

Estadística no Paramétrica: Proyecto 1

Paula Rodriguez y Felipe González Casabianca

20 de marzo de 2017

Resumen

Este documento cuenta como entregable para el el proyecto 1 el curso: *Estadística no Paramétrica y y remuestreo*, dictado en el primer semestre del año 2017 por Adolfo Quiroz en la Universidad de los Andes.

Capítulo 1

Aproximación normal a la distribución de W^+

1.1. Introducción

En este capítulo se pretende evaluar la velocidad de convergencia de la aproximación normal de la distribución del estadístico Wilcoxon Signado W^+ cuando la muestra considerada es simétrica. Para esto se tiene en cuenta el *Teorema del Límite Central para U-estadísticos*. Note que para una muestra X_1, \dots, X_n el estadístico de Wilcoxon Signado está dado por:

$$\begin{aligned} W^+ &= \sum_{i=1}^n \Psi_i R_i^+ \\ &= \sum_{0 \leq i \leq j \leq n} \Psi(X_{(i)} - X_{(j)}) \\ &= \sum_{0 \leq i \leq j \leq n} \Psi(X_i - X_j) \\ &= \sum_{j=1}^n \Psi(2X_j) \sum_{0 \leq i < j \leq n} \Psi(X_i - X_j) \\ &= nU(X_1, \dots, X_n) + \binom{n}{2} U'(X_1, \dots, X_n) \end{aligned}$$

Donde U y U' son U-estadísticos que estiman $P(X_1 > 0)$ y $P(X_1 + X_2 > 0)$ respectivamente. Por lo tanto, por el *Teorema del Límite Central para U-estadísticos*, tenemos que la aproximación normal a considerar está dada por

$$\frac{\sqrt{3n}(W^+ - EW^+)}{\binom{n}{2}} \longrightarrow N(0, 1)$$

cuando $n \longrightarrow \infty$. Para evaluar su velocidad de convergencia se calculará el valor máximo de la diferencia entre la distribución real del estadístico estandarizado

y una normal estándar. Con esto se observará el tamaño de la muestra a partir del cuál esta diferencia es menor a 0.005.

1.2. Implementación

Como fue sugerido, la implementación se realizó en R (versión 3.3.2). No se usó ninguna librería adicional a la descarga estándar.

Se incluye el método que calcula la función de distribución acumulativa del estadístico propuesto, el código completo se encuentra en Anexos.¹

La función de distribución acumulativa, fue calculada entre: $[-6, 6]$ con un paso de 0,25. Note que a partir de las siguientes operaciones, se puede obtener la distribución de este estadístico a partir de la fda del estadístico W^+ :

$$\begin{aligned} F(t) &= P\left(\frac{\sqrt{3n}(W^+ - EW^+)}{\binom{n}{2}} \leq t\right) \\ &= P\left(W^+ \leq \frac{t\binom{n}{2} + EW^+}{\sqrt{3n}}\right) \\ &= F_{W^+}\left(\frac{t\binom{n}{2} + EW^+}{\sqrt{3n}}\right) \\ &= \sum_{i=0}^K \frac{c_n(i)}{2^n} \end{aligned}$$

donde K , es la parte entera de:

$$\frac{t\binom{n}{2} + EW^+}{\sqrt{3n}}$$

Por lo tanto, a partir de la función $c_n(k)$ dada, es posible obtener la distribución del estadístico propuesto.

```

1
2 #Calculates the proposed statistic acumulative probability function
3 #the vector will be from -6 to 6 steps of 0.25
4 proposed_statistic_acum_prob = function(n, prob_vec = NULL)
5 {
6   #PARAMETER
7   #-----
8   #n integer: correspondes to the number of samples taken into account
9   #           into the Signed Wilcoxon statistic
10  #prob_vec vector: the probability vector of the Wilcoxon Signed statistic
11  #                 so that it does not need to be coputed again
12
13  #calculates the distribution as a matrix multiplication

```

¹ El código está comentado en Inglés para facilitar la resolución de problemas en blogs de programación en la web

```

14 #recall that the distribution is the position sum of the probability vector
15 if(is.null(prob_vec))
16 {
17   prob_vec = prob_w_plus(n)
18 }
19
20 #gets the distribution of W+
21 d_w = acum_prob_w_plus(n, prob_vec)
22
23 #calculates the corresponding coordinates
24 N = series_sum(n)
25
26 index = give_quantiles()
27
28 #Transforms the indexes to fit the new distribution
29 index = index*(choose(n,2)/(sqrt(3*n)))
30 index = index + N/2
31 index = floor(index)
32
33 #shifts the indexes since R starts at 1
34 index = index + 1
35
36 #Calculates the distribution for the acceptable indexes
37 # The out of range indexes are left as zero or one (accordingly)
38 distribution = rep(0,length(index))
39 selected_index = which(index >0 & index <= length(d_w))
40 distribution[selected_index] = d_w[index[selected_index]]
41 distribution[index > length(d_w)] = 1
42
43 return(distribution)
44
45 }
46
47 #end of proposed_statistic_acum_prob

```

1.3. Resultados

A continuación se encuentran los resultados del experimento. Se incluye el gráfico de convergencia del estadístico propuesto y algunos gráficos Prob - Prob para valores ascendentes de la muestra.

Se probó la convergencia del estadístico propuesto a la distribución normal, iterando desde una muestra de tamaño 15 hasta una de tamaño 200. Dicho procedimiento tomó 6 min y 32s en un computador con un procesador de 3.1 GHz y dos núcleos.

Aproximation Convergence

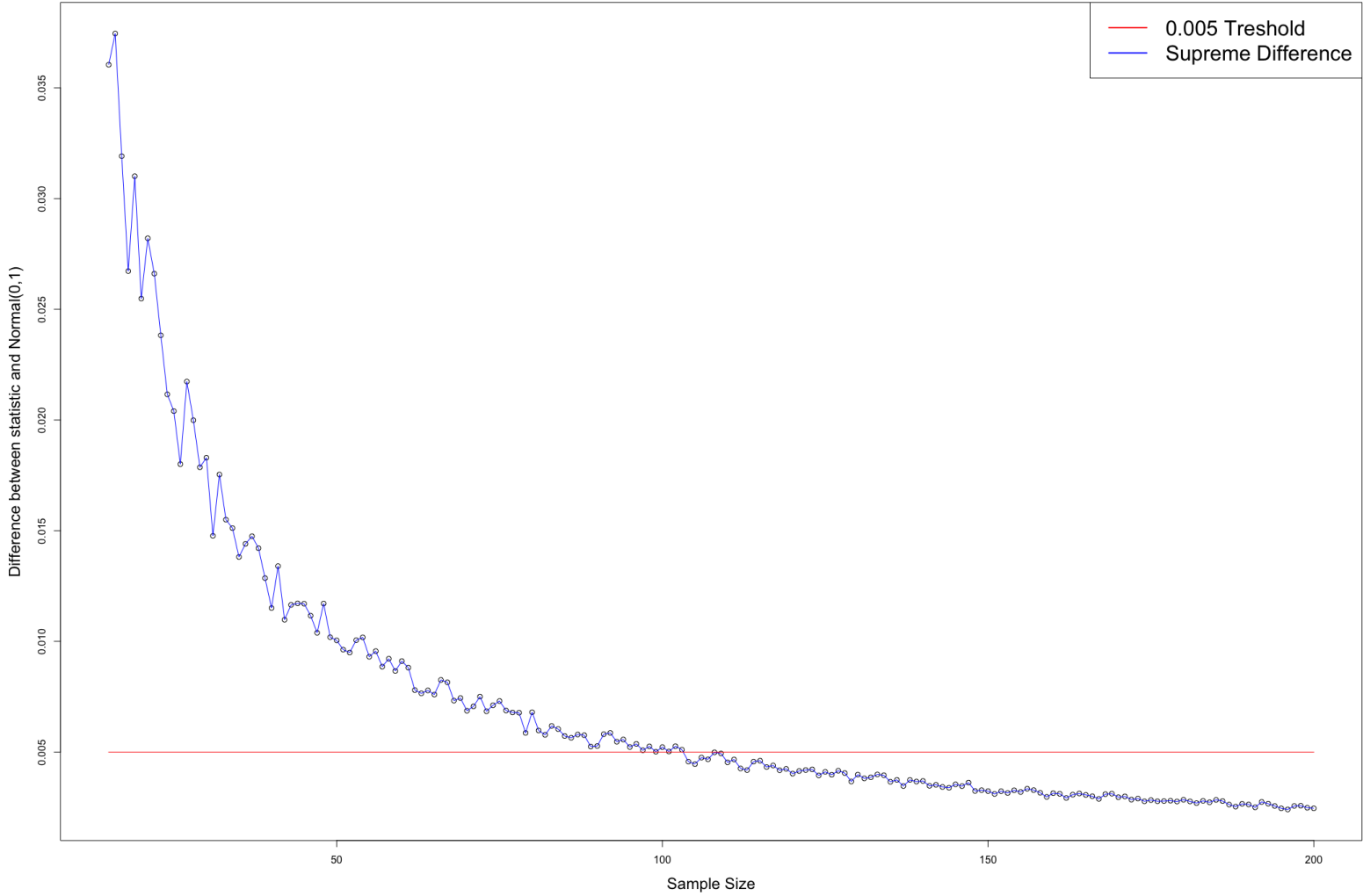


Figura 1.1: Convergencia del estadístico propuesto

Al graficar la diferencia máxima entre la distribución real del estadístico estandarizado y una normal estándar para distintos tamaños de muestra se puede observar que para tamaños de muestra mayores a 100 la diferencia máxima es menor al 0.5 %. Los siguientes prob-prob plots muestran la aproximación para cada uno de los tamaños de muestra considerados.

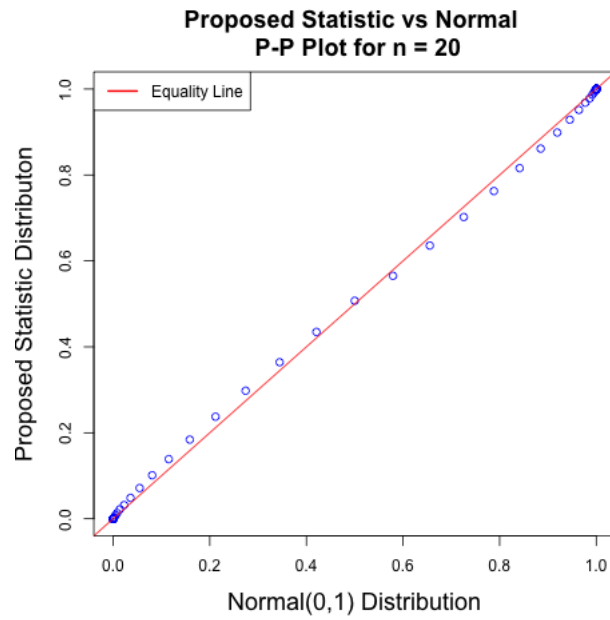


Figura 1.2: Prob-Prob entre las dos distribuciones para $n = 20$

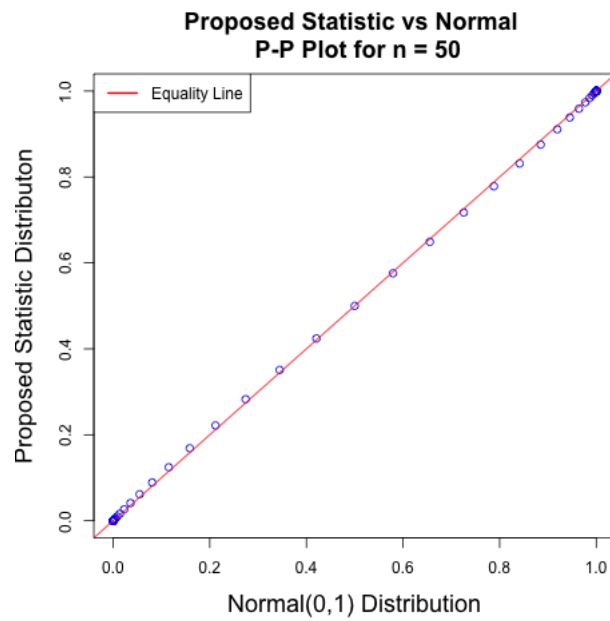


Figura 1.3: Prob-Prob entre las dos distribuciones para $n = 50$

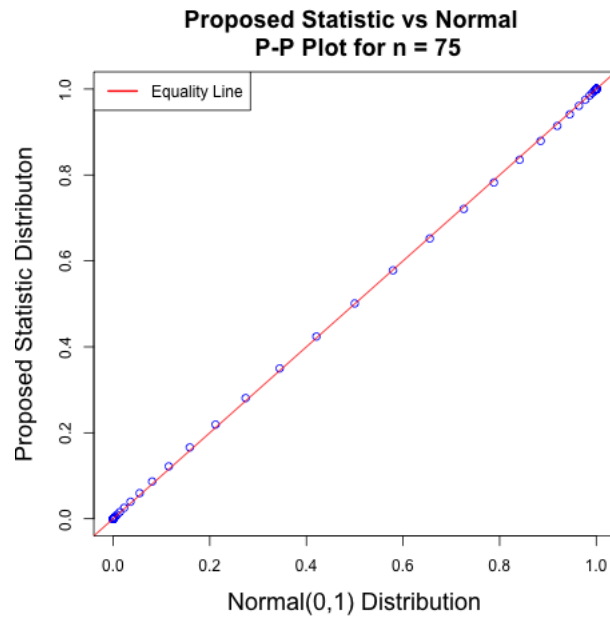


Figura 1.4: Prob-Prob entre las dos distribuciones para $n = 75$

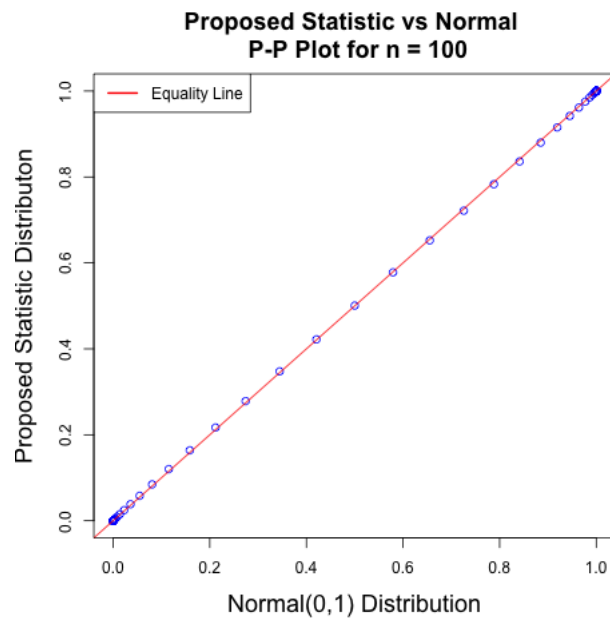


Figura 1.5: Prob-Prob entre las dos distribuciones para $n = 100$

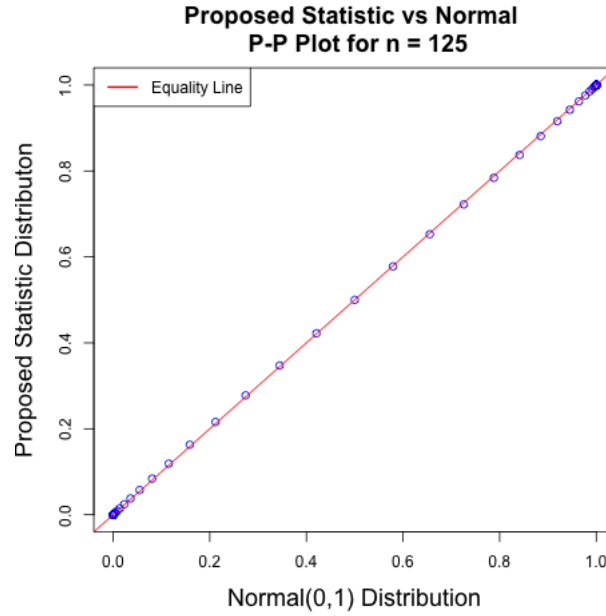


Figura 1.6: Prob-Prob entre las dos distribuciones para $n = 125$

Se puede ver que a partir de 50 datos la aproximación es bastante buena, pues los gráficos muestran que la probabilidad dada por la distribución del estadístico estandarizado se asemeja a la probabilidad dada por una normal estándar.²

1.4. Conclusiones

A partir del experimento realizado se puede concluir que la distribución del estadístico de Wilcoxon Signado estandarizado tiene una alta velocidad de convergencia a una normal estándar puesto que para muestras de tamaño mayor a 100 su diferencia es menor al 0.5 %. Incluso para muestras de tamaño mayor a 20 datos la diferencia es menor al 5 %, lo cuál evidencia una muy buena y rápida aproximación del estadístico.

²No es muy claro por qué se presenta el comportamiento de la prob-prob plot para la muestra de tamaño 20. Para esto se puede consultar el artículo *Goodness-of-fit tests based on P-P probability plots* de Gan y Koehler.

Capítulo 2

Eficiencia de estadísticos para pruebas sobre el centro de simetría

2.1. Introducción

En este capítulo se busca estudiar la eficiencia de diferentes estadísticos para una muestra con una distribución dada. Teniendo una muestra Z_1, \dots, Z_n con distribución simétrica centrada en θ , se evaluará la hipótesis nula $\theta = \theta_0$ contra la hipótesis nula $\theta \neq \theta_0$ utilizando el estadístico t student, el estadístico de signos y el estadístico de rangos signados W^+ . Esto, con el propósito de observar cuál estadístico es preferible para la prueba de hipótesis considerada teniendo en cuenta distintas distribuciones simétricas y tamaños de muestra.

2.2. Implementación

Al igual que el ítem anterior, la implementación se realizó en R (versión 3.3.2). Se adicionaron las librerías:

- **pracma**: Para comparación de Strings (Lectura de parámetros)
- **rmuti**: Para utilizar la distribución Laplace
- **ggplot2**: Para graficar con más colores que magenta

Para ejecutar el experimento se construyeron 3 métodos auxiliares:

- ***generateSample(n, dist, θ)*** : este método devuelve una muestra de tamaño n centrada en θ con distribución: *dist*. Donde este último debe ser alguno de los siguientes *Strings* :

$\{normal, uniform, laplace, student, cauchy\}$

- ***getStatistic*(*Z*, *statName*, θ_0)**: este método devuelve un estadístico: *statName* sobre la muestra *Z*, utilizando el valor de θ_0 . Aca, *statName* puede ser:

$\{student, signed, rank\}$

- ***getPValue*(ω , *n*, *statName*)**: este método devuelve la probabilidad de que el estadístico *statName* sea mayor o igual a ω , en una muestra de tamaño *n* y asumiendo la hipótesis nula ($\theta = \pi$)

Luego, se construyó un *data.frame* con todas las combinaciones posibles, usando la función: *expand.grid* de los siguientes parámetros:

- $\theta = \{\pi + 0,1, \pi + 0,25, \pi + 0,5, \pi + 1\}$
- $n = \{20, 30, 50, 100\}$
- *statName* = $\{student, signed, rank\}$
- *dist* = $\{normal, uniform, laplace, student, cauchy\}$

Así, con la ayuda de los tres métodos auxiliares, se recorrió cada fila de esta estructura ejecutando 500 veces el caso correspondiente y calculando la potencia promedio. Esta forma de ejecutar el procedimiento (a diferencia de uno anidado) quita legibilidad al código pero permite tener un resultado unificado de las potencias (en el mismo *data.frame* de las combinaciones).

Se incluye el código del procedimiento central del experimento. El resto del código se incluye en Anexos.

```

1 #The main method for item 2
2 test_hypothesis = function(num_ite = 500,
3                             alpha = 0.05,
4                             thetas = c(pi + 0.1, pi + 0.25, pi + 0.5, pi + 1),
5                             sample_sizes = c(20,30,50,100))
6 {
7   #PARAMETER
8   #-----
9   # num_ite numeric: the number of iterations for each scenario
10  # alpha numeric: the desired confidence for the hypothesis test
11  # thetas numeric vector: a numeric vector with the desired centralities to
12    check. The method
13  # assumes the null hypothesis with theta = pi
14  # sample_size numeric vecor: a numeric vector with the desired sample sizes
15    to test
16
17  sample_types = c('normal','uniform','laplace','student','cauchy')
18  statistics = c('student','signed','rank')
19
20  #creates all the possible combinations
21  result = expand.grid(sample_sizes,thetas, sample_types,statistics)
22  #assignes colnames
23  colnames(result) = c('sample_size','theta','sample_dist','stat')
24  # Adds the power columns
25  result$power = rep(0,nrow(result))
26  #Adds the shift that the theta received (for reading purposes)
27  result$shift = result$theta - pi

```

```

26
27
28
29 #iterates over each row and excecutes experiment
30 for(i in 1:nrow(result))
31 {
32   row = result[i,]
33   print(row)
34
35   discards = 0
36
37   for(j in 1:num_ite)
38   {
39     #Sample
40     Z = generate_sample(n = row$sample_size, distribution = toString(row$
41       sample_dist), theta = row$theta)
42     #Statistic
43     omega = get_statistic(Z = Z, stat_name = toString(row$stat), theta_0 =
44       pi)
45     #P-Value
46     p_value = get_p_value(omega = omega, n = row$sample_size, stat_name =
47       toString(row$stat))
48     #Checks if the null hipotheiss is discarded
49     if(p_value <= alpha)
50     {
51       discards = discards + 1
52     }
53   }
54
55   #Assignes power
56   result[i,]$power = discards/num_ite
57 }
58
59 return(result)
60 }
61 #end of test_hipotesis

```

2.3. Resultados

La ejecución completa de todos los casos del experimento, tomó 42s en un computador con un procesador de 3.1 GHz y dos núcleos.

A continuación se presentan gráficas comparativas para los distintos experimentos. En éstas se puede ver la potencia de cada estadístico para cada distribución considerada y tamaño de muestra.

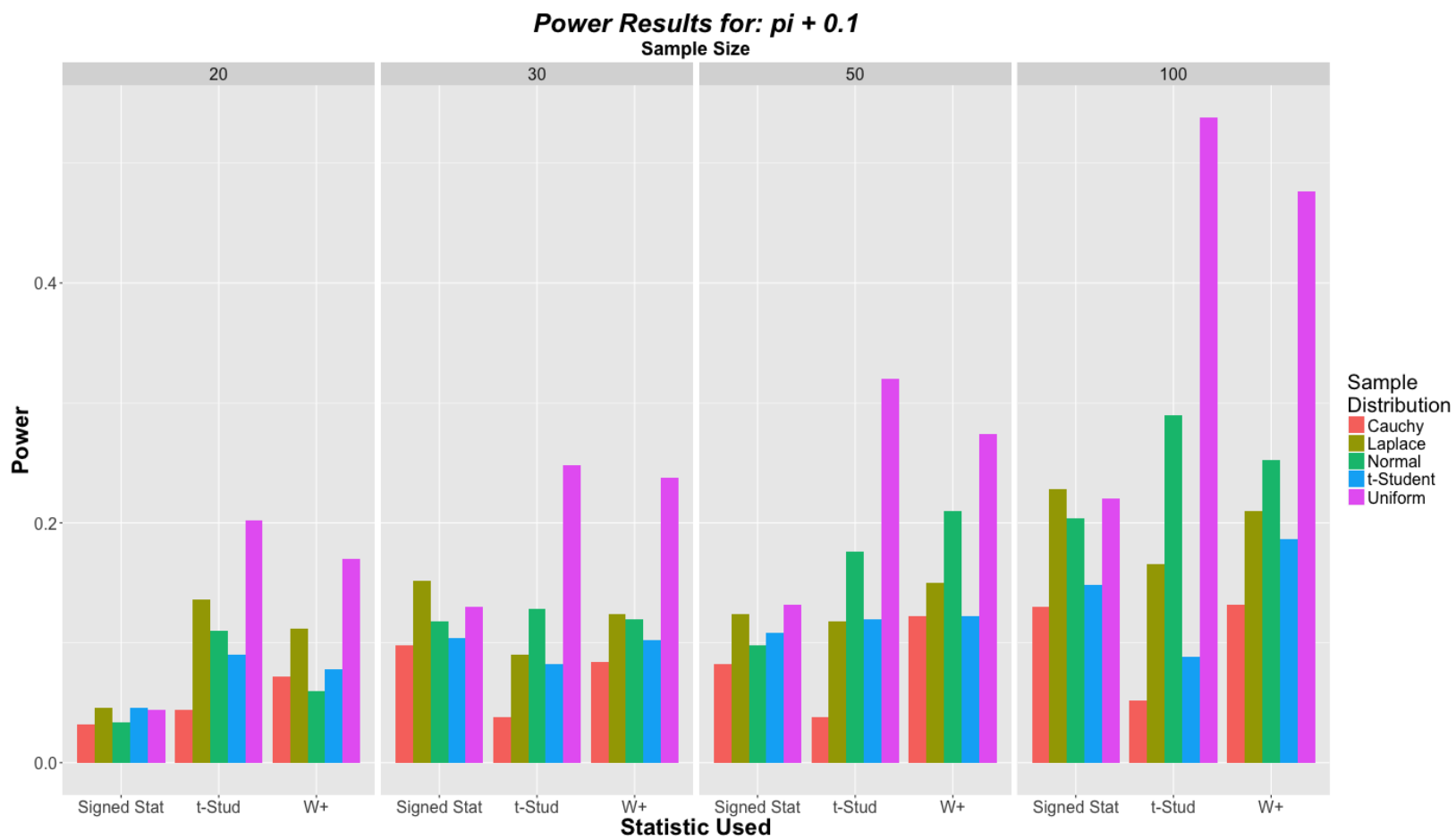


Figura 2.1: Comparación de los resultados para: $\pi + 0,1$

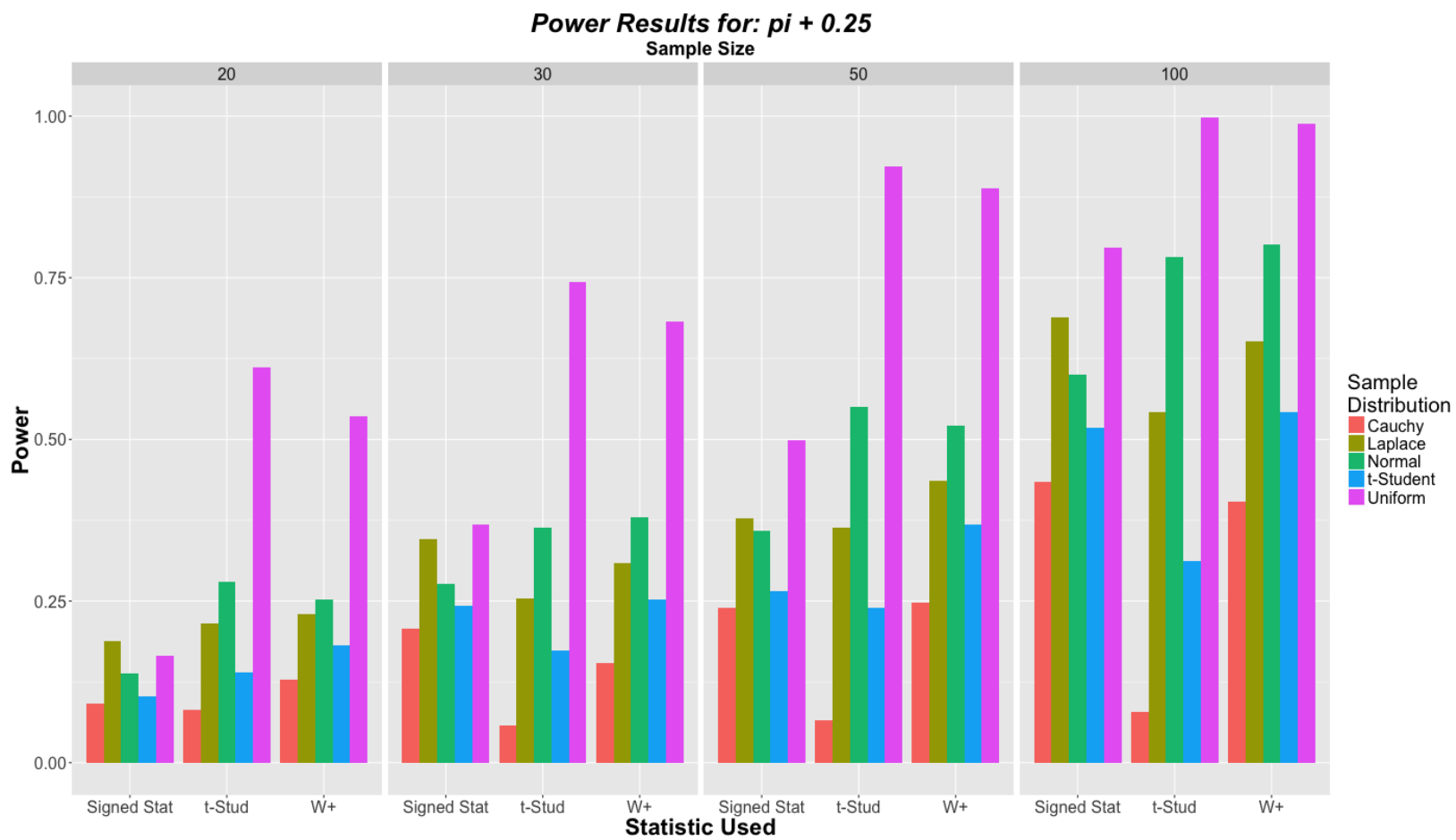


Figura 2.2: Comparación de los resultados para: $\pi + 0,25$

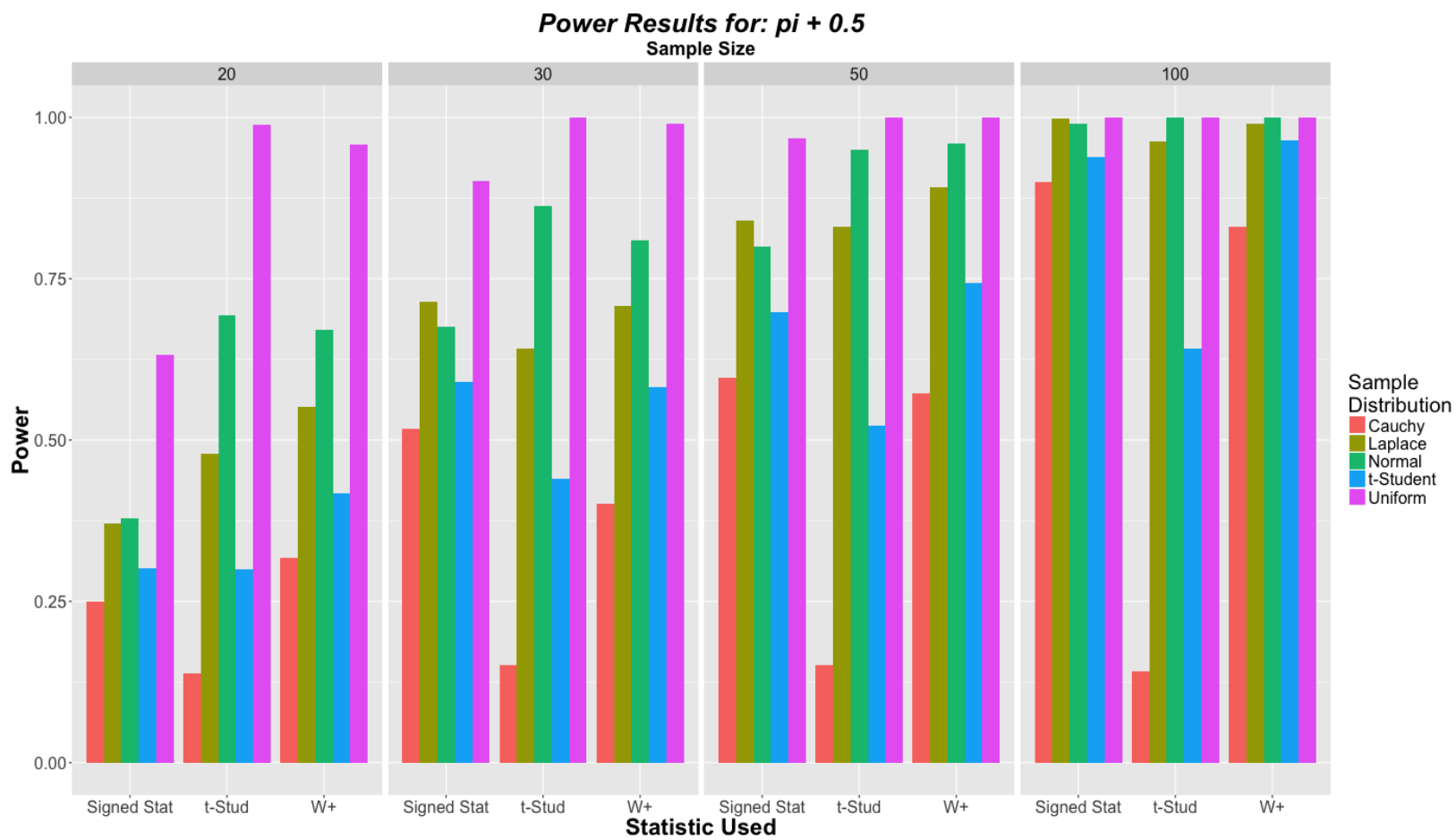


Figura 2.3: Comparación de los resultados para: $\pi + 0,5$

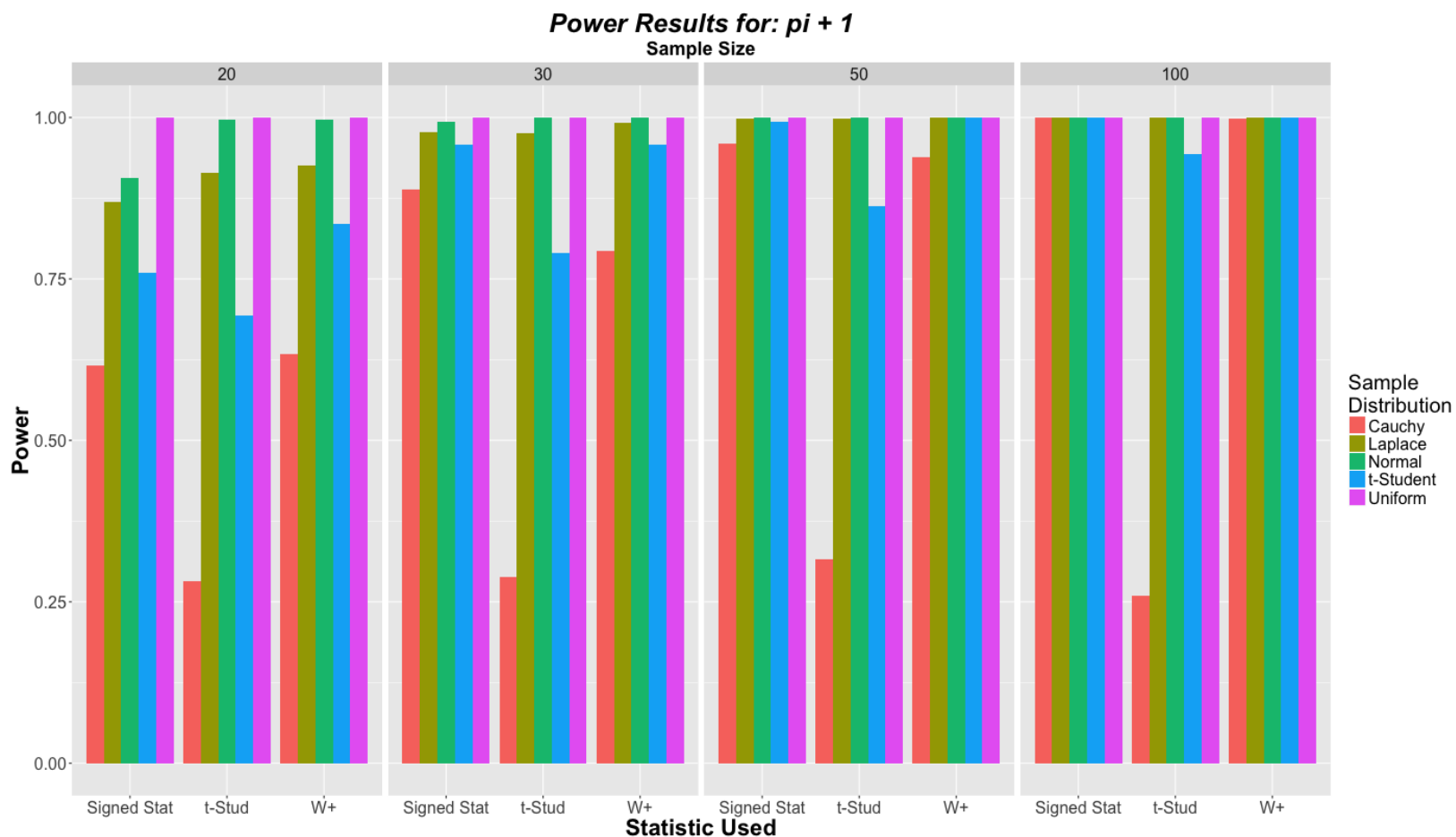


Figura 2.4: Comparación de los resultados para: $\pi + 1$

2.4. Conclusiones

Según los resultados presentados anteriormente se pudo ver que el estadístico t de Student es más efectivo para muestras con distribución uniforme y distribución normal. Esto sucede puesto que para tamaños grandes estas distribuciones tienden a cumplir con las condiciones necesarias para la efectividad en el uso del estadístico t de Student. También es importante ver que es un mal estadístico para la prueba cuando la muestra tiene distribución Cauchy; esto porque la distribución no tiende a tener las condiciones necesarias sobre normalidad y existencia de varianza necesarias.

En cuanto a los estadísticos no paramétricos se tuvo que el estadístico de wilcoxon signado tuvo mejor desempeño, pero es importante tener en cuenta que este

mismo fue más costoso computacionalmente. Su desempeño es mejor puesto que utiliza más información sobre las distribuciones de la muestra, más específicamente sus rangos. Para muestras y traslados de simetría grandes el desempeño del estadístico de signos y el wilcoxon signado es bastante similar. Por lo tanto, en tales casos, se preferiría usar el estadístico de signos por tiempo computacional. Para tamaños de muestra y traslados más pequeños es preferible utilizar el estadístico de wilcoxon signado dado que presenta mayor potencia.

Capítulo 3

Anexos

Por completitud, se incluye el código completo de las implementaciones. También se encuentra público en:

https://github.com/minigonche/no_parametrica_proj_1

```
1 #This script is intended as a solution for the first part of the first
  project
2 #of the course: non-parametric statistics, offered in the first semester of
  2017
3 #at Universidad de los Andes by Adolfo Quiroz
4
5 # Used libraries
6 # For string comparisson
7 library('pracma')
8 #for the laplace distribution
9 library('rmutil')
10 #library for ggplot
11 library('ggplot2')
12 #-----
13 #The following Script will simulate the convergence of a Normal Aproximation
14 #Usinig a random variable based on the wilcoxon statistic, and will compare
15 # the strength of different hipotesis test using diffent statistics
16 #-----
17
18 #-----
19 #----- All Around Funcions -----
20 #-----
21
22
23 #Since the Wilcoxonstatistic is based on the partition formula, here's a
  recursive
24 #Formula to calculate its exact value. (Implementation given by the Professor
  )
25
26 #Function that gives the partitions for a given natural number.
27 #The number of subsets that of {1,2,...,n} that sum the position of the
28 #returned vector
29 cnk= function(n){
30   N=1+n*(n+1)/2
31   vector0=vector1=rep(0,N)
32   vector0[1:2]=1
33   for(m in 2:n){
34     vector1=vector0
35     vector1[(m+1):N]=vector1[(m+1):N]+vector0[1:(N-m)]
```

```

36     vector0=vector1
37 } # fin del for m
38 return(vector1)
39 }
40 #End of cnk
41
42
43 #series sum
44 series_sum = function(n)
45 {
46     return(n*(n+1)/2)
47 }
48
49
50 #Global list with the probabilities of the W+ statistic (so they only need to
    be computed once)
51 global_prob_w_plus = list()
52
53 #Probability function of the Wilcoxon plus statistic
54 #Calculates a vector which the entry k corresponds to:
55 #P(W+ = k)
56 #Note that the length of the returned vector is: (n*(n+1)/2) + 1
57 #Since n*(n+1)/2 is the maximum value, it also includes zero
58 prob_w_plus = function(n)
59 {
60     #PARAMETERS
61     #-----
62     #n integer: correspondes to the number of samples taken into account
63     #             into the Signed Wilcoxon statistic
64
65     #checks if it has been assigned to the global enviornment
66     if(is.null(global_prob_w_plus[n][[1]]))
67     {
68         #Extracts the partition formula
69         probability = cnk(n)
70         #Divides under all possible subsets
71         probability = probability/(2**n)
72         #Assignes the corresponding probability
73         global_prob_w_plus[[n]] <- probability
74         print(paste('W+ Probability Calculated:',n))
75     }
76 }
77
78 return(global_prob_w_plus[[n]])
79
80 }
81 #end of prob_w_plus
82
83
84 #Calculates the cumulative probability function of the Wilcoxon Signed
    statistic
85 acum_prob_w_plus = function(n, prob_vec = NULL)
86 {
87     #PARAMETER
88     #-----
89     #n integer: correspondes to the number of samples taken into account
90     #             into the Signed Wilcoxon statistic
91     #prob_vec vector: the probability vector of the Wilcoxon Signed statistic
92     #             so that it does not need to be computed again
93
94     #check the local and global enviornment to see if the w+ probability vector
        is already computed
95     if(is.null(prob_vec))
96     {
97         prob_vec = prob_w_plus(n)
98     }
99
100     #calculates the distribution as a matrix multiplication

```

```

101 #recall that the distribution is the position sum of the probability vector
102 # that is the dot product with a lower triangular matrix
103 m = matrix(1, length(prob_vec), length(prob_vec))
104 m = lower.tri(m, diag = TRUE)
105 response = m %*% prob_vec
106 return(response)
107 }
108 }
109 #end_w_plus
110
111 #function that gives the probability that the wilcoxon signed rank statistic
112 # is greater or equal
113 # to the recieved parameter
114 prob_greater_w_plus = function(t, n)
115 {
116   prob = prob_w_plus(n)
117   initial_index = ceiling(t)
118   #the plus one if because R starts at 0
119   if(initial_index + 1 > length(prob) )
120   {
121     return(0)
122   }
123   return(sum(prob[initial_index:length(prob)]))
124 }
125
126
127 #Gives the sequence of quantiles that will be calculated
128 #unified method
129 give_quantiles = function()
130 {
131   max_index = 7
132   step = 0.2
133
134   #initial indexes of distribution
135   return(seq(-max_index,max_index,step))
136 }
137
138
139 #-----
140 #----- First Problem -----
141 #-----
142
143
144
145 #Calculates the proposed statistic acumulative probability function
146 #the vector will be from -6 to 6 steps of 0.25
147 proposed_statistic_acum_prob = function(n, prob_vec = NULL)
148 {
149   #PARAMETER
150   #-----
151   #n integer: correspondes to the number of samples taken into account
152   #           into the Signed Wilcoxon statistic
153   #prob_vec vector: the probability vector of the Wilcoxon Signed statistic
154   #               so that it does not need to be coputed again
155
156   #calculates the distribution as a matrix multiplication
157   #recall that the distribution is the position sum of the probability vector
158   if(is.null(prob_vec))
159   {
160     prob_vec = prob_w_plus(n)
161   }
162
163   #gets the distribution of W+
164   d_w = acum_prob_w_plus(n, prob_vec)
165
166   #calculates the corresponding coordinates
167   N = series_sum(n)

```

```

168 index = give_quantiles()
169
170
171 #Transforms the indexes to fit the new distribution
172 index = index*(choose(n,2)/(sqrt(3*n)))
173 index = index + N/2
174 index = floor(index)
175
176 #shifts the indexes since R starts at 1
177 index = index + 1
178
179 #Calculates the distribution for the acceptable indexes
180 # The out of range indexes are left as zero or one (accordingly)
181 distribution = rep(0,length(index))
182 selected_index = which(index > 0 & index <= length(d_w))
183 distribution[selected_index] = d_w[index[selected_index]]
184 distribution[index > length(d_w)] = 1
185
186 return(distribution)
187
188 }
189 #end of proposed_statistic_acum_prob
190
191 #function that plots the metric of conertion versus the size of the sample
192 plot_convergence = function(max_sample = 200, min_sample = 15, location = "~/
193   Dropbox/Universidad/Materias/Estadistica No parametrica/Proyecto1/plots"
194   )
195 {
196   differences = rep(0,max_sample-min_sample)
197   x_coor = min_sample:max_sample
198
199   #calculates the differences
200   for(n in x_coor)
201   {
202     #extracts the proposed statistic
203     proposed = proposed_statistic_acum_prob(n)
204     #extracts the normal distribution
205     normal = pnorm(give_quantiles())
206
207     differences[n - (min_sample-1)] = max(abs(proposed - normal))
208   }
209
210   png(paste(location,'/',n,'.png', sep = '' ) , width = 600, height = 600)
211   plot(x = x_coor,
212        y = differences,
213        main = 'Aproximation Convergence',
214        cex.main = 3,
215        xlab = 'Sample Size',
216        ylab = 'Difference between statistic and Normal(0,1)',
217        cex.lab = 2)
218
219   for(j in 1:length(x_coor)-1)
220   {
221     lines(x = c(x_coor[j],x_coor[j+1]), y = c(differences[j], differences[j
222       +1]), col = 'blue')
223   }
224
225   #adds the treshold line
226   lines(x = c(min_sample,max_sample), y = c(0.005,0.005), col = 'red')
227   #adds legend
228   legend('topright', c('0.005 Treshold', 'Supreme Difference'),
229         lty=c(1,1), lwd=c(2,2),col=c('red','blue'), cex = 2) # gives the
230     legend lines the correct color and width
231   dev.off()

```

```

232 }
233
234 #Function that plots the different prob-prob graphs
235 plot_prob_prob = function(n_values = c(15),
236                           location = "~/Dropbox/Universidad/Materias/
                           Estadistica No parametrica/Proyecto1/plots")
237 {
238   #PARAMETER
239   #-----
240   # n_values vector: a vector with the n values that will be plotted
241   # title_values vector: a string vector with the different titles for the
242   # plots
243   # it should have the same size as the n_values vector
244   #location String: the location where the plots will be saved
245   quantils = give_quantils()
246   for(i in 1:length(n_values))
247   {
248     #plots the proposed statistic
249     n = n_values[i]
250     title = paste('Proposed Statistic vs Normal \n P-P Plot for n =',n)
251     p_dist = proposed_statistic_acum_prob(n)
252     n_dist = pnorm(quantils)
253
254     #Plots the normal distribution
255     png(paste(location,'/',n,'.png', sep = '' ))
256
257     plot(x = n_dist, y = p_dist, col = 'blue', xlab = 'Normal(0,1)
258           Distribution', ylab = 'Proposed Statistic Distributon', cex = 1.5)
259
260     #Adds the lines
261     lines(x = c(-0.2,1.2), y = c(-0.2,1.2), col = 'red')
262
263     #Adds title and legend
264     title(main = title, cex.main = 1.5)
265     legend('topleft', c("Equality Line"),
266           lty=c(1,1), lwd=c(2,2),col=c('red'), cex = 1) # gives the legend
267           lines the correct color and width
268
269     dev.off()
270   }
271 }
272
273 #-----
274 #----- Second Problem -----
275 #-----
276
277 #generates the desired samples from the symetric distributions for the
278 #hipotesis testing
279 generate_sample = function(n = 20, distribution = 'normal', theta = pi)
280 {
281   #PARAMETER
282   #-----
283   # n integer: the size of teh sample
284   # distribution String: a string with the desired distribution.
285   # can be one of the following:
286   # - normal: Normal(theta,1)
287   # - unif: Unif(theta - 1,theta + 1)
288   # - laplace: Double laplace distribution
289   # - student: A shifted t student distribution with two freedom
290   # grades
291   # - cauchy: a cauchy distribution
292   # theta numeric: the center for the distributions
293
294   # Normal dsitribution
295   if(strcmpi(distribution,'normal'))

```

```

295 {
296   return(rnorm(n, mean = theta, sd = 1))
297 }
298 # Uniform distribution
299 if(strcmpi(distribution,'uniform'))
300 {
301   return(runif(n, min = theta-1, max = theta+1))
302 }
303 # laplace distribution
304 if(strcmpi(distribution,'laplace'))
305 {
306   return(rlaplace(n, m=theta, s=1))
307 }
308
309 # t-student
310 if(strcmpi(distribution,'student'))
311 {
312   return(rt(n, df = 2) + theta)
313 }
314 # cauchy
315 if(strcmpi(distribution,'cauchy'))
316 {
317   return(rcauchy(n, location = theta, scale = 1))
318 }
319
320 stop(paste('The received parameter is not recognized or supported:',
321           distribution))
322 }
323
324 #Calculates the desired statistic, given the sample of Z
325 get_statistic = function(Z, stat_name = 'student', theta_0 = pi)
326 {
327   #PARAMETER
328   #-----
329   # Z vector: the sample vector
330   # stat_name String: a string with the desired statistic
331   #           can be one of the following:
332   #           - student: The t-student statistic
333   #           - signed: The signed statistic
334   #           - rank: the signed rank statistic
335   # theta_0 numeric: the theta for the null hypothesis
336
337   n = length(Z)
338   # t-student
339   if(strcmpi(stat_name,'student'))
340   {
341     temp = sqrt(n)*(mean(Z) - theta_0)/std(Z)
342     return(temp)
343   }
344   #signed statistic
345   if(strcmpi(stat_name,'signed'))
346   {
347     temp = Z - theta_0
348     return(length(which(temp > 0)))
349   }
350
351   #signed rank statistic (wilcoxon)
352   if(strcmpi(stat_name,'rank'))
353   {
354     temp = Z - theta_0
355     #orders them
356     temp = temp[order(abs(temp))]
357     #sums the ranks
358     return(sum(which(temp > 0)))
359   }
360
361   stop(paste('The received parameter is not recognized or supported:', stat_

```

```

        name))
362 }
363 # end of get_statistic
364
365 #function that gets the p-value for the given statistic
366 get_p_value = function(omega,n, stat_name = 'student')
367 {
368   #PARAMETER
369   #-----
370   # omega numeric: the value of the statistic
371   # stat_name String: a string with the desired statistic
372   #           can be one of the following:
373   #           - student: The t-student statistic
374   #           - signed: The signed statistic
375   #           - rank: the signed rank statistic
376
377
378   # t-student
379   if(strcmpi(stat_name,'student'))
380   {
381     #dsitributes t-student
382     return(1 - pt(omega,n-1))
383   }
384   #signed statistic
385   if(strcmpi(stat_name,'signed'))
386   {
387     #distributes binomial with p = 0.5
388     return(1 - pbinom(ceiling(omega) - 1, n, 0.5))
389   }
390
391   #signed rank statistic (wilcoxon)
392   if(strcmpi(stat_name,'rank'))
393   {
394     return(prob_greater_w_plus(omega,n))
395   }
396
397   stop(paste('The received parameter is not recognized or supported:', stat_
398             name))
399 }
400 # end of get_statistic
401
402 #The main method for item 2
403 test_hypothesis = function(num_ite = 500,
404                             alpha = 0.05,
405                             thetas = c(pi + 0.1, pi + 0.25, pi + 0.5, pi + 1),
406                             sample_sizes = c(20,30,50,100))
407 {
408   #PARAMETER
409   #-----
410   # num_ite numeric: the number of iterations for each scenario
411   # alpha numeric: the desired confidence for the hypothesis test
412   # thetas numeric vector: a numeric vector with the desired centralities to
413   #           check. The method
414   #           assumes the null hypothesis with theta = pi
415   # sample_size numeric vecor: a numeric vector with the desired sample sizes
416   #           to test
417
418   sample_types = c('normal','uniform','laplace','student','cauchy')
419   statistics = c('student','signed','rank')
420
421   #creates all the possible combinations
422   result = expand.grid(sample_sizes,thetas, sample_types,statistics)
423   #assignes colnames
424   colnames(result) = c('sample_size','theta','sample_dist','stat')
425   # Adds the power columns
426   result$power = rep(0,nrow(result))
427   #Adds the shift that the theta received (for reading purposes)
428   result$shift = result$theta - pi

```



```

426
427
428
429 #iterates over each row and excecutes experiment
430 for(i in 1:nrow(result))
431 {
432   row = result[i,]
433   print(row)
434
435   discards = 0
436
437   for(j in 1:num_ite)
438   {
439     #Sample
440     Z = generate_sample(n = row$sample_size, distribution = toString(row$
         sample_dist), theta = row$theta)
441     #Statistic
442     omega = get_statistic(Z = Z, stat_name = toString(row$stat), theta_0 =
         pi)
443     #P-Value
444     p_value = get_p_value(omega = omega, n = row$sample_size, stat_name =
         toString(row$stat))
445     #Checks if the null hipotheiss is discarded
446     if(p_value <= alpha)
447     {
448       discards = discards + 1
449     }
450   }
451
452   #Assignes power
453   result[i,]$power = discards/num_ite
454 }
455
456
457 return(result)
458 }
459 #end of test_hipotesis
460
461 #plots the results for the experiment using ggplot
462
463 plot_results = function(experiment_results, location = "~/Dropbox/Universidad
         /Materias/Estadistica No parametrica/Proyecto1/plots" )
464 {
465   #PARAMETER
466   #-----
467   # experiment_results data.frame: a data frame with the results of the
         experiments. Must contain a specific
468   #         columns with the shifts
469   # location String: the location where the plots will be saved
470
471   for(shift in unique(experiment_results$shift))
472   {
473
474     results = experiment_results[which(experiment_results$shift == shift),]
475
476     title = paste("Power Results for: pi +", shift)
477
478     #edits the data-frame so it is in readable format
479     #sample distributions
480     results$sample_dist = gsub('normal','Normal', results$sample_dist)
481     results$sample_dist = gsub('uniform','Uniform', results$sample_dist)
482     results$sample_dist = gsub('laplace','Laplace', results$sample_dist)
483     results$sample_dist = gsub('cauchy','Cauchy', results$sample_dist)
484     results$sample_dist = gsub('student','t-Student', results$sample_dist)
485
486     #Statistic
487     results$stat = gsub('student','t-Stud', results$stat)
488     results$stat = gsub('signed','Signed Stat', results$stat)

```

```

489 results$stat = gsub('rank','W+', results$stat)
490
491 q = ggplot(results, aes(x = stat, y = power, fill = sample_dist) ) +
      geom_bar(position="dodge", stat="identity")
492 q = q + facet_grid(~sample_size)
493 q = q + scale_fill_discrete(name = "Sample\nDistribution")
494 q = q + labs(title = title, subtitle = "Sample Size")
495 q = q + xlab("Statistic Used") + ylab("Power")
496 q = q + theme(text = element_text(size=20),
497               plot.title = element_text( size=25, face="bold.italic",
498               hjust = 0.5),
499               plot.subtitle = element_text( size=18, face="bold", hjust =
500               0.5),
501               axis.title.x = element_text( size=22, face="bold"),
502               axis.title.y = element_text( size=22, face="bold"))
503
504 png(paste(location,'/power_',shift,'.png', sep = '' ), width = 1400,
505     height = 800)
506 q
507 dev.off()
508
509 }
510 #end of plot_results
511 res = test_hipotesis(num_ite = 500)

```

Bibliografía

Ronald H Randles and Douglas A Wolfe. *Introduction to the theory of nonparametric statistics*, volume 1. Wiley New York, 1979.