 : Peso inicial.
- X_4 : Volúmen en centímetros cúbicos.
- X_5 : Tiempo que duró el huevo en el agua.
- X_6 : Peso final.
- X_7 : Dureza: 1= Duro, 0= Blando.

Luego de realizar la validación de normalidad bivariada en los datos (Largo y Ancho), se obtuvo el coeficiente de correlación muestral y los respectivos intervalos de confianza al 95 % para ρ , el coeficiente de correlación verdadero entre Largo y Ancho de los huevos.

Los resultados son los siguientes:

r = 0.851471, n = 27			
Método	Lím. Inferior	Lím. Superior	Long. Intervalo
Bootstrap	0.71040	0.93288	0.22247
Arctanh	0.69698	0.93043	0.23345
L.R	0.69235	0.92694	0.23460
Jeyaratnam	0.69647	0.93056	0.23409
Z1	0.69816	0.93012	0.23196
Z2	0.69806	0.93014	0.23208
Z3	0.69917	0.92985	0.23068
Z4	0.69909	0.92987	0.23077
P.G	0.68828	0.92745	0.23917
Apriori McCullagh	0.53535	0.87879	0.34343
Apriori QMP	0.69697	0.91919	0.22222

Tabla 6-1: Intervalo de confianza para los diferentes métodos considerados en la Aplicación

Se observa un comportamiento similar para la mayoría de los intervalos en cuanto a longitud del mismo. El único que ofrece una “mala” calidad es el obtenido mediante la Apriori de McCullagh. Por lo anterior se sugiere emplear el método Bootstrap o la Apriori Quantile Matching Priors.

6.2. Aplicación Base de Datos Vinos

Se considera para esta aplicación la base de datos en la cual se miden ciertas características en un vino llamado Pinot Noir, base de datos que consta de 38 registros en cada una de las siguientes variables:

	Claridad	Aroma	Cuerpo	Sabor	Sabor a madera	Calidad	Región
1	1.00	3.30	2.80	3.10	4.10	9.80	1.00
2	1.00	4.40	4.90	3.50	3.90	12.60	1.00
3	1.00	3.90	5.30	4.80	4.70	11.90	1.00
4	1.00	3.90	2.60	3.10	3.60	11.10	1.00
5	1.00	5.60	5.10	5.50	5.10	13.30	1.00
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
35	0.80	5.00	5.70	5.50	4.80	13.50	1.00
36	0.80	3.50	4.70	4.20	3.30	12.20	1.00
37	0.80	4.30	5.50	3.50	5.80	10.30	1.00
38	0.80	5.20	4.80	5.70	3.50	13.20	1.00

Base de datos extraída de: Montgomery, D.C., Peck, E.A., and Vining, C.G. (2001) Introduction to Linear Regression Analysis. 3rd Edition, John Wiley and Sons.

Se toman las variables Cuerpo y Aroma para el respectivo análisis, en el cual se prueba normalidad bivariada, y se construyen los intervalos de confianza para el coeficiente de correlación con todas las metodologías por la facilidad que tiene todos en su construcción y a manera de ilustración:

r = 0.5489102, n = 38			
Método	Lím. Inferior	Lím. Superior	Long. Intervalo
Bootstrap	0.30270	0.74280	0.44010
Arctanh	0.27801	0.73893	0.46093
L.R	0.27650	0.73284	0.45634
Jeyaratnam	0.27746	0.73920	0.46174
Z1	0.27918	0.73835	0.45917
Z2	0.27911	0.73839	0.45928
Z3	0.28023	0.73783	0.45760
Z4	0.28018	0.73786	0.45769
P.G	0.27763	0.73608	0.45844
Apriori McCullagh	0.19192	0.65657	0.46465
Apriori QMP	0.27273	0.73737	0.46465

Tabla 6-2: Intervalo de confianza para los diferentes métodos considerados en la Aplicación

Según los resultados de la tabla anterior, se aprecia una gran similitud entre las longitudes de cada intervalo, siendo mejor el intervalo Bootstrap y le seguiría en desempeño el de Razón de Verosimilitud.

Por tanto para la estimación por intervalo del coeficiente de correlación entre el Cuerpo y Aroma de un vino, se recomienda usar estos dos métodos.

7 Conclusiones

Los procedimientos que se tuvieron en cuenta para la construcción de intervalos de confianza para el coeficiente de correlación en una distribución normal bivariada, ρ difieren en calidad dependiendo del tamaño de muestra empleado para ello. Se pudo observar que, en el caso de muestras pequeñas ($n = 5$ y $n = 10$), los que mejor desempeño tuvieron fueron los de Razón de Verosimilitud y Bootstrap. El método de Razón de Verosimilitud ofrece longitudes cortas de intervalo y nivel de confianza real más cercano al nivel nominal establecido en este estudio de simulación, 95 % para correlaciones desde 0 hasta 0.7. En los casos restantes, es decir, correlaciones de 0.8 y 0.9, el método Bootstrap supera al anterior.

Cuando $n = 20$, el método que se comporta mejor en la gran mayoría de las correlaciones consideradas es el de Pivote Generalizado. Le siguen en orden de calidad el intervalo obtenido por el método de Razón de Verosimilitud. Cuando $n = 50$ y $n = 100$, todos los métodos tienen un comportamiento casi igual. Las longitudes de intervalo para cada uno de estos dos casos ($n = 50$ y $n = 100$) disminuyen considerablemente a lo observado en el caso anterior, es decir, en comparación con $n = 20$, como se puede evidenciar en las tablas.

En la mayoría de los casos, el método de Bootstrap ofrece niveles de confianza reales menores que los de los demás métodos. Se observa un detalle muy particular con respecto a esta característica y todos los métodos a excepción del Pivote Generalizado: los niveles de confianza reales son muy parecidos entre sí cuando el tamaño de muestra es muy grande ($n=100$).

En la parte Bayesiana se puede notar que en la mayoría de los casos, los resultados obtenidos por la distribución apriori de McCullagh, cuando se tienen muestras pequeñas son los mejores si se reúnen los dos criterios de decisión en uno solo. Cuando las muestras tienen tamaños de 50 y 100, no existe mucha diferencia entre los intervalos construidos por McCullagh y Quantile Matching priors. De esto último se destaca que las longitudes de los intervalos disminuyen drásticamente al aumentar los tamaños muestrales, y que por ejemplo, no se observa diferencia entre los intervalos a partir de $n = 100$ y coeficientes de correlación de 0.5 en adelante.

8 Recomendaciones

Para un trabajo futuro con relación a este tema se recomienda explorar con más distribuciones a priori para ρ como complemento de la parte bayesiana debido a que en la literatura existen muchas más distribuciones para ρ que por cuestión de tiempo y recursos no fue posible trabajar. Además, se sugiere ampliar el rango de valores posibles para dicho parámetro con el propósito de evaluar más ampliamente la calidad de los Intervalos de Credibilidad que se construyan.

Por otro lado, como aporte adicional, se podría analizar también la distribución de probabilidad que siguen las longitudes de los intervalos de confianza y credibilidad para cada escenario que se desee considerar en un futuro trabajo sobre el tema.

Con respecto a la aplicaciones en la vida real se recomienda emplear el método de razón de verosimilitud que es el que en situaciones globales es el que mejores resultados presenta, tanto en longitud como en precisión.

9 Anexos

9.1. Tablas de Resultados Metodología Clásica

Las siguientes tablas presentan los resultados de las simulaciones realizadas para $n = 20, 50, 100$:

$n = 20$										
ρ	Bootstrap	Arctanh	L.R	Jeyaratnam	Z1	Z2	Z3	Z4	P.G	
0.0	Longitud	0.8696	0.84080	0.87200	0.86320	0.86390	0.85800	0.85850	0.8395	
	Nivel	0.950	0.963	0.963	0.958	0.959	0.956	0.956	0.957	
0.1	Longitud	0.85240	0.83480	0.86550	0.85670	0.85740	0.85150	0.85200	0.8351	
	Nivel	0.933	0.933	0.936	0.928	0.928	0.926	0.926	0.945	
0.2	Longitud	0.84160	0.8492	0.82210	0.85160	0.84360	0.83780	0.83830	0.8209	
	Nivel	0.931	0.940	0.937	0.940	0.939	0.937	0.937	0.949	
0.3	Longitud	0.81010	0.8173	0.79270	0.81970	0.81180	0.80620	0.80670	0.7870	
	Nivel	0.932	0.942	0.944	0.943	0.942	0.941	0.941	0.952	
0.4	Longitud	0.75510	0.7659	0.74510	0.76820	0.75990	0.76060	0.75510	0.7500	
	Nivel	0.939	0.948	0.943	0.949	0.946	0.946	0.946	0.944	
0.5	Longitud	0.68150	0.6900	0.67440	0.69210	0.68440	0.68500	0.68000	0.6877	
	Nivel	0.937	0.944	0.938	0.944	0.942	0.942	0.941	0.951	
0.6	Longitud	0.58140	0.5936	0.58370	0.59550	0.58860	0.58910	0.58450	0.5870	
	Nivel	0.940	0.948	0.944	0.948	0.946	0.946	0.944	0.944	
0.7	Longitud	0.46910	0.4850	0.48040	0.48660	0.48060	0.48110	0.47720	0.4810	
	Nivel	0.954	0.960	0.953	0.961	0.959	0.959	0.958	0.949	
0.8	Longitud	0.33913	0.3565	0.35636	0.35779	0.35316	0.35351	0.35046	0.3549	
	Nivel	0.937	0.945	0.944	0.945	0.943	0.943	0.942	0.951	
0.9	Longitud	0.18235	0.1935	0.19580	0.19424	0.19155	0.19176	0.18999	0.19315	
	Nivel	0.943	0.948	0.943	0.948	0.947	0.947	0.947	0.945	

Tabla 9-1: Longitud y nivel de confianza de los intervalos. Tamaño de muestra 20

$n = 50$									
ρ	Bootstrap	Arctanh	L.R	Jeyaratnam	Z1	Z2	Z3	Z4	P.G
0.0	Longitud	0.55190	0.54410	0.55260	0.55040	0.55050	0.54910	0.54910	0.5408
	Nivel	0.953	0.951	0.953	0.952	0.952	0.951	0.951	0.948
0.1	Longitud	0.54830	0.54060	0.54900	0.54680	0.54690	0.54550	0.54550	0.5389
	Nivel	0.954	0.954	0.954	0.954	0.954	0.954	0.954	0.953
0.2	Longitud	0.53430	0.52720	0.53500	0.53290	0.53300	0.53160	0.53160	0.5271
	Nivel	0.958	0.958	0.958	0.958	0.958	0.957	0.957	0.937
0.3	Longitud	0.51140	0.50500	0.51200	0.51000	0.51000	0.50870	0.50870	0.5010
	Nivel	0.949	0.949	0.950	0.949	0.949	0.949	0.949	0.934
0.4	Longitud	0.46960	0.46460	0.47020	0.46830	0.46830	0.46710	0.46720	0.4686
	Nivel	0.941	0.938	0.941	0.939	0.939	0.939	0.939	0.939
0.5	Longitud	0.42400	0.42030	0.42450	0.42280	0.42290	0.42170	0.42180	0.4187
	Nivel	0.943	0.943	0.943	0.942	0.942	0.941	0.941	0.935
0.6	Longitud	0.36230	0.36010	0.36270	0.36120	0.36130	0.36030	0.36030	0.3594
	Nivel	0.957	0.958	0.957	0.956	0.956	0.956	0.956	0.953
0.7	Longitud	0.28690	0.28620	0.28730	0.28610	0.28610	0.28530	0.28540	0.2942
	Nivel	0.950	0.950	0.950	0.950	0.950	0.949	0.950	0.955
0.8	Longitud	0.20962	0.20984	0.20990	0.20901	0.20904	0.20840	0.20850	0.20769
	Nivel	0.947	0.943	0.947	0.947	0.947	0.946	0.946	0.94
0.9	Longitud	0.11260	0.11323	0.11279	0.11228	0.11230	0.11197	0.11198	0.11267
	Nivel	0.955	0.953	0.955	0.955	0.955	0.953	0.953	0.95

Tabla 9-2: Longitud y nivel de confianza de los intervalos. Tamaño de muestra 50

		$n = 100$									
ρ		Bootstrap	Arctanh	L.R	Jeyaratnam	Z1	Z2	Z3	Z4	P.G	
0.0	Longitud	0.38740	0.39120	0.38820	0.39140	0.39070	0.39070	0.39070	0.39020	0.3858	
	Nivel	0.947	0.953	0.951	0.953	0.953	0.953	0.953	0.952	0.933	
0.1	Longitud	0.38460	0.38870	0.38580	0.38900	0.38820	0.38820	0.38770	0.38770	0.3824	
	Nivel	0.951	0.953	0.951	0.953	0.952	0.952	0.952	0.952	0.946	
0.2	Longitud	0.37630	0.37880	0.37600	0.37900	0.37830	0.37830	0.37780	0.37780	0.3735	
	Nivel	0.947	0.951	0.950	0.951	0.950	0.950	0.949	0.950	0.932	
0.3	Longitud	0.35580	0.35950	0.35710	0.35970	0.35900	0.35900	0.35860	0.35860	0.3554	
	Nivel	0.942	0.945	0.944	0.945	0.943	0.943	0.943	0.943	0.943	
0.4	Longitud	0.32970	0.33140	0.32950	0.33160	0.33090	0.33100	0.33050	0.33050	0.3291	
	Nivel	0.949	0.949	0.948	0.949	0.949	0.949	0.949	0.949	0.945	
0.5	Longitud	0.29520	0.29700	0.29560	0.29720	0.29660	0.29660	0.29620	0.29620	0.2943	
	Nivel	0.950	0.956	0.954	0.956	0.956	0.956	0.956	0.956	0.956	
0.6	Longitud	0.25300	0.25520	0.25440	0.25540	0.25490	0.25490	0.25460	0.25460	0.2518	
	Nivel	0.948	0.953	0.953	0.953	0.951	0.951	0.951	0.951	0.95	
0.7	Longitud	0.20180	0.20410	0.20370	0.20420	0.20380	0.20380	0.20360	0.20360	0.2023	
	Nivel	0.942	0.951	0.949	0.951	0.951	0.951	0.951	0.951	0.955	
0.8	Longitud	0.14204	0.14403	0.14410	0.14413	0.14383	0.14384	0.14365	0.14365	0.14418	
	Nivel	0.935	0.938	0.941	0.939	0.938	0.938	0.938	0.938	0.953	
0.9	Longitud	0.07503	0.07618	0.07639	0.07623	0.07607	0.07607	0.07597	0.07597	0.07722	
	Nivel	0.938	0.942	0.942	0.942	0.941	0.941	0.941	0.941	0.946	

Tabla 9-3: Longitud y nivel de confianza de los intervalos. Tamaño de muestra 100

9.2. Inferencias para el coeficiente de correlación ρ

9.2.1. Verosimilitud Simplificada

Se tiene que la expresión de la verosimilitud es la siguiente:

$$L(\mu_1, \mu_2, \sigma_1^2, \sigma_2^2, \rho) = \left(\frac{1}{2\pi}\right)^n \left(\frac{1}{\sigma_1}\right)^n \left(\frac{1}{\sigma_2}\right)^n \left(\frac{1}{1-\rho^2}\right)^{\frac{n}{2}} \exp \left\{ -\frac{1}{2(1-\rho^2)} \sum_{i=1}^n \left[\left(\frac{x_i - \mu_1}{\sigma_1}\right)^2 - 2\rho \left(\frac{x_i - \mu_1}{\sigma_1}\right) \left(\frac{y_i - \mu_2}{\sigma_2}\right) + \left(\frac{y_i - \mu_2}{\sigma_2}\right)^2 \right] \right\} \quad (9-1)$$

Si llamamos A la parte de la sumatoria:

$$A = \sum_{i=1}^n \left[\left(\frac{x_i - \mu_1}{\sigma_1}\right)^2 - 2\rho \left(\frac{x_i - \mu_1}{\sigma_1}\right) \left(\frac{y_i - \mu_2}{\sigma_2}\right) + \left(\frac{y_i - \mu_2}{\sigma_2}\right)^2 \right] \quad (9-2)$$

Y a su vez, repartiendo la sumatoria en cada término:

$$\frac{1}{\sigma_1^2} \sum_{i=1}^n (x_i - \mu_1)^2 - \frac{2\rho}{\sigma_1\sigma_2} \sum_{i=1}^n (x_i - \mu_1)(y_i - \mu_2) + \frac{1}{\sigma_2^2} \sum_{i=1}^n (y_i - \mu_2)^2 \quad (9-3)$$

Digamos que:

$$B = \sum_{i=1}^n (x_i - \mu_1)^2 \quad D = \sum_{i=1}^n (x_i - \mu_1)(y_i - \mu_2) \quad C = \sum_{i=1}^n (y_i - \mu_2)^2 \quad (9-4)$$

Resolviendo B:

$$\sum_{i=1}^n (x_i - \mu_1)^2 = \sum_{i=1}^n (x_i - \bar{x} + \bar{x} - \mu_1)^2 = \sum_{i=1}^n [(x_i - \bar{x}) + (\bar{x} - \mu_1)]^2 \quad (9-5)$$

$$= \sum_{i=1}^n [(x_i - \bar{x})^2 + 2(x_i - \bar{x})(\bar{x} - \mu_1) + (\bar{x} - \mu_1)^2] \quad (9-6)$$

$$= \sum_{i=1}^n (x_i - \bar{x})^2 + 2(\bar{x} - \mu_1) \sum_{i=1}^n (x_i - \bar{x}) + n(\bar{x} - \mu_1)^2 \quad (9-7)$$

$$= (n-1)S_1^2 + 2(\bar{x} - \mu_1) \sum_{i=1}^n (x_i - n\bar{x}) + n(\bar{x} - \mu_1)^2 \quad (9-8)$$

$$= (n-1)S_1^2 + 2(\bar{x} - \mu_1)(n\bar{x} - n\bar{x}) + n(\bar{x} - \mu_1)^2 \quad (9-9)$$

$$= (n-1)S_1^2 + n(\bar{x} - \mu_1)^2 \quad (9-10)$$

Y haciendo el mismo procedimiento para C, se tiene que:

$$\sum_{i=1}^n (y_i - \mu_2)^2 = (n-1)S_2^2 + n(\bar{y} - \mu_2)^2 \quad (9-11)$$

Simplificando D:

$$\sum_{i=1}^n (x_i - \mu_1)(y_i - \mu_2) = \sum_{i=1}^n (x_i - \bar{x} + \bar{x} - \mu_1)(y_i - \bar{y} + \bar{y} - \mu_2) \quad (9-12)$$

$$= \sum_{i=1}^n [(x_i - \bar{x}) + (\bar{x} - \mu_1)][(y_i - \bar{y}) + (\bar{y} - \mu_2)] \quad (9-13)$$

$$= \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) + \sum_{i=1}^n (x_i - \bar{x})(\bar{y} - \mu_2) \quad (9-14)$$

$$+ \sum_{i=1}^n (y_i - \bar{y})(\bar{x} - \mu_1) + \sum_{i=1}^n (\bar{x} - \mu_1)(\bar{y} - \mu_2) \quad (9-15)$$

$$= \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) + (\bar{y} - \mu_2) \sum_{i=1}^n (x_i - \bar{x}) \quad (9-16)$$

$$+ (\bar{x} - \mu_1) \sum_{i=1}^n (y_i - \bar{y}) + n(\bar{x} - \mu_1)(\bar{y} - \mu_2) \quad (9-17)$$

$$= \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) + (\bar{y} - \mu_2)(n\bar{x} - n\bar{x}) \quad (9-18)$$

$$+ (\bar{x} - \mu_1)(n\bar{y} - n\bar{y}) + n(\bar{x} - \mu_1)(\bar{y} - \mu_2) \quad (9-19)$$

$$= (n-1)S_{12} + n(\bar{x} - \mu_1)^2(\bar{y} - \mu_2)^2 \quad (9-20)$$

Por tanto la función de verosimilitud queda de la siguiente manera:

$$L(\mu_1, \mu_2, \sigma_1^2, \sigma_2^2, \rho) = \left(\frac{1}{2\pi}\right)^n \left(\frac{1}{\sigma_1}\right)^n \left(\frac{1}{\sigma_2}\right)^n \left(\frac{1}{1-\rho^2}\right)^{\frac{n}{2}} \exp\left\{-\frac{1}{2(1-\rho^2)}A\right\} \quad (9-21)$$

donde A es:

$$A = \frac{1}{\sigma_1^2} [(n-1)S_1^2 + n(\bar{x} - \mu_1)^2] - \frac{2\rho}{\sigma_1\sigma_2} [(n-1)S_{12} + n(\bar{x} - \mu_1)(\bar{y} - \mu_2)] + \frac{1}{\sigma_2^2} [(n-1)S_2^2 + n(\bar{y} - \mu_2)^2] \quad (9-22)$$

9.2.2. Distribuciones condicionales

Para μ_1

Solo se considera la parte que contiene a μ_1 en A:

$$\begin{aligned} & \frac{n(\bar{x} - \mu_1)^2}{\sigma_1^2} - \frac{2\rho}{\sigma_1\sigma_2} n(\bar{x} - \mu_1)(\bar{y} - \mu_2) = \\ & \frac{n}{\sigma_1^2} \left[(\bar{x} - \mu_1)^2 - \frac{2\rho\sigma_1}{\sigma_2} (\bar{x} - \mu_1)(\bar{y} - \mu_2) \right] = \\ \frac{n}{\sigma_1^2} \left[(\bar{x} - \mu_1)^2 - \frac{2\rho\sigma_1}{\sigma_2} (\bar{x} - \mu_1)(\bar{y} - \mu_2) + \left[\frac{\rho\sigma_1}{\sigma_2} (\bar{y} - \mu_2) \right]^2 - \left[\frac{\rho\sigma_1}{\sigma_2} (\bar{y} - \mu_2) \right]^2 \right] &= \\ & \frac{n}{\sigma_1^2} \left[(\bar{x} - \mu_1)^2 - \frac{2\rho\sigma_1}{\sigma_2} (\bar{x} - \mu_1)(\bar{y} - \mu_2) + \left[\frac{\rho\sigma_1}{\sigma_2} (\bar{y} - \mu_2) \right]^2 \right] = \\ & \frac{n}{\sigma_1^2} \left[(\bar{x} - \mu_1) - \frac{\rho\sigma_1}{\sigma_2} (\bar{y} - \mu_2) \right]^2 = \\ & \frac{n}{\sigma_1^2} \left[- \left(\mu_1 - \bar{x} + \frac{\rho\sigma_1}{\sigma_2} (\bar{y} - \mu_2) \right) \right]^2 = \\ & \frac{n}{\sigma_1^2} \left[\mu_1 - \left(\bar{x} - \frac{\rho\sigma_1}{\sigma_2} (\bar{y} - \mu_2) \right) \right]^2 \end{aligned}$$

Esto último es el kernel de una distribución normal.

$$p(\mu_1 | \mu_2, \sigma_1^2, \sigma_2^2, \rho, \text{Datos}) \propto \exp\left\{-\frac{1}{2(1-\rho^2)} \frac{n}{\sigma_1^2} \left[\mu_1 - \left(\bar{x} - \rho \frac{\sigma_1}{\sigma_2} (\bar{y} - \mu_2) \right) \right]^2\right\} \quad (9-23)$$

Luego:

$$(\mu_1 | \mu_2, \sigma_1^2, \sigma_2^2, \rho, \text{Datos}) \sim N\left(\bar{x} - \rho \frac{\sigma_1}{\sigma_2} (\bar{y} - \mu_2), \frac{\sigma_1^2(1-\rho^2)}{n}\right) \quad (9-24)$$

Para σ_1^2

Observando la parte A, y descartando el resto, que no está en términos de σ_1^2 :

$$\begin{aligned} \frac{1}{\sigma_1^2} [(n-1)S_1^2 + n(\bar{x} - \mu_1)^2] - \frac{2\rho}{\sigma_1\sigma_2} [(n-1)S_{12} + n(\bar{x} - \mu_1)(\bar{y} - \mu_2)] &= \\ \frac{1}{\sigma_1^2} [(n-1)S_1^2 + n(\bar{x} - \mu_1)^2] - \frac{1}{\sigma_1} \left[\frac{2\rho}{\sigma_2} [(n-1)S_{12} + n(\bar{x} - \mu_1)(\bar{y} - \mu_2)] \right] &= \end{aligned}$$

Digamos que:

$$B = (n-1)S_1^2 + n(\bar{x} - \mu_1)^2 \quad C = \left[\frac{2\rho}{\sigma_2} [(n-1)S_{12} + n(\bar{x} - \mu_1)(\bar{y} - \mu_2)] \right] \quad (9-25)$$

Por tanto:

$$p(\sigma_1^2 | \mu_2, \mu_1, \sigma_2^2, \rho, \text{Datos}) \propto \frac{1}{\sigma_1^{n+1}} \exp \left\{ -\frac{1}{2(1-\rho^2)} \left[\frac{1}{\sigma_1^2} B - \frac{1}{\sigma_1} C \right] \right\} \quad (9-26)$$

10 Códigos en R

10.1. Intervalos de Confianza parte clásica

```
# Programa para comparar la eficiencia de nueve intervalos de confianza para
# el coeficiente de correlación
# Asintótico
# Bootstrap
# Exacto
# Transformaciones
# Jayaratman
# Pivote Generalizado
# Este programa necesita la librería mvtnorm para generar
# los datos de la Normal multivariada
library(mvtnorm)

calcule.cor<-function(datos)cor(datos)[1,2]

#Intervalo Razón de Verosimilitud
intervaloLR<-function(r,n){
log.L<-function(rho,r,n){

#Cálculo Integral de l(rho) Observemos que el parámetro toma el valor de rho
res1<-integrate(function(x,rho1=rho,r1=r,n1=n)
(cosh(x)-rho1*r1)^(-(n1-1)),0,Inf)$value

#Cálculo Integral de l(r) Observemos que el parámetro toma el valor de r
res2<-integrate(function(x,rho1=r,r1=r,n1=n)
(cosh(x)-rho1*r1)^(-(n1-1)),0,Inf)$value

#Aplicamos lognatural a la función de verosimilitud
temp<-(n-1)/2*(log(1-rho^2)-log(1-r^2))+log(res1)-log(res2)

#R(rho)>=0.147, despeje y aplique ln. Queda esta expresión.
temp<-temp-log(0.147)
temp
}

#Raíz más pequeña en donde r(rho) es > 0.147
LI<-uniroot(log.L,c(-0.99999999,r),r=r,n=n)$root

#Raíz más grande en donde r(rho) es > 0.147
LS<-uniroot(log.L,lower=r,upper=0.999999,r=r,n=n)$root
res<-c(LI,LS)
res
}

#Intervalo Bootstrap
#Función para calcular el Intervalo Bootstrap
intervalo.bootstrap<-function(n,medias,varcov,Nboot=1000){
```

```

#Se generan 1000 muestras de una normal bivariada de tamaño n
temp<-rmvnorm(n*Nboot,medias,varcov)
res<-NA
for(j in 1:Nboot){
  indices<-(n*(j-1)+1):(j*n)
  res<-c(res,calcule.cor(temp[indices,]))
}
res<-res[-1]

#Intervalo Bootstrap, cuantiles del histograma
res<-quantile(res,probs=c(0.025,0.975))
res
}

#Intervalo Transformación Arcotangente Original
#Función para calcular el I.C Transformación Arco Tangente Hiperb.
intervalo.aprox<-function(n,r){
  res<-c(tanh(atanh(r)-1.96/sqrt(n-3)),tanh(atanh(r)+1.96/sqrt(n-3)))
  res
}

#Intervalo Transformación Arcotangente 2
intervalo.aprox.z1<-function(n,r){
  z<-atanh(r)
  LI.1<-z-(((7*z)+r)/(8*(n-1)))-(1.96/sqrt(n-1))
  LS.1<-z-(((7*z)+r)/(8*(n-1)))+(1.96/sqrt(n-1))
  LI.2<-((8*LI.1*(n-1)+r)/((8*(n-1))-7))
  LS.2<-((8*LS.1*(n-1)+r)/((8*(n-1))-7))
  LI<-tanh(LI.2)
  LS<-tanh(LS.2)
  res<-c(LI,LS)
  res
}

#Intervalo Transformación Arcotangente 3
intervalo.aprox.z3<-function(n,r){
  z<-atanh(r)
  LI.1<-z-(((3*z)+r)/(4*(n-1)))-(1.96/sqrt(n-1))
  LS.1<-z-(((3*z)+r)/(4*(n-1)))+(1.96/sqrt(n-1))
  LI.2<-((4*LI.1*(n-1)+r)/((4*(n-1))-3))
  LS.2<-((4*LS.1*(n-1)+r)/((4*(n-1))-3))
  LI<-tanh(LI.2)
  LS<-tanh(LS.2)
  res<-c(LI,LS)
  res
}

#Intervalo Transformación Arcotangente 4
intervalo.aprox.z2<-function(n,r){
  z<-atanh(r)
  z2<-z-(((7*z)+r)/(8*(n-1)))-(((119*z)+(57*r)+(3*(r^2)))/(384*((n-1)^2)))
  LI.1<-z2-(1.96/sqrt(n-1))
  LS.1<-z2+(1.96/sqrt(n-1))
  a1<-((384*((n-1)^2)*LI.1)+(57*r)+(3*(r^2)))
  a2<-((384*((n-1)^2)*LS.1)+(57*r)+(3*(r^2)))
  b<-((3072*((n-1)^2))-((2688*(n-1))-952))
  LI.2<-((8*a1)+(384*(n-1)*r))/b
  LS.2<-((8*a2)+(384*(n-1)*r))/b
  LI<-tanh(LI.2)
  LS<-tanh(LS.2)
  res<-c(LI,LS)
}

```

```

res
}

#Intervalo Transformación Arcotangente 5
intervalo.aprox.z4<-function(n,r){
z<-atanh(r)
z4<-z-(((3*z)+r)/(4*(n-1)))-(((23*z)+(33*r)-(5*(r^2)))/(96*((n-1)^2)))
LI.1<-z4-(1.96/sqrt(n-1))
LS.1<-z4+(1.96/sqrt(n-1))
a1<-(384*((n-1)^2)*LI.1)+(96*(n-1)*r)+(132*r)-(20*(r^2))
a2<-(384*((n-1)^2)*LS.1)+(96*(n-1)*r)+(132*r)-(20*(r^2))
b<-(384*((n-1)^2)-(288*(n-1))-92
LI.2<-a1/b
LS.2<-a2/b
LI<-tanh(LI.2)
LS<-tanh(LS.2)
res<-c(LI,LS)
res
}

#Intervalo Jayaratnam
intervalo.J<-function(n,r){
t.alfa<-qt(0.975,n-2)
w<-(t.alfa/sqrt(n-2))/sqrt(1+((t.alfa)^2/(n-2)))
res<-c((r-w)/(1-(r*w)),(r+w)/(1+(r*w)))
res
}

#Intervalo Generalizado
intervalo.G<-function(n,ro){
r<-ro/sqrt(1-(ro^2))
resultado<-NULL
for(i in 1:1000){
z<-rnorm(1)
u1<-rchisq(1,n-1)
u2<-rchisq(1,n-2)
qi<-((r*sqrt(u2))-z)/sqrt(((r*sqrt(u2))-z)^2+u1)
resultado<-c(resultado,qi)
}
#Intervalo, cuantiles del histograma
resultado<-quantile(resultado,probs=c(0.025,0.975))
resultado
}

simulacion<-function(ro,n,Nsim=1000){
resultado<-rep(NA,18)

for(i in 1:Nsim){

#Se generan datos de una normal biv. con matriz de var-cov especificada
datos<-rmvnorm(n,c(0,0),matrix(c(1,ro,ro,1),ncol=2))

#Cálculo del vector de medias de los datos aleatorios
medias<-unlist(apply(datos,2,mean))

#Matriz de var-cov de los datos aleatorios
varcov<-cov(datos)

#Intervalo Bootstrap
res1<-intervalo.bootstrap(n,medias,varcov)

```

```
#Longitud del intervalo
L1<-res1[2]-res1[1]

#Verificación para el nivel de confianza
if(ro<res1[2] & ro>res1[1]) cae1<-1
else cae1<-0

r<-calcule.cor(datos)

#Intervalo Aproximado
res2<-intervalo.aprox(n,r)

#Longitud del intervalo
L2<-res2[2]-res2[1]
if(ro<res2[2] & ro>res2[1]) cae2<-1
else cae2<-0

#Intervalo LikelihoodRatio
res3<-intervaloLR(r,n)

#Longitud del intervalo
L3<-res3[2]-res3[1]
if(ro<res3[2] & ro>res3[1]) cae3<-1
else cae3<-0

#Intervalo Jayaratnam
res4<-intervalo.J(n,r)

#Longitud del intervalo
L4<-res4[2]-res4[1]
if(ro<res4[2] & ro>res4[1]) cae4<-1
else cae4<-0

#Intervalo Aproximado Transformación 1
res5<-intervalo.aprox.z1(n,r)

#Longitud del intervalo
L5<-res5[2]-res5[1]
if(ro<res5[2] & ro>res5[1]) cae5<-1
else cae5<-0

#Intervalo Aproximado Transformación 2
res6<-intervalo.aprox.z2(n,r)

#Longitud del intervalo
L6<-res6[2]-res6[1]
if(ro<res6[2] & ro>res6[1]) cae6<-1
else cae6<-0

#Intervalo Aproximado Transformación 3
res7<-intervalo.aprox.z3(n,r)

#Longitud del intervalo
L7<-res7[2]-res7[1]
if(ro<res7[2] & ro>res7[1]) cae7<-1
```

```

else cae7<-0

#Intervalo Aproximado Transformación 4
res8<-intervalo.aprox.z4(n,r)

#Longitud del intervalo
L8<-res8[2]-res8[1]
if(ro<res8[2] & ro>res8[1]) cae8<-1
else cae8<-0

#Intervalo Generalizado
res9<-intervalo.G(n,r)

#Longitud del intervalo
L9<-res9[2]-res9[1]

#Verificación para el nivel de confianza
if(ro<res9[2] & ro>res9[1]) cae9<-1
else cae9<-0

resultado<-rbind(resultado,c(L1,cae1,L2,cae2,L3,cae3,L4,cae4,
L5,cae5,L6,cae6,L7,cae7,L8,cae8,L9,cae9))

#fin del for
}
resultado<-resultado[-1,]
resultado
}

imprimir.linea<-function(x,ro){
l1<-substr(x[3,1],9,16)
l2<-substr(x[3,3],9,16)
l3<-substr(x[3,5],9,16)
l4<-substr(x[3,7],9,16)
l5<-substr(x[3,9],9,16)
l6<-substr(x[3,11],9,16)
l7<-substr(x[3,13],9,16)
l8<-substr(x[3,15],9,16)
l9<-substr(x[3,17],9,16)
n1<-substr(x[4,2],9,16)
n2<-substr(x[4,4],9,16)
n3<-substr(x[4,6],9,16)
n4<-substr(x[4,8],9,16)
n5<-substr(x[4,10],9,16)
n6<-substr(x[4,12],9,16)
n7<-substr(x[4,14],9,16)
n8<-substr(x[4,16],9,16)
n9<-substr(x[4,18],9,16)

cat(" \multirow{2}{*}{ } ",ro,' & ','Longitud',' & ',l1,' & ',l2,' & ',l3,' & ',l4,
', ' & ',l5,' & ',l6,' & ',l7,' & ',l8,' & ',l9, " \\ ",'\n',' & ','Nivel',' & ',n1,
' & ',n2,' & ',n3,' & ',n4,' & ',n5,' & ',n6,' & ',n7,' & ',n8,' & ',n9," \\ \hline ",'\n')
}

```


10.2. Error cuadrático medio de los estimadores de ρ

```

#Para sacar los MSE de los estimadores de rho
Error.cuad.medio1<-function(ro,n,Nsim=1000){
  resultado<-NULL
  for(i in 1:Nsim){
    datos<-rmvnorm(n,c(0,0),matrix(c(1,ro,ro,1),ncol=2))
    r<-calcule.cor(datos)
    resultado<-c(resultado,r)
  }
  valor.esp<-mean(resultado)
  varianza<-var(resultado)
  B<-valor.esp-ro
  ECM<-varianza+(B^2)
  return(ECM)
}

Error.cuad.medio2<-function(ro,n,Nsim=1000){
  UNO<-gamma((n-2)/2)
  DOS<-gamma(1/2)
  TRES<-gamma((n-3)/2)
  resultado<-NULL
  for(i in 1:Nsim){
    datos<-rmvnorm(n,c(0,0),matrix(c(1,ro,ro,1),ncol=2))
    r<-calcule.cor(datos)
    integrand <- function(x,n,r) {((x^(-1/2))*((1-x)^((n-5)/2)))/sqrt(1-(x*(1-(r^2))))}
    integ<-integrate(integrand,lower=0,upper=1,n,r)$value
    umvue<-r*(UNO/(DOS*TRES))*integ
    resultado<-c(resultado,umvue)
  }
  valor.esp<-mean(resultado)
  varianza<-var(resultado)
  B<-valor.esp-ro
  ECM<-varianza+(B^2)
  return(ECM)
}

#FUNCIÓN DE GRÁFICAS Y RESULTADOS
#Muestras de tamaño 5,10,20,30
#i y j son los elementos del vector de rho a graficar

graficas.ECM<-function(n,ro,i,j){
  par(mfrow=c(1,2))
  for(k in i:j){
    ECM1<-NULL
    ECM2<-NULL
    for(v in 1:4){
      ecm1<-Error.cuad.medio1(ro[k],n[v])
      ECM1<-c(ECM1,ecm1)
      ecm2<-Error.cuad.medio2(ro[k],n[v])
      ECM2<-c(ECM2,ecm2)
    }
    plot(n,ECM1,type='l',col=3,ylab='ECM',xlab='n',lwd=2,ylim=c(0.0,0.35))
    lines(n,ECM2,col='blue',lwd=2)
    legend("topright",c("R", "UMVUE"),lty=c(1,1),col=c(3,4),lwd=c(2,2))
    title(substitute(rho == valores, list(valores=ro[k])))
    print(as.matrix(c(ECM1[1],ECM2[1],ECM1[2],ECM2[2],ECM1[3],ECM2[3],ECM1[4],ECM2[4])))
  }
}

```

```
#Ensayo Función de gráficas
muestras<-c(5,10,20,30)
rhos<-seq(0.0,0.9,by=0.1)

graficas.ECM(muestras,rhos,1,2)
```

10.3. Intervalo de confianza Bayesiano con la apriori 2

```
#Aposteriori para rho utilizando apriori de la forma: 1/[(Sigma1^2)(Sigma2^2)(1-rho^2)]
```

```
#Datos
library(mvtnorm)
#####Correlación muestral#####
calcule.cor<-function(datos)cor(datos)[1,2]

#####Función para generar Mu1#####
muestra.mu1<-function(mu2,xbarra,ybarra,sigma1,sigma2,ro,n){
media.normal.1<-xbarra-(ro*(sigma1/sigma2)*(ybarra-mu2))
var.normal.1<-(sigma1^2)*(1-ro^2)/n
mu1<-rnorm(1,mean=media.normal.1,sd=sqrt(var.normal.1))
mu1
}

#####3#Función para generar Mu2#####
muestra.mu2<-function(xbarra,ybarra,sigma1,sigma2,ro,mu1,n){
media.normal.2<-ybarra-(ro*(sigma2/sigma1)*(xbarra-mu1))
var.normal.2<-(sigma2^2)*(1-ro^2)/n
mu2<-rnorm(1,mean=media.normal.2,sd=sqrt(var.normal.2))
mu2
}

#####Función para generar Sigma1^2#####
muestra.sigma1<-function(n,ro,xbarra,ybarra,mu1,mu2,s11,s12,sigma2){
sigmas1<-matrix(seq(0,5,length=100),ncol=1)
probco1<-matrix(0,ncol=1,nrow=length(sigmas1))
operación<-function(sigmas1,n,ro,xbarra,ybarra,mu1,mu2,s11,s12,sigma2) {
B<-((n-1)*s11)+(n*(xbarra-mu1)^2)
C<-((2*ro/sigma2)*(((n-1)*s12)+(n*(xbarra-mu1)*(ybarra-mu2))))
dis<-((1/(sigmas1^(n+1)))*exp(-(1/(2*(1-ro^2)))))*(((1/(sigmas1^2))*B)-((1/sigmas1)*C)))
return(dis)}
probco1<-apply(sigmas1,1,operación,n,ro,xbarra,ybarra,mu1,mu2,s11,s12,sigma2)
proba<-ifelse(is.na(probco1),0,probco1)
res0<-sample(sigmas1,1,prob=proba)
res0 }

#####Función para generar Sigma2^2#####
muestra.sigma2<-function(n,ro,xbarra,ybarra,mu1,mu2,s22,s12,sigma1){
sigmas2<-matrix(seq(0,6,length=100),ncol=1)
probco1<-matrix(0,ncol=1,nrow=length(sigmas2))
operación<-function(sigmas2,n,ro,xbarra,ybarra,mu1,mu2,s22,s12,sigma2) {
D<-((n-1)*s22)+(n*(ybarra-mu2)^2)
E<-((2*ro/sigma1)*(((n-1)*s12)+(n*(xbarra-mu1)*(ybarra-mu2))))
dis<-((1/(sigmas2^(n+1)))*exp(-(1/(2*(1-ro^2)))))*(((1/(sigmas2^2))*D)-((1/sigmas2)*E)))
return(dis)}
probco1<-apply(sigmas2,1,operación,n,ro,xbarra,ybarra,mu1,mu2,s22,s12,sigma1)
proba<-ifelse(is.na(probco1),0,probco1)
res0<-sample(sigmas2,1,prob=proba)
res0 }
```

```
#####Función para generar rho#####
muestra.ro<-function(n,xbarra,ybarra,mu1,mu2,s22,s12,s11,sigma1,sigma2){
ros<-matrix(seq(-1,1,length=100),ncol=1)
probco1=matrix(0,ncol=1,nrow=length(ros))
operación<-function(ros,n,xbarra,ybarra,mu1,mu2,s22,s12,s11,sigma1,sigma2) {
A<-((1/(sigma1^2))*(((n-1)*s11)+(n*(xbarra-mu1)^2))) - ((2*ros)/(sigma1*sigma2))
*(((n-1)*s12)+(n*(xbarra-mu1)*(ybarra-mu2)))) + ((1/(sigma2^2))*(((n-1)*s22)+
(n*((ybarra-mu2)^2))))
dis<-(1/((1-(ros^2))^(n/2+1)))*exp(-1/(2*(1-(ros^2))))*A
return(dis)}
probco1<-apply(ros,1,operación,n,xbarra,ybarra,mu1,mu2,s22,s12,s11,sigma1,sigma2)
proba<-ifelse(is.na(probco1),0,probco1)
res0<-sample(ros,1,prob=proba)
res0 }
#####

#####Una cadena del Gibbs#####

paso<-function(xbarra,ybarra,s11,s22,s12,n,mu1,mu2,sigma1,sigma2,r){
resultados<-c(xbarra,ybarra,s11,s22,s12,n)

mu1.n<-muestra.mu1(mu2,xbarra,ybarra,sigma1,sigma2,r,n)
resultados<-c(resultados,mu1.n)

mu2.n<-muestra.mu2(xbarra,ybarra,sigma1,sigma2,r,mu1.n,n)
resultados<-c(resultados,mu2.n)

sigma1.n<-muestra.sigma1(n,r,xbarra,ybarra,mu1.n,mu2.n,s11,s12,sigma2)
resultados<-c(resultados,sigma1.n)

sigma2.n<-muestra.sigma2(n,r,xbarra,ybarra,mu1.n,mu2.n,s22,s12,sigma1.n)
resultados<-c(resultados,sigma2.n)

ro.n<-muestra.ro(n,xbarra,ybarra,mu1.n,mu2.n,s22,s12,s11,sigma1.n,sigma2.n)
resultados<-c(resultados,ro.n)

return(resultados)
}

#####Cálculo de un intervalo, con una muestra de tamaño 7000#####
#Esto me da una muestra de 1000 valores de los parámetros.
intervalo.apriori2<-function(n,ro){

#Se generan datos de una normal biv. con matriz de var-cov especificada
datos<-rmvnorm(n,c(0,0),matrix(c(1,ro,ro,1),ncol=2))

#Cálculo del vector de medias de los datos aleatorios (medias muestrales)
xbarra<-unlist(apply(datos,2,mean))[1]

#Cálculo del vector de medias de los datos aleatorios (medias muestrales)
ybarra<-unlist(apply(datos,2,mean))[2]

#Cálculo de la matriz de var-cov muestral
```

```

varcov<-cov(datos)
r<-calcule.cor(datos)

sigma1<-1
sigma2<-1
mu1<-0
mu2<-0
s11<-varcov[1,1]
s22<-varcov[2,2]
s12<-varcov[1,2]

#En particular se obtiene un solo IC para rho.
x<-c(xbarra,ybarra,s11,s22,s12,n,mu1,mu2,sigma1,sigma2,r)
muestras<-sapply(1:1000,function(i)
{x<-paso(x[1],x[2],x[3],x[4],x[5],x[6],x[7],x[8],x[9],x[10],x[11])})
muestras.organizado<-matrix(data=muestras,ncol=11,byrow=T)
quemado<-muestras.organizado[-(1:3000), ]
mu1.sim<-quemado[,7]
mu2.sim<-quemado[,8]
sigma1.sim<-quemado[,9]
sigma2.sim<-quemado[,10]
ro.sim<-quemado[,11]

#Intervalo, cuantiles del histograma
resultado<-quantile(ro.sim,probs=c(0.025,0.975))
return(resultado)

} #Fin de una sola corrida (una sola muestra de 7000 ros, cálculo de un solo intervalo)

#####Simulación de los 1000 intervalos, obteniendo las longitudes y las coberturas#####

simulacion<-function(ro,n,Nsim){
res10<-NULL
resultado<-NULL
enes<-matrix(data=rep(n,Nsim),ncol=1)
res10<-apply(enes,1,intervalo.apriori2,ro)

#Longitud del intervalo
long<-function(x,ro){
long<-x[2]-x[1]
long }

#Cobertura del intervalo
cobert<-function(x,ro){
if(ro<x[2] & ro>x[1]) cae10<-1
else cae10<-0
cae10 }

res1<-apply(res10,2,long,ro)
res2<-apply(res10,2,cobert,ro)

resultado<-cbind(res1,res2)
return(resultado) }

#Para organizar los resultados
imprimir.linea<-function(x){

```

```
l10<-substr(x[3,1],9,16)
n10<-substr(x[4,2],9,16)
res<-c(l10,n10)
res
}
```

```
#Función para calcular intervalos con diferentes escenarios
intervalos<-function(rho,n){
res<-imprimir.linea(summary(simulacion(rho,n,1000)))
return(res)}
```

```
#Ensayo con los rhos que se estudian y muestra tamaño 5
```

```
rhos<-seq(0,0.9,by=0.1)
```

```
muestra5<-sapply(rhos,intervalos,5)
muestra5
```

10.4. Índices de resumen

```
#Función para realizar los cálculos del Índice de resumen
suma<-function(x,y){(2-x)*(y/0.95)}
```

```
#Función para organizar directamente una matriz de índices a partir de la
#lectura de las matrices de longitud y nivel real en cada tamaño de muestra.
indice<-function(x,y){
matriz<-matrix(data=mapply(suma,x,y),ncol=10,byrow=F)
return(matriz)
}
```

```
i1<-t(indice(l1,m1))
colnames(i1)<-c("Bootstrap", "ArcTanh", "LR", "Jayaratnam", "Z1", "Z2", "Z3", "Z4", "P.G")
rownames(i1)<-c(0.0,0.1,0.2,0.3,0.4,0.5,0.6,0.7,0.8,0.9)
i1
```

```
i2<-t(indice(l2,m2))
colnames(i2)<-c("Bootstrap", "ArcTanh", "LR", "Jayaratnam", "Z1", "Z2", "Z3", "Z4", "P.G")
rownames(i2)<-c(0.0,0.1,0.2,0.3,0.4,0.5,0.6,0.7,0.8,0.9)
i2
```

```
i3<-t(indice(l3,m3))
colnames(i3)<-c("Bootstrap", "ArcTanh", "LR", "Jayaratnam", "Z1", "Z2", "Z3", "Z4", "P.G")
rownames(i3)<-c(0.0,0.1,0.2,0.3,0.4,0.5,0.6,0.7,0.8,0.9)
i3
```

```
i4<-t(indice(l4,m4))
colnames(i4)<-c("Bootstrap", "ArcTanh", "LR", "Jayaratnam", "Z1", "Z2", "Z3", "Z4", "P.G")
rownames(i4)<-c(0.0,0.1,0.2,0.3,0.4,0.5,0.6,0.7,0.8,0.9)
i4
```

```
i5<-t(indice(l5,m5))
colnames(i5)<-c("Bootstrap", "ArcTanh", "LR", "Jayaratnam", "Z1", "Z2", "Z3", "Z4", "P.G")
rownames(i5)<-c(0.0,0.1,0.2,0.3,0.4,0.5,0.6,0.7,0.8,0.9)
i5
```

10.5. Construcción de Clústers

Donde cada matriz “i” es la matriz de índices construidos para cada tamaño de muestra y método de Intervalo de Confianza.

```
#Clusters

distancia1<-dist(t(i1),method="euclidean")
cluster1<-hclust(distancia1)
plot(cluster1,main="Clusters para n=5")

distancia2<-dist(t(i2),method="euclidean")
cluster2<-hclust(distancia2)
plot(cluster2,main="Clusters para n=10")

distancia3<-dist(t(i3),method="euclidean")
cluster3<-hclust(distancia3)
plot(cluster3,main="Clusters para n=20")

distancia4<-dist(t(i4),method="euclidean")
cluster4<-hclust(distancia4)
plot(cluster4,main="Clusters para n=50")

distancia5<-dist(t(i5),method="euclidean")
cluster5<-hclust(distancia5)
plot(cluster5,main="Clusters para n=100")
```

Bibliografía

- J.M. Bernardo and A.F.M. Smith. *Bayesian Theory*. John Wiley & Sons, Ltd, 2000.
- B. Efron. Computers and theory of statistics: Thinking the unthinkable. *SIAM Review*, 21(4):460–480, 1979.
- R. Falk and A.D. Well. Many faces of the correlation coefficient. *Journal of Statistics Education*, 5(3), 1997.
- R.A Fisher. Frequency distribution of the values of the correlation coefficient in samples from an indefinitely large population. *Biometrika*, 10(4):507–521, 1915.
- R.A Fisher. On the “probable error” of a coefficient of correlation deduced from a small sample. *Metron*, 1:3–32, 1921.
- M. Ghosh, B. Mukherjee, U. Santra, and D. Kim. Bayesian and likelihood-based inference for the bivariate normal correlation coefficient. *Journal of Statistical Planning and Inference*, 140(6):1410–1416, 2010.
- F.A. Graybill. *Theory and Application of the Linear Model*. Duxbury Press: Boston, 1976.
- H. Hotelling. New light on the correlation coefficient and its transforms. *Journal of the Royal Statistical Society*, 15(2):193–232, 1953.
- S. Jackman. *Bayesian Analysis for the Social Sciences*. John Wiley & Sons, Ltd, 2009.
- R.A. Johnson and D.W. Wichern. *Applied Multivariate Statistical Analysis*. Pearson Education, Inc. Pearson Prentice Hall, 2007.
- J.G. Kalbfleish. *Probability and Statistical Inference*. Springer-Verlag: New York, 1985.
- I.K. Krishnamoorthy and Y. Xia. Inferences on correlation coefficients: One-sample, independent and correlated cases. *Journal of Statistical Planning and Inference*, 137(7): 2362–2379, 2007.
- P. McCullagh. Some statistical properties of a family of continuous univariate distributions. *Journal of the American Statistician Association*, 84(405):125–129, 1989.
- Y. Pawitan. *In All Likelihood*. Clarendon Press: Oxford, 2001.

-
- Q. Zheng and J.H. Matis. Correlation coefficient revisited. *The American Statistician*, 48 (3):240–241, 1994.