

1. Introducción

El siguiente documento fue compilado usando **knitr** un paquete diseñado para la generación *dinámica* de reportes para **R** y *Python*. Esto se hace con el propósito de tener un documento consistente y replicable. Normalmente las entregas tienen un pantallazo del código y de la salida, sin embargo, esto además de verse muy *feo* no tiene alguna manera sencilla de hacer la replicación de la simulación. Por eso mismo, en esta entrega se utiliza este paquete y facilitar manipulación, lectura del código, además de inmediatamente mostrar la salida de los códigos.

En el archivo comprimido `.7z` se encuentran varios archivos; `notas.tex`, `notas.Rnw`, `notas.pdf`, `referencias.bib`, `run.R`, `body_fat.csv`, `loans_income.csv`, `renv.lock` y unos directorios `img/` `renv/`. Aquí interesan dos archivos; `notas.Rnw` se encuentra todo el código fuente del documento (\LaTeX y **R**), en `run.R` se encuentra el *script* para poder compilar el documento. Si se desea hacer uso del *script* se necesitan unas librerías, paquetes y programas ya instalados en su máquina local, estos son; \LaTeX , se recomienda hacer la instalación local dependiendo de su sistema operativo; **R**, lenguaje de programación para estadística; **knitr**, para poder compilar el `.Rnw` y que genere un PDF. Con estas dependencias ya instaladas, tiene que ejecutar el archivo `run.R` (en sistemas UNIX como GNU/Linux, MacOSX, FreeBSD, etc. Se puede correr como cualquier *script* de *bash*, si cuenta con otro sistema operativo tiene que ejecutar el *script* dentro de un entorno para correr código en **R**). Esto va a generar un PDF y de ahí puede ejecutar cuantas veces desee el *script* y ver cómo cambian los valores del PDF dinámicamente. Si no desea hacer todo ese proceso de compilación del `.Rnw` puede ejecutar directamente el PDF que viene en el comprimido, siendo este la última salida de la compilación de los autores.

Por último, todo el presente trabajo fue hecho en un entorno virtual de **R** para evitar conflictos entre librerías, se trabajó con el paquete **renv** para manejar el ambiente local junto con sus librerías. Puede revisar la documentación [aquí](#) para tener más información de cómo manejar los ambientes virtuales y tener una mayor estabilidad en su sistema, así evita lo que le pasó a uno de los autores, que decidió instalar una librería en el ambiente global de *Python* y dañó todo el sistema base.

2. Stopping generating new simulation data

Write a program to generate standard normal random variables until you have generated n of them, where $n \geq 100$ is such that $S/\sqrt{n} < 0.01$, where S is the sample standard deviation of the n data values. Note that this is the “Method for Determining When to Stop Generating New Data”. Also, answer the following questions:

Debido a que se va a generar variables aleatorias **normales estándar** X_i , es decir, $X_i \sim \mathcal{N}(0, 1)$ y se debe parar de generar cuando $\frac{S}{\sqrt{n}} < 0.01$.

```
x <- numeric()
s <- 1
n <- 0
```

```
while(n < 100 || s / sqrt(n) >= 0.01) {  
  x <- c(x, rnorm(1))  
  n <- length(x)  
  s <- sd(x)  
}  
  
cat("Cantidad de Xi generados", n, "\n")  
  
## Cantidad de Xi generados 10198  
  
cat("Desviación estándar", round(s, 5), "\n")  
  
## Desviación estándar 1.00985  
  
cat("Error estándar S/sqrt(n)", round(s/sqrt(n), 5), "\n")  
  
## Error estándar S/sqrt(n) 0.01
```

En el código se declara una variable s inicializada en 1 para que luego sea nuevamente computada a la desviación estándar de x que es un vector que va a almacenar todas las variables normales aleatorias generadas. Según el enunciado tenemos una condición de que al menos deben haber 100 variables aleatorias generadas ($n \geq 100$) y que el error estándar sea menor a 0.01 ($S/\sqrt{n} < 0.01$). En el ciclo `while`, por lo tanto, tiene sentido que el ciclo continúe si alguna de las afirmaciones anteriores son falsas, por eso queda la condición de esa manera. Se hace uso de la función `rnorm(1)` para que genere una variable aleatoria normal con los valores por defecto (promedio 0 y desviación 1).

2.1. How many normals do you think will be generated? Give an analytic estimate.

2.1.1. Respuesta

Aunque ya se tengan los resultados de la simulación, se puede hacer un estimado analítico con la condición $S/\sqrt{n} < 0.01$ donde se puede despejar n para saber cuántas variables se necesitan para parar el criterio.

Se tiene en primer lugar la inecuación:

$$\frac{S}{\sqrt{n}} < 0.01$$

Se eleva ambas partes con menos 1:

$$\left(\frac{S}{\sqrt{n}}\right)^{-1} > 0.01^{-1}$$
$$\frac{\sqrt{n}}{S} > 100$$

Multiplicando ambas partes por S :

$$\begin{aligned} \mathcal{J} \times \frac{\sqrt{n}}{\mathcal{J}} &> 100 \times S \\ \sqrt{n} &> 100 \times S \end{aligned}$$

Ahora se cancela la raíz elevando ambas partes al cuadrado:

$$(\sqrt{n})^2 > (100 \times S)^2$$

Quedando entonces:

$$n > 10000 \times S^2$$

Para hacer el ejercicio se utilizó en R la función `rnorm(1, mean = 0, sd = 1)`, de esta manera genera únicamente 1 valor con media 0 y desviación 1, por lo tanto, podemos hacer $S = 1$ para estimar cuántos n necesitamos para que se cumpla condición y deje de generar variables aleatorias, por lo tanto:

$$\begin{aligned} n &> 10000 \times (1)^2 \\ n &> 10000 \end{aligned}$$

Este valor se acerca bastante al que se imprime en la simulación.

2.2. How many normals did you generate?

2.2.1. Respuesta

Se han generado `n= 10198` normales.

2.3. What is the sample mean of all the normals generated?

2.3.1. Respuesta

Utilizando el comando `mean()` al vector `x` se obtiene 0.0173852.

2.4. What is the sample variance?

2.4.1. Respuesta

La varianza muestral se calcula como:

$$S^2 = \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{n - 1}$$

Esto se puede sacar con la función `var()` de R, por lo tanto, aplicando eso al vector `x` se obtiene 1.0197956.

2.5. Comment on the results of (2.3) and (2.4). Were they surprising?

2.5.1. Respuesta

El resultado de 2.3 es la media de \bar{x} que es 0.0173852 y el de 2.4 es la varianza muestral de \bar{x} dando 1.0197956. Estos valores muy sorprendentes no fueron... Se espera que al simular muchas veces una distribución normal, esta nos dé los valores de la media y la desviación estándar al cuadrado. Pero sí es interesante que dado este algoritmo para detenerse después de cruzar un S/\sqrt{n} nos dé valores muy cercanos a una distribución normal con $\mu = 0$ y $\sigma = 1$.

3. Gaining confidence with confidence intervals

We know that the $\mathcal{U}(-1, 1)$ r.v. has mean 0. Use a sample of size 1000 to estimate the mean and give a 95 % confidence interval (CI). Does the CI contain 0? Repeat the above a large number of times (≥ 100). What percentage of time does the CI contain 0? Write your code so that it produces output similar to the following:

```
Number of trials: 10

Sample mean  lower bound  upper bound  contains mean?
-0.0733      -0.1888      0.0422      1
-0.0267      -0.1335      0.0801      1
-0.0063      -0.1143      0.1017      1
-0.0820      -0.1869      0.0230      1
-0.0354      -0.1478      0.0771      1
-0.0751      -0.1863      0.0362      1
-0.0742      -0.1923      0.0440      1
 0.0071      -0.1011      0.1153      1
 0.0772      -0.0322      0.1867      1
-0.0243      -0.1370      0.0885      1

100 percent of CI's contained the mean
```

3.1. Respuesta

Se sabe que para tener un intervalo de confianza de 95 %, el límite inferior y superior serán, respectivamente:

$$\left(\bar{X} - 1.96 \frac{S}{\sqrt{n}}, \bar{X} + 1.96 \frac{S}{\sqrt{n}} \right)$$

El siguiente código permite sacar el intervalo de confianza del 95 % y también el estimado de la media.

```
n <- 1000
x <- runif(n, -1, 1)
xmean <- mean(x)
S <- sd(x)
L <- xmean - 1.96 * S/sqrt(n)
U <- xmean + 1.96 * S/sqrt(n)
cat("El estimado es", xmean, "\n")

## El estimado es -0.01663151

cat("95% está entre (", L, ", ", U, ") \n", sep="")

## 95% está entre (-0.05192202, 0.01865899)
```

Aquí se estimó la media de la distribución uniforme, dando un resultado de -0.0166315 , muy cercano a 0. El intervalo de confianza está entre -0.051922 y 0.018659 , en este caso sí contiene a 0.¹ Ahora, se va a realizar el proceso unas 200 veces para conocer el porcentaje de intervalos de confianza que contienen 0.

```
trials <- 200
n <- 1000
true_mean <- 0

results <- data.frame(
  sample_mean = numeric(trials),
  lower_bound = numeric(trials),
  upper_bound = numeric(trials),
  contains_mean = integer(trials)
)

for (i in 1:trials){
  x <- runif(n, -1, 1)
  xmean <- mean(x)
  S <- sd(x)
  L <- xmean - 1.96 * S/sqrt(n)
  U <- xmean + 1.96 * S/sqrt(n)
  contains <- as.integer(L <= true_mean & true_mean <= U)

  results[i, ] <- c(xmean, L, U, contains)
}

cat("Número de intentos:", trials, "\n")
```

¹Aquí hay que hacer una aclaración, los autores utilizaron knitr, una herramienta para la generación de reportes dinámicos en R, entonces cada vez que se compile este archivo .Rnw va a salir una salida diferente y puede ser que ya el intervalo no contenga 0.

```
## Número de intentos: 200
```

```
print(results, digits= 4)
```

##	sample_mean	lower_bound	upper_bound	contains_mean
## 1	9.114e-03	-0.0259926	0.044220	1
## 2	1.883e-02	-0.0165825	0.054250	1
## 3	1.106e-02	-0.0251099	0.047224	1
## 4	-5.354e-03	-0.0408235	0.030116	1
## 5	1.880e-03	-0.0324504	0.036210	1
## 6	2.786e-02	-0.0078839	0.063600	1
## 7	5.987e-03	-0.0302786	0.042253	1
## 8	-6.234e-03	-0.0418871	0.029419	1
## 9	-1.246e-02	-0.0492623	0.024349	1
## 10	-6.149e-03	-0.0410561	0.028758	1
## 11	-5.132e-03	-0.0410393	0.030775	1
## 12	-1.109e-02	-0.0465200	0.024337	1
## 13	1.308e-02	-0.0228377	0.049004	1
## 14	-2.340e-03	-0.0383425	0.033662	1
## 15	7.955e-03	-0.0287899	0.044700	1
## 16	-1.146e-02	-0.0466822	0.023762	1
## 17	-2.122e-02	-0.0575987	0.015158	1
## 18	-2.200e-02	-0.0568909	0.012886	1
## 19	1.686e-02	-0.0183366	0.052048	1
## 20	2.216e-03	-0.0337857	0.038218	1
## 21	3.978e-03	-0.0325484	0.040504	1
## 22	1.064e-02	-0.0240120	0.045291	1
## 23	-1.655e-02	-0.0517918	0.018701	1
## 24	-2.363e-02	-0.0592356	0.011980	1
## 25	-6.399e-03	-0.0419704	0.029172	1
## 26	-2.379e-02	-0.0590977	0.011512	1
## 27	-8.181e-03	-0.0437806	0.027418	1
## 28	-1.450e-02	-0.0499634	0.020964	1
## 29	-1.214e-02	-0.0477870	0.023497	1
## 30	-1.620e-02	-0.0518954	0.019491	1
## 31	1.943e-02	-0.0161756	0.055034	1
## 32	8.595e-03	-0.0273695	0.044560	1
## 33	2.043e-02	-0.0153104	0.056164	1
## 34	-3.439e-02	-0.0706630	0.001877	1
## 35	1.966e-02	-0.0163470	0.055657	1
## 36	1.222e-02	-0.0234325	0.047874	1
## 37	-1.430e-02	-0.0499019	0.021304	1
## 38	-2.051e-02	-0.0561805	0.015166	1
## 39	5.031e-03	-0.0307468	0.040809	1
## 40	6.933e-04	-0.0351831	0.036570	1

## 41	1.019e-02	-0.0251782	0.045548	1
## 42	4.521e-02	0.0090214	0.081405	0
## 43	4.542e-03	-0.0306436	0.039728	1
## 44	1.801e-06	-0.0357592	0.035763	1
## 45	-1.129e-02	-0.0466401	0.024066	1
## 46	8.932e-03	-0.0262723	0.044137	1
## 47	-1.937e-02	-0.0545625	0.015820	1
## 48	-4.169e-03	-0.0402089	0.031871	1
## 49	-2.265e-02	-0.0582527	0.012943	1
## 50	1.216e-02	-0.0232116	0.047528	1
## 51	3.689e-02	0.0007917	0.072989	0
## 52	-2.001e-03	-0.0381907	0.034189	1
## 53	-1.378e-02	-0.0496764	0.022112	1
## 54	-4.266e-03	-0.0407150	0.032182	1
## 55	-9.889e-03	-0.0452566	0.025478	1
## 56	1.095e-02	-0.0261696	0.048062	1
## 57	1.921e-02	-0.0171574	0.055573	1
## 58	-1.178e-02	-0.0474970	0.023939	1
## 59	-8.726e-03	-0.0450663	0.027615	1
## 60	-2.824e-02	-0.0641998	0.007728	1
## 61	-4.769e-03	-0.0399173	0.030380	1
## 62	3.001e-02	-0.0052267	0.065240	1
## 63	8.890e-03	-0.0265681	0.044349	1
## 64	3.453e-02	-0.0014190	0.070471	1
## 65	2.600e-02	-0.0094840	0.061479	1
## 66	-3.690e-03	-0.0404914	0.033111	1
## 67	2.815e-02	-0.0084581	0.064755	1
## 68	2.314e-03	-0.0333165	0.037945	1
## 69	7.242e-03	-0.0297992	0.044284	1
## 70	-2.589e-02	-0.0620649	0.010289	1
## 71	1.599e-02	-0.0191483	0.051128	1
## 72	7.368e-03	-0.0277467	0.042482	1
## 73	-4.265e-03	-0.0393907	0.030860	1
## 74	2.556e-02	-0.0106175	0.061728	1
## 75	2.561e-02	-0.0108899	0.062110	1
## 76	-3.112e-04	-0.0360654	0.035443	1
## 77	2.778e-03	-0.0329106	0.038466	1
## 78	-4.311e-03	-0.0406582	0.032037	1
## 79	3.943e-03	-0.0315251	0.039412	1
## 80	-1.886e-02	-0.0553230	0.017595	1
## 81	-3.333e-04	-0.0352785	0.034612	1
## 82	6.353e-03	-0.0289269	0.041632	1
## 83	1.781e-02	-0.0182680	0.053897	1
## 84	2.454e-02	-0.0103537	0.059443	1

## 85	-3.148e-02	-0.0679586	0.005002	1
## 86	1.273e-02	-0.0239809	0.049445	1
## 87	-1.474e-03	-0.0373964	0.034448	1
## 88	-7.918e-03	-0.0434727	0.027636	1
## 89	-2.758e-02	-0.0621813	0.007021	1
## 90	1.792e-02	-0.0180313	0.053880	1
## 91	-3.395e-02	-0.0699297	0.002037	1
## 92	-1.127e-02	-0.0458387	0.023306	1
## 93	-1.240e-02	-0.0479299	0.023123	1
## 94	-4.132e-02	-0.0765216	-0.006127	0
## 95	2.708e-02	-0.0084340	0.062600	1
## 96	1.249e-02	-0.0229147	0.047891	1
## 97	-2.190e-02	-0.0576468	0.013841	1
## 98	-6.273e-03	-0.0413547	0.028809	1
## 99	6.786e-03	-0.0288930	0.042464	1
## 100	-1.658e-02	-0.0524504	0.019283	1
## 101	-2.047e-03	-0.0378436	0.033750	1
## 102	4.558e-03	-0.0302380	0.039354	1
## 103	-1.707e-02	-0.0532893	0.019150	1
## 104	6.857e-03	-0.0294823	0.043196	1
## 105	3.683e-02	0.0011470	0.072508	0
## 106	1.173e-02	-0.0253368	0.048799	1
## 107	4.933e-02	0.0131008	0.085552	0
## 108	3.401e-02	-0.0022675	0.070293	1
## 109	-1.937e-03	-0.0386107	0.034737	1
## 110	3.432e-02	-0.0007225	0.069357	1
## 111	-7.145e-03	-0.0430597	0.028769	1
## 112	-1.238e-02	-0.0482504	0.023490	1
## 113	2.687e-02	-0.0094595	0.063201	1
## 114	8.631e-03	-0.0268261	0.044087	1
## 115	1.375e-02	-0.0221151	0.049624	1
## 116	-9.904e-03	-0.0456088	0.025801	1
## 117	8.005e-03	-0.0280895	0.044100	1
## 118	7.882e-03	-0.0275387	0.043303	1
## 119	-2.355e-02	-0.0591300	0.012037	1
## 120	-3.051e-02	-0.0667530	0.005724	1
## 121	-1.161e-02	-0.0473919	0.024182	1
## 122	-1.112e-03	-0.0358799	0.033656	1
## 123	3.149e-02	-0.0044495	0.067437	1
## 124	1.106e-02	-0.0246188	0.046730	1
## 125	4.471e-03	-0.0317559	0.040698	1
## 126	-7.970e-03	-0.0432378	0.027299	1
## 127	2.025e-02	-0.0151102	0.055606	1
## 128	1.498e-02	-0.0215990	0.051559	1

## 129	3.055e-02	-0.0044786	0.065574	1
## 130	-9.231e-03	-0.0447234	0.026260	1
## 131	-3.236e-02	-0.0681695	0.003456	1
## 132	2.110e-02	-0.0147306	0.056939	1
## 133	-1.650e-03	-0.0377649	0.034465	1
## 134	-1.203e-03	-0.0367873	0.034382	1
## 135	3.991e-02	0.0036337	0.076193	0
## 136	1.605e-02	-0.0205843	0.052679	1
## 137	1.738e-02	-0.0184283	0.053190	1
## 138	4.354e-02	0.0076194	0.079455	0
## 139	5.087e-03	-0.0306671	0.040841	1
## 140	3.743e-03	-0.0305178	0.038004	1
## 141	-4.356e-02	-0.0797014	-0.007415	0
## 142	2.252e-03	-0.0343834	0.038887	1
## 143	2.094e-02	-0.0148651	0.056739	1
## 144	1.626e-02	-0.0192407	0.051764	1
## 145	1.384e-02	-0.0214362	0.049108	1
## 146	1.370e-02	-0.0225243	0.049933	1
## 147	-1.086e-02	-0.0472441	0.025515	1
## 148	4.450e-03	-0.0312760	0.040176	1
## 149	6.852e-03	-0.0286596	0.042364	1
## 150	-2.081e-02	-0.0566383	0.015019	1
## 151	-3.608e-03	-0.0392265	0.032011	1
## 152	-5.167e-03	-0.0419837	0.031649	1
## 153	5.650e-02	0.0215266	0.091465	0
## 154	-8.887e-03	-0.0452865	0.027513	1
## 155	5.030e-03	-0.0310133	0.041072	1
## 156	2.137e-02	-0.0152626	0.057994	1
## 157	9.127e-03	-0.0267913	0.045046	1
## 158	-1.020e-02	-0.0456897	0.025298	1
## 159	-2.372e-03	-0.0380900	0.033347	1
## 160	2.810e-03	-0.0327780	0.038399	1
## 161	-1.833e-02	-0.0537669	0.017114	1
## 162	-2.113e-02	-0.0571027	0.014852	1
## 163	-2.870e-03	-0.0387719	0.033033	1
## 164	-8.018e-03	-0.0442808	0.028246	1
## 165	-2.240e-02	-0.0592792	0.014478	1
## 166	-2.698e-02	-0.0624579	0.008493	1
## 167	2.005e-02	-0.0159614	0.056053	1
## 168	-1.821e-02	-0.0546434	0.018215	1
## 169	3.898e-03	-0.0328464	0.040642	1
## 170	-2.703e-02	-0.0625422	0.008486	1
## 171	2.961e-02	-0.0055887	0.064807	1
## 172	-3.368e-04	-0.0365006	0.035827	1

```
## 173  1.089e-02 -0.0248871  0.046663  1
## 174  1.562e-02 -0.0192045  0.050435  1
## 175 -1.637e-02 -0.0524892  0.019757  1
## 176 -8.749e-03 -0.0442979  0.026800  1
## 177 -2.440e-02 -0.0607984  0.011996  1
## 178  1.583e-02 -0.0201671  0.051828  1
## 179 -4.655e-02 -0.0819730 -0.011118  0
## 180 -2.044e-02 -0.0570372  0.016164  1
## 181  3.853e-03 -0.0310463  0.038752  1
## 182 -2.604e-02 -0.0608398  0.008759  1
## 183  2.156e-02 -0.0143727  0.057485  1
## 184 -1.560e-03 -0.0366291  0.033510  1
## 185  8.378e-03 -0.0268459  0.043603  1
## 186 -1.737e-03 -0.0372613  0.033787  1
## 187  5.647e-04 -0.0357757  0.036905  1
## 188  2.505e-02 -0.0105823  0.060689  1
## 189  2.617e-03 -0.0339877  0.039222  1
## 190 -9.039e-03 -0.0447826  0.026704  1
## 191 -1.868e-03 -0.0382030  0.034467  1
## 192 -2.100e-02 -0.0579210  0.015919  1
## 193 -1.943e-02 -0.0556271  0.016762  1
## 194  3.817e-03 -0.0322500  0.039884  1
## 195  2.727e-02 -0.0092706  0.063820  1
## 196 -1.014e-02 -0.0455805  0.025307  1
## 197 -1.143e-03 -0.0361035  0.033817  1
## 198  1.072e-02 -0.0252266  0.046661  1
## 199 -1.388e-02 -0.0496213  0.021868  1
## 200  4.852e-03 -0.0313237  0.041028  1

porcentaje <- mean(results$contains_mean) * 100

cat("\n", porcentaje, "% de los intervalos de",
    "confianza contienen la media real\n")

##
## 95 % de los intervalos de confianza contienen la media real
```

En conclusión, el 95 % de los intervalos generados incluyen al 0.

4. Standard deviation of a proportion

Assume a manager is using the sample proportion \hat{p} to estimate the proportion p of a new shipment of computer chips that are defective. He doesn't know p for this shipment, but in previous shipments it has been close to 0.01, that is 1 % of chips have been defective.

4.1. If the manager wants the standard deviation of \hat{p} to be about 0.02, how large a sample should she take based on the assumption that the rate of defectives has not changed dramatically?

En clase se vieron los estimadores de probabilidad, donde se quiere estimar

$$p = (PX \in A)$$

Donde A es el subconjunto del espacio de estados de Ω de X . Es decir, para nuestro problema, este subconjunto de espacios en la muestra que tomó la administradora del cargamento de chips defectuosos. Se puede definir la variable indicadora Z como:

$$Z = \begin{cases} 1, & X \in A \\ 0, & X \notin A \end{cases}$$

1 significa que sí está defectuoso y 0 que no lo está. Se puede escribir el estimador como:

$$p = E[Z]$$

En clase se realizó una demostración de como el valor esperado de Z se le puede asignar a p . Se define la varianza de Z como:

$$\begin{aligned} \text{Var}(Z) &= E[Z^2] - E[Z]^2 \\ &= (1^2 \times P(X \in A) + 0^2 \times P(X \notin A)) - (P(X \in A))^2 \\ &= p - p^2 \end{aligned}$$

Aplicando factorización:

$$\text{Var}(Z) = p(1 - p)$$

Estimar p se puede realizar mediante el promedio muestral de Z , es decir, se puede estimar la proporción \hat{p} con el promedio de los chips que son defectuosos:

$$\hat{p} = \frac{1}{n} \sum_{i=1}^n Z_i$$

Se puede calcular entonces, la varianza del estimador:

$$\text{Var}(\hat{p}) = \frac{1}{n} \text{Var} \left(\sum_{i=1}^n Z_i \right)$$

Por independencia de los valores del cargamento de chips Z_i se puede meter la varianza de estos en la suma:

$$\begin{aligned}\text{Var}(\hat{p}) &= \frac{1}{n^2} \sum_{i=1}^n \text{Var}(Z_i) \\ &= \frac{1}{n^2} \sum_{i=1}^n p(1-p) = \frac{1}{n} p(1-p)\end{aligned}$$

Por lo tanto, una aproximación es:

$$\text{Var}(\hat{p}) \approx \frac{1}{n} \hat{p}(1-\hat{p})$$

Teniendo en cuenta que la desviación estándar es el cuadrado de la varianza, se puede despejar:

$$\sigma_{\hat{p}} \approx \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

El problema nos dice que la desviación estándar de \hat{p} ($\sigma_{\hat{p}}$) debe ser sobre 0.02 y nos piden hallar cuántas muestras debe tomar la administradora, por lo tanto, en la parte izquierda de la ecuación reemplazamos por 0.02, se despeja para n y la proporción de defectuosos se toma como 0.01 debido a que en el problema se nos dice que no ha cambiado tanto.

$$\begin{aligned}\sigma_{\hat{p}} &\approx \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \\ 0.02 &\approx \sqrt{\frac{0.1(1-0.1)}{n}} \\ 0.02 &\approx \sqrt{\frac{0.0099}{n}} \\ (0.02)^2 &\approx \left(\sqrt{\frac{0.0099}{n}} \right)^2 \\ 0.0004 &\approx \frac{0.0099}{n} \\ n &\approx \frac{0.0099}{0.0004} = 24.75\end{aligned}$$

Por lo tanto, la administradora debe tomar una muestra de por lo menos 25 chips para tener una desviación estándar de al menos 0.02.

4.2. Now suppose something went wrong with the production run and the actual proportion of defectives in the shipment is 0.3, that is 30 % are defective. Now what would be the actual standard deviation of \hat{p} for the sample size you choose in a)?

Teniendo 25 chips y cambiando el valor de la proporción de defectuosos, se obtiene:

$$\begin{aligned}
\sigma_{\hat{p}} &\approx \sqrt{\frac{\hat{p}(1 - \hat{p})}{n}} \\
&\approx \sqrt{\frac{0.3(1 - 0.3)}{25}} \\
&\approx 0.091
\end{aligned}$$

Es decir, si la proporción de detección es del 30 %, la desviación estándar de \hat{p} con una cantidad de chips de 25, será aproximadamente 0.091.

5. Bets

You pay 10.000 pesos to participate in a bet game, which consists in tossing two coins together. If two heads fall, you earn 15.000 pesos. If one head and one tail fall you earn 5.000 pesos. In any other case you earn nothing. Let X the random variable of your profit.

5.1. Analytically find the probability mass function p_X , the mean $E[X]$, and variance $\text{Var}(X)$ of X .

Dado que X es la variable aleatoria de la ganancia, podemos listar los posibles valores que puede tomar de acuerdo a los lanzamientos:

Salida moneda	Ganancia	Ganancia X	Probabilidad
HH	15.000	+5.000	0.25
HT-TH	5.000	-5.000	0.5
TT	0	-10.000	0.25

Cuadro 1: Posibles salidas y ganancias al tirar dos monedas, (H) cara y (T) sello.

Del cuadro 1 se toma que los posibles valores para las ganancias son:

$$X \in \{-5000, -10000, 5000\}$$

Teniendo en cuenta los valores que puede tomar X , la función probabilidad de masa debe retornar la probabilidad (valga la redundancia) de obtener alguno de los valores de X , en el cuadro 1 se pueden ver dichas probabilidades, estas fueron calculadas de acuerdo a contar cuántas veces pueden salir los diferentes valores para la ganancia de X .

$$p_X(x) = \begin{cases} 0.25, & x = 5000 \\ 0.5, & x = -5000 \\ 0.25, & x = -10000 \\ 0, & \text{En otro caso} \end{cases}$$

Debido a que estamos con un conjunto continuo para X , el valor esperado se multiplica la probabilidad por cada valor de este:

$$E[X] = \sum_{I=1}^3 x \dot{p}_X(x) = 5000(0.25) + (-5000)(0.5) + (-10000)(0.25)$$

$$E[X] = 1250 - 2500 - 2500 = -3750$$

Es decir, en promedio se pierde 3.750 pesos por jugada.

Para calcular la varianza se puede utilizar la fórmula:

$$\text{Var}(X) = E[X^2] - (E[X])^2$$

Se computa $E[X^2]$:

$$\begin{aligned} E[X^2] &= 5000^2(0.25) + (-5000)^2(0.5) + (-10000)^2(0.25) \\ &= 25000000(0.25) + 25000000(0.5) + 100000000(0.25) \\ &= 6250000 + 12500000 + 25000000 = 43750000 \end{aligned}$$

Ahora se le resta el cuadrado de la media:

$$\text{Var}(X) = 43750000 - (-3750)^2 = 43750000 - 14062500 = 29687500$$

5.2. Write a code that simulates the r.v. X using the command `sample`. Generate an *iid* `sample` $\{X_i\}$ of size $n = 10^5$.

```
outcomes <- c(-10000, -5000, 5000)
probabilities <- c(0.25, 0.5, 0.25)

n <- 10^5

X <- sample(outcomes, size = n, replace = TRUE, prob = probabilities)

length(X)

## [1] 100000

summary(X)

##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## -10000   -5000   -5000   -3769   -5000    5000
```

En este fragmento de código se realiza la creación del vector X de tamaño 10^5 creado con la función `sample` que toma las salidas que se explicaron anteriormente, el tamaño de 10^5 , que se reemplace y la probabilidad.

- 5.3. Modify your code to calculate: (i) the estimated mean profit \bar{X}_j for each sample subsequence: $\{X_1, X_2, \dots, X_j\}$, $j = 2, 3, \dots, n$, (ii) the 95 % CI's of each estimated mean profit \bar{X}_j .**

```
mean_estimates <- numeric(n)
lower_CI <- numeric(n)
upper_CI <- numeric(n)
ie <- 1.96

for (j in 2:n) {
  x_sub <- X[1:j]
  mean_j <- mean(x_sub)
  sd_j <- sd(x_sub)
  se_j <- sd_j / sqrt(j)

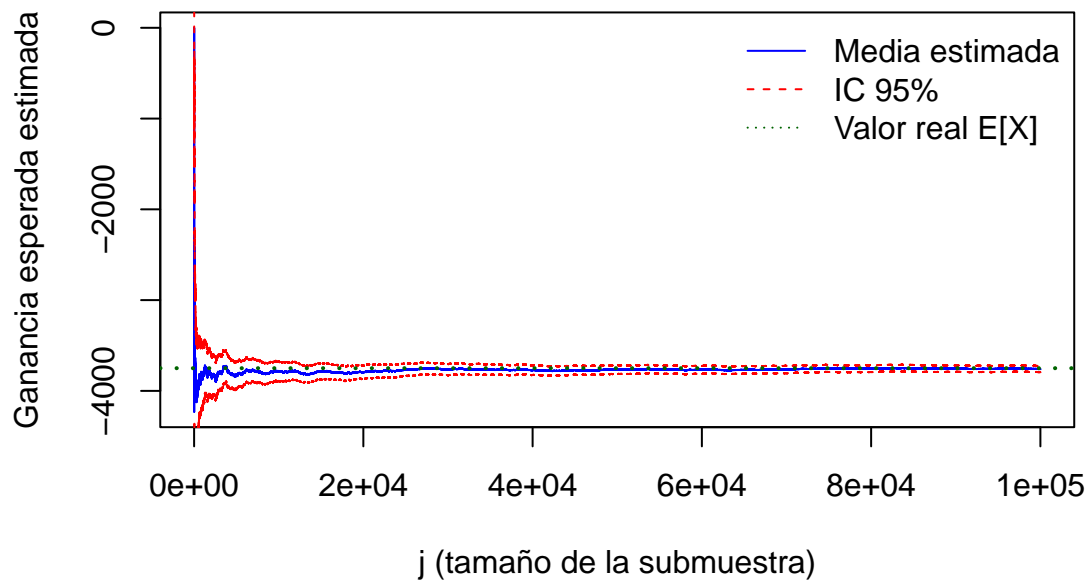
  mean_estimates[j] <- mean_j
  lower_CI[j] <- mean_j - ie * se_j
  upper_CI[j] <- mean_j + ie * se_j
}
```

En esta parte del código se crean tres vectores `mean_estimates`, `lower_CI`, `upper_CI` de tamaño 10^5 y también una variable `ie`. El bucle en cada iteración calcula la media muestral de los valores \bar{X}_j , la desviación estándar `sd_j` y el error estándar, estos son necesarios para luego calcular el intervalo de confianza junto con la media.

- 5.4. Plot \bar{X}_j and their 95 % CI's in terms of $j = 2, \dots, 10^5$. Add an horizontal line corresponding to the actual value $E[X]$.**

```
# Your plotting code remains the same
plot(2:n, mean_estimates[2:n], type = "l", col = "blue", lwd = 1,
     xlab = "j (tamaño de la submuestra)",
     ylab = "Ganancia esperada estimada",
     main = "Evolución de la media muestral y su IC del 95%")
lines(2:n, lower_CI[2:n], col = "red", lty = 2)
lines(2:n, upper_CI[2:n], col = "red", lty = 2)
abline(h = -3750, col = "darkgreen", lty = 3, lwd = 2)
legend("topright", legend = c("Media estimada", "IC 95%", "Valor real E[X]"),
     col = c("blue", "red", "darkgreen"), lty = c(1, 2, 3), bty = "n")
```

Evolución de la media muestral y su IC del 95%



Se puede ver que a medida que se aumenta la submuestra j , la media estimada \bar{X} se acerca más al valor esperado $E[X]$ y también, el intervalo de confianza se va acortando más.

5.5. Repeat c) and d) to estimate the probabilities $p_X(x)$, their 95% CI's, and their plots for each $j = 2, \dots, 10^5$, adding the actual values.

```
# Posibles valores de X
outcomes <- c(-10000, -5000, 5000)
probs_true <- c(0.25, 0.5, 0.25)

# Inicializar matrices para almacenar las estimaciones y CIs
p_estimates <- matrix(0, nrow = n, ncol = 3)
p_lower <- matrix(0, nrow = n, ncol = 3)
p_upper <- matrix(0, nrow = n, ncol = 3)

z <- 1.96 # Nivel de confianza del 95%

for (j in 2:n) {
  x_sub <- X[1:j]
  for (i in 1:3) {
    val <- outcomes[i]
    p_hat <- mean(x_sub == val)
```



```
se <- sqrt(p_hat * (1 - p_hat) / j)

p_estimates[j, i] <- p_hat
p_lower[j, i] <- max(0, p_hat - z * se)
p_upper[j, i] <- min(1, p_hat + z * se)
}
}
```

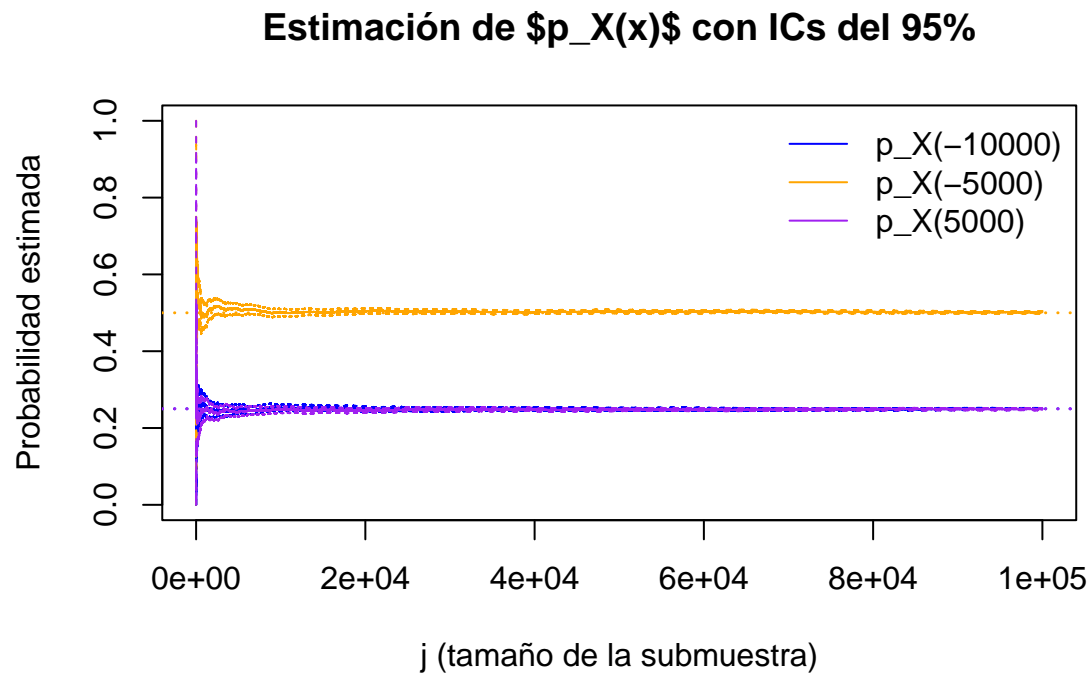
Gráfico de las probabilidades estimadas con sus ICs para cada posible valor de X :

```
colors <- c("blue", "orange", "purple")
labels <- c("p_X(-10000)", "p_X(-5000)", "p_X(5000)")
true_values <- c(0.25, 0.5, 0.25)

plot(2:n, p_estimates[2:n,1], type = "l", col = colors[1], ylim = c(0,1),
     xlab = "j (tamaño de la submuestra)", ylab = "Probabilidad estimada",
     main = "Estimación de $p_X(x)$ con ICs del 95%")

for (i in 1:3) {
  lines(2:n, p_estimates[2:n,i], col = colors[i], lwd = 1)
  lines(2:n, p_lower[2:n,i], col = colors[i], lty = 2)
  lines(2:n, p_upper[2:n,i], col = colors[i], lty = 2)
  abline(h = true_values[i], col = colors[i], lty = 3, lwd = 1.5)
}

legend("topright", legend = labels, col = colors, lty = 1:1, bty = "n")
```



Se observa cómo las probabilidades estimadas convergen a los valores reales a medida que aumenta el tamaño de muestra j , y cómo los intervalos de confianza se vuelven más estrechos.

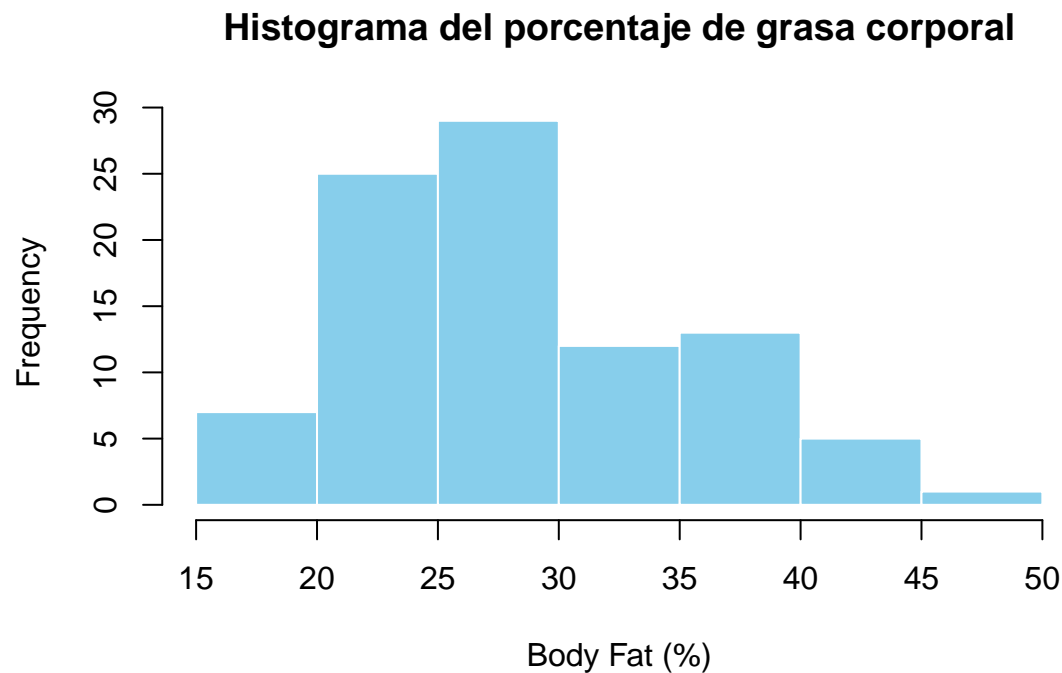
6. Bootstrap

This exercise is based on the article *Introduction to Bootstrapping in Statistics with an Example* and the dataset `body fat.csv` that contains the body fat percentages of 92 adolescent girls. Generate a program that gives:

6.1. An histogram of the sample data.

```
# Leer el archivo
fat_data <- read.csv("body_fat.csv" )

# Histograma
hist(fat_data[[colnames(fat_data)[1]]],
     main = "Histograma del porcentaje de grasa corporal",
     xlab = "Body Fat (%)",
     col = "skyblue",
     border = "white")
```



Realizar el histograma en R es bastante sencillo, únicamente tener el archivo `.csv` en la misma ruta local que donde se está trabajando, utilizar la función `read.csv()` para leer este archivo y luego con la función `hist()` se colocan los datos de la primera columna y otras opciones para el estilo del gráfico.

Se muestra que la mayor frecuencia de grasa corporal se encuentra entre 20-30 %.

6.2. The 95 % confidence interval of the mean of the data from the traditional method (i.e., via the Central Limit Theorem).

```
# Estadísticas necesarias
media <- mean(fat_data[[colnames(fat_data)[1]]])
std <- sd(fat_data[[colnames(fat_data)[1]]])
n <- length(fat_data[[colnames(fat_data)[1]]])

# Cálculo del error estándar
se <- std / sqrt(n)

# Intervalo de confianza del 95%
c <- 1.96
li <- media - c*se
ls <- media + c*se

# Mostrar resultado
```

```
cat(sprintf("The 95%% confidence interval for the mean is: [%.2f, %.2f]",
           li, ls))

## The 95% confidence interval for the mean is: [27.14, 29.99]
```

Como se puede apreciar en el programa, hay un 95% de confianza para que la media esté en el intervalo [27.1370094, 29.9934253]. Las funciones de R ayudan mucho a hacer los cálculos del promedio, la desviación estándar y el valor de n se sacaron con la ayuda de estas funciones. Los intervalos de confianza se sacaron como ya se viene haciendo, utilizando un $c = 1.96$.

6.3. A number of 500 bootstrapped samples from the original dataset, with 92 observations each.

```
original_data <- fat_data[[colnames(fat_data)[1]]]
n <- length(original_data)

# Generar 500 muestras bootstrap
bootstrap_samples <- replicate(500, sample(original_data, size = n,
                                           replace = TRUE), simplify = FALSE)
```

Este código utiliza la función `replicate()` para repetir la operación 500 veces, la función `sample()` toma una muestra con reemplazo del tamaño de 92. Esto termina generando que `bootstrap_samples` tenga 500 vectores, cada uno con 92 observaciones.

6.4. An histogram of the means of each bootstrapped sample.

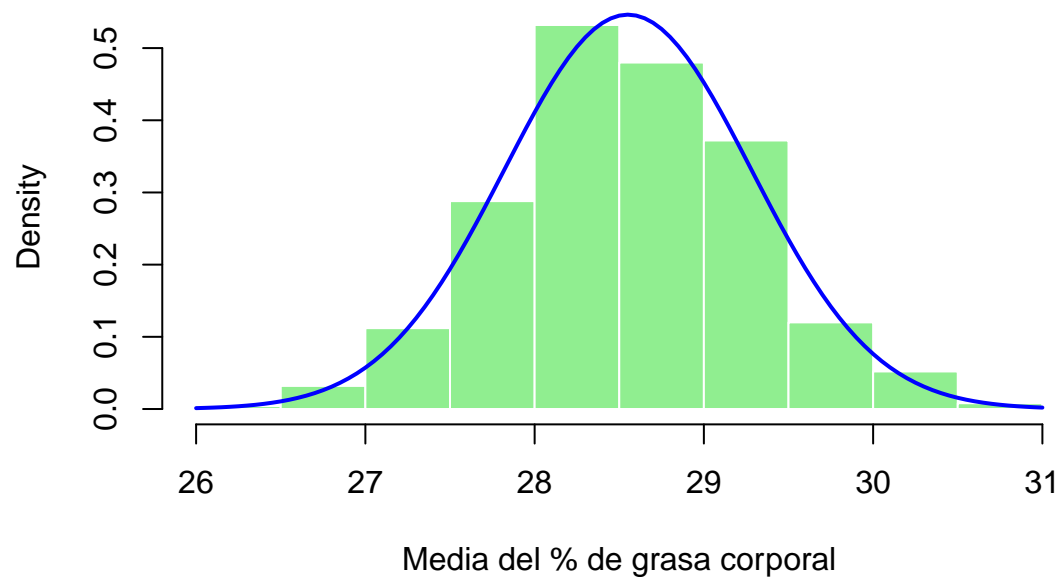
```
# Calcular la media de cada muestra bootstrap
bootstrap_means <- sapply(bootstrap_samples, mean)

# Histograma de las medias
hist(bootstrap_means,
     main = "Histograma de medias de muestras bootstrap",
     xlab = "Media del % de grasa corporal",
     col = "lightgreen",
     border = "white",
     probability = TRUE) # Importante: escala el eje y para densidades

# Parámetros de la normal
mean_boot <- mean(bootstrap_means)
sd_boot <- sd(bootstrap_means)
```

```
# Añadir curva de distribución normal
curve(dnorm(x, mean = mean_boot, sd = sd_boot),
      col = "blue", lwd = 2, add = TRUE)
```

Histograma de medias de muestras bootstrap



Con la función `apply()` le aplica a los elementos en `bootstrap_samples` la función `mean()`. Después se crea el histograma y se dibuja una curva de tipo campana para una distribución normal, reflejando la similitud de los promedios de los *bootstraps*.

6.5. A 95 % bootstrapped confidence interval of the mean.

```
# Calcular intervalo de confianza bootstrap del 95%
bootstrap_ci <- quantile(bootstrap_means, probs = c(0.025, 0.975))

bootstrap_ci

##      2.5%      97.5%
## 27.14875 30.03261
```

Se utiliza la función `quantile()` para obtener los intervalos exactos, para posteriormente poder hacer la comparación de ambos resultados.

Entonces, el 95 % del intervalo de confianza para la media es:

[27.15, 30.03]

6.6. A comparison of both confidence intervals

Con el método tradicional del teorema central del límite:

$$[27.14, 29.99]$$

Con el método *bootstrap*

$$[27.15, 30.03]$$

7. Reliability of a system

Suppose a 3-out-of-4 system where each component is functioning with probabilities $p = \{p_i\}$:

7.1. Write the structure function of the system $\phi(x)$.

La estructura funciona si al menos 3 de los 4 componentes están estructurados. Dado que hay dos estados; funciona ($\phi = 1$) o lo contrario ($\phi = 0$). Se puede tener un estado de los componentes $x = (x_1, x_2, x_3, x_4)$, donde cada $x_i \in \{0, 1\}$ representa si el componente i está funcionando.

La estructura entonces es:

$$\phi(x_1, x_2, x_3, x_4) = \begin{cases} 1 & \text{si } x_1 + x_2 + x_3 + x_4 \geq 3 \\ 0 & \text{si } x_1 + x_2 + x_3 + x_4 < 3 \end{cases}$$

7.2. Deduce analytically the reliability function $R(\mathbf{p})$. Evaluate it when $\{p_i\} = \{0.9, 0.5, 0.2, 0.1\}$

Se tomó gran parte de la información del libro *Introduction to Probability Models* de Sheldon Ross [2].

Suponiendo que X_i , el estado del i -ésimo componente, es una variable aleatoria tal que:

$$P\{X_i = 1\} = p_i = 1 - P\{X_i = 0\}$$

El valor p_i es la probabilidad de que el i -ésimo componente está funcionando, esto se llama la confiabilidad del i -ésimo componente. Si definimos r como:

$$r = P\{\phi(X) = 1\}, \quad \text{donde } X = (X_1, \dots, X_n)$$

r es llamada la confiabilidad del sistema. Cuando los componentes, es decir, las variables aleatorias $X_i, i = 1, \dots, n$, son independientes, se expresa r como una función de componentes de confiabilidad. Es decir:

$$r = r(\mathbf{p}), \quad \text{donde } \mathbf{p} = (p_1, \dots, p_n)$$

Esta función $r(\mathbf{p})$ es llamada la función de confiabilidad, que es la que se busca encontrar para este sistema 3-out-of-4.

Dado que para este sistema se necesita que por lo menos 3 de los 4 componentes sirvan, uno puede pasar a la función ϕ unos parámetros de $(1, 1, 1, 0)$, simbolizando que los componentes x_1, x_2 y x_3 están funcionando. La función de confiabilidad está dada entonces por todas las posibles combinaciones de los elementos que sumadas den mayor o igual a 3, es decir:

$$r(\mathbf{p}) = P \left\{ \sum_{i=1}^n X_i \geq 3 \right\}$$

Esto se puede expandir como:

$$r(\mathbf{p}) = P\{X = (1, 1, 1, 1)\} + P\{X = (1, 1, 1, 0)\} + P\{X = (1, 1, 0, 1)\} + P\{X = (1, 0, 1, 1)\} \\ + P\{X = (0, 1, 1, 1)\}$$

Teniendo que p_i simboliza la probabilidad de que el i -ésimo componente está funcionando y $1 - p_i$ es que no está funcionando el componente, lo anterior se puede escribir como:

$$r(\mathbf{p}) = p_1 p_2 p_3 p_4 + p_1 p_2 p_3 (1 - p_4) + p_1 p_2 (1 - p_3) p_4 + p_1 (1 - p_2) p_3 p_4 + (1 - p_1) p_2 p_3 p_4$$

Los casos para los cuales el sistema cuenta con un componente no funcionando, es decir, $1 - p_i$, se pueden agrupar en uno mismo, quedando:

$$r(\mathbf{p}) = p_1 p_2 p_3 p_4 + p_1 p_2 p_3 + p_1 p_2 p_4 \\ + p_1 p_3 p_4 + p_2 p_3 p_4 - 4 p_1 p_2 p_3 p_4$$

Sustrayendo con el término positivo $p_1 p_2 p_3 p_4$:

$$r(\mathbf{p}) = p_1 p_2 p_3 + p_1 p_2 p_4 + p_1 p_3 p_4 + p_2 p_3 p_4 - 3 p_1 p_2 p_3 p_4$$

Ahora, podemos evaluar cuando $\mathbf{p} = \{0.9, 0.5, 0.2, 0.1\}$:

$$r(\mathbf{p}) = 0.9 \times 0.5 \times 0.2 + 0.9 \times 0.5 \times 0.1 + 0.9 \times 0.2 \times 0.1 \\ + 0.5 \times 0.2 \times 0.1 - 3(0.9 \times 0.5 \times 0.2 \times 0.1)$$

Esto se puede colocar en una calculadora y va a dar el resultado de **0.136**, se puede construir una pequeña función en R para ahorrar calculos en el caso de que se den más conjuntos para \mathbf{p} :

```
reli_fun <- function(p1, p2, p3, p4){
  return (p1*p2*p3 + p1*p2*p4 + p1*p3*p4 + p2*p3*p4 - 3 * (p1 * p2 * p3 * p4))
}

reli_fun(0.9, 0.5, 0.2, 0.1)
```

```
## [1] 0.136
```

7.3. Estimate via simulation, for each $n = 10^2, 10^3, 10^4, 10^5$ realizations, the reliability of the system $R(p)$. Take $\{p_i\} = \{0.9, 0.5, 0.2, 0.1\}$.

```
# Given probabilities
p <- c(0.9, 0.5, 0.2, 0.1)

# Function to calculate system reliability for one trial
calculate_reliability <- function(p) {
  # Simulate the components as working (1) or failing (0) based on probabilities
  x <- rbinom(4, 1, p)

  # Apply the reliability function r(p) as given
  r <- sum(x[1]*x[2]*x[3]) + sum(x[1]*x[2]*x[4]) + sum(x[1]*x[3]*x[4]) +
    sum(x[2]*x[3]*x[4]) - 3 * prod(x)

  return(r)
}

# Function to run simulation for different number of realizations n
simulate_system_reliability <- function(n, p) {
  reliabilities <- numeric(n)

  for (i in 1:n) {
    reliabilities[i] <- calculate_reliability(p)
  }
  return(mean(reliabilities))
}

# Run simulation for each n: 10^2, 10^3, 10^4, 10^5 realizations
results <- sapply(c(10^2, 10^3, 10^4, 10^5),
  function(n) simulate_system_reliability(n, p))

# View results
names(results) <- c("n=10^2", "n=10^3", "n=10^4", "n=10^5")
print(results)

## n=10^2 n=10^3 n=10^4 n=10^5
## 0.12000 0.12600 0.14050 0.13702
```

En primer lugar se tiene un vector \mathbf{x} que simula cada componente del sistema retornando 1 o 0 con una probabilidad que es el conjunto para p_i que se define en el enunciado del problema.

Después se pasa a la función de confiabilidad. Seguido, se pasa a la función `simulate_system_reliability()` donde se simula la confiabilidad n veces y se guarda en un vector `reliabilities`. Por último, se corre la simulación para diferentes valores de n , guardando todo en `results`. Los resultados se presentan en el cuadro 2 y se ve que a medida que se hacen más simulaciones, la probabilidad de que el sistema esté funcionando se va acercando al valor obtenido analíticamente de 0.136.

n	Confiabilidad de que el sistema está funcionando.
10^2	0.12
10^3	0.126
10^4	0.1405
10^5	0.13702

Cuadro 2: Resultados de la función confiabilidad para diferentes valores de n .

7.4. Find the standard deviation of the reliability estimation for each n .

Se hace una modificación a la función `simulate_system_reliability()` para que también entregue la desviación estándar de cada estimación de la confiabilidad.

```
simulate_system_reliability <- function(n, p) {
  reliabilities <- numeric(n)

  for (i in 1:n) {
    reliabilities[i] <- calculate_reliability(p)
  }

  sd_val <- sd(reliabilities)

  return(c(SD = sd_val))
}
```

Se deja tal cual estaba la función, pero ahora retorna la desviación estándar de esa estimación para ese n en específico. Por lo tanto, ahora queda volver a correr la simulación para los diferentes n .

```
results <- sapply(c(10^2, 10^3, 10^4, 10^5), function(n) simulate_system_reliability(n,
names(results) <- c("n=10^2", "n=10^3", "n=10^4", "n=10^5")

# Print results
print(results)

##      n=10^2      n=10^3      n=10^4      n=10^5
## 0.3015113 0.3471607 0.3409914 0.3428752
```

Los resultados se muestran en el cuadro 3 reflejando cómo los sistemas varían entre ellos.

n	Desviación estándar
10^2	0.3015113
10^3	0.3471607
10^4	0.3409914
10^5	0.3428752

Cuadro 3: Desviación estándar para cada n .

7.5. Find the 95 % CI's for each n . Do the confidence intervals contain the actual reliability value from item 7.2?

Teniendo en cuenta:

$$CI = \bar{X} \pm c \times \frac{s}{\sqrt{n}}$$

Se puede volver a modificar nuevamente la función `simulate_system_reliability()` para que también retorne el intervalo superior e inferior, esto teniendo en cuenta la formula que se viene manejando durante todo el taller, por lo tanto, el código sería:

```
simulate_system_reliability <- function(n, p) {
  reliabilities <- numeric(n)
  for (i in 1:n) {
    reliabilities[i] <- calculate_reliability(p)
  }

  mean_val <- mean(reliabilities)
  sd_val <- sd(reliabilities)
  error_margin <- 1.96 * (sd_val / sqrt(n)) # 95% CI

  return(c(mean = mean_val,
           lower = mean_val - error_margin,
           upper = mean_val + error_margin))
}
```

Se ha agregado una variable de `mean_val` y el `error_margin` para luego entre ellas restar y sumar, respectivamente, para obtener el intervalo de confianza completo. Ahora, se puede volver a correr la simulación e imprimir la salida.

```
# Run simulation for each n
results <- t(sapply(c(10^2, 10^3, 10^4, 10^5),
                  function(n) simulate_system_reliability(n, p)))
```

```
rownames(results) <- c("n=10^2", "n=10^3", "n=10^4", "n=10^5")

# Print result as a table
print(round(results, 5))

##           mean   lower  upper
## n=10^2 0.15000 0.07966 0.22034
## n=10^3 0.12200 0.10170 0.14230
## n=10^4 0.13630 0.12957 0.14303
## n=10^5 0.13569 0.13357 0.13781
```

Recordando que el valor de la sección 7.2 es 0.136, y el intervalo quedó como:

$$[0.1335674, 0.1378126]$$

Por lo tanto, el valor que se encontró anteriormente sí cae en el intervalo de confianza.

8. Reading assignment

Se optó por buscar una opción gratuita del libro en PDF, esto debido a que la versión que se encuentra en la Tadeo es muy incómoda por leer, el PDF se puede encontrar libre de descargar **aquí**. Se han tomado las siguientes ideas mientras se hacía la lectura del libro:

- Las muestras que se deben elegir de una población deben ser representativas y aleatorias, porque pueden pasar casos como el de las elecciones de Landon y Roosevelt, creando un sesgo para el *dataset*.
- Lo que se entiende por un muestreo estratificado es seleccionar nuestros datos por grupos. Por ejemplo, si quiero saber los ingresos de los habitantes de Bogotá, no me puedo tomar solo muestras de la gente estrato 6, porque me quedará sesgado el análisis y voy a concluir que todos los habitantes de Bogotá ganan mucho dinero. En cambio, tengo que reconocer cómo están segmentados los estratos sociales en la ciudad; 10 % para estrato 6, 20 % para el 5, etc.² De esos grupos ahora se toman muestras aleatorias y así hemos reducido el sesgo del experimento.
- Las diferencias entre \bar{x} y μ están en que la información sobre muestras es observada, y la información sobre poblaciones grandes es a menudo inferida de muestras más pequeñas.
- “Si toturas a los datos por un tiempo largo, tarde o temprano confesarán”.
- Especifica una hipótesis y recolectar datos seguido de una aleatoriedad y principios de muestreo aleatorio asegura reducir los sesgos.

²Datos inventados por los autores

- El teorema central del límite tiene mucha importancia en los libros clásicos de estadística pero no es tan *central* en la práctica de la ciencia de datos.
- El error estándar es una métrica que mide la variabilidad de la distribución muestreada para una estadística. A medida que el tamaño de la muestra aumenta, el error estándar disminuye.
- La desviación estándar mide la variabilidad de puntos de datos individuales y el error estándar mide la variabilidad de una muestra.
- Una muestra *bootstrap* es aquella que se toma con reemplazo de un conjunto de datos observado.
- Remuestro es el proceso de tomar repetidamente muestras de los datos observados, incluyendo procedimientos como *bootstrap* y permutaciones.
- Pasos para un algoritmo *Bootstrap*:
 1. Obtener un valor de muestra, guardarlo, y luego reemplazarlo.
 2. Repetir n veces.
 3. Obtener la media de los n valores remuestreados.
 4. Repetir los pasos 1 a 3 R veces.
 5. Usar los R resultados para:
 - a) Calcular la desviación estándar.
 - b) Producir un histograma o un *boxplot*.
 - c) Encontrar el intervalo de confianza.
- El término remuestreo se refiere a hacer procedimientos de permutaciones, donde múltiples muestras son combinadas y el muestro puede ser hecho sin reemplazo. El término *bootstrapping* siempre implica hacer muestro con reemplazo.
- La mayoría de los datos no se distribuyen de manera normal y la distribución Gaussiana está sobrevalorada.
- La distribución t es como una distribución normal pero con colas más anchas.
- La distribución binomial es muy usada cuando se trata de una salida de sí/no. Se puede adaptar a varios conceptos; funciona/no funciona, comprar/no comprar, etc.

8.1. Include in your submission an exploration of the dataset *loans_income.csv* by replicating some of the codes of the sections: *Sampling Distribution of a Statistic* , *Bootstrap*, or *Confidence Intervals*.

8.1.1. *Sampling Distribution onf a Statistic*

```
library(ggplot2)

# Lectura del dataset
loans_income <- read.csv("loans_income.csv")
loans_income <- loans_income[[colnames(loans_income)]]

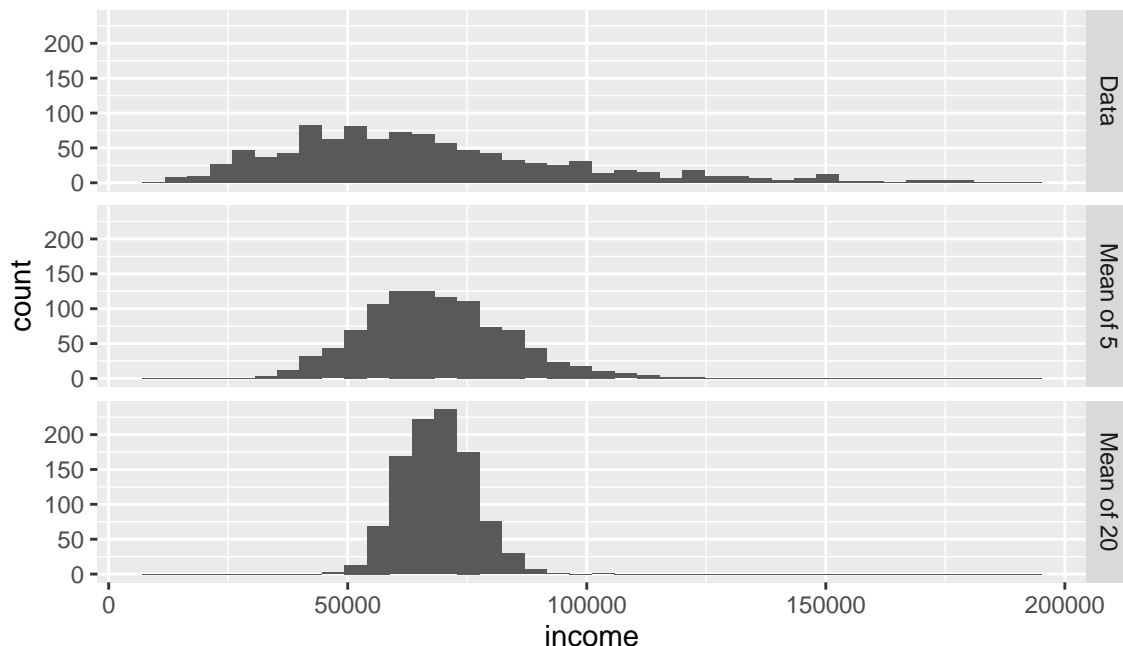
# Tomar una muestra de una variable aleatoria simple.
samp_data <- data.frame(income=sample(loans_income, 1000),
                        type = 'data_dist')

# Tomar una muestra de medias de 5 valores
samp_mean_05 <- data.frame(income = tapply(sample(loans_income, 1000*5),
                                             rep(1:1000, rep(5, 1000)), FUN=mean), type='mean_of_5')

# Tomar una muestra de medias de 20 valores
samp_mean_20 <- data.frame(income = tapply(sample(loans_income, 1000*20),
                                             rep(1:1000, rep(20, 1000)), FUN=mean), type='mean_of_20')

# Unir los data frames y convertir la columna type como factor, es decir,
# una categoría
income <- rbind(samp_data, samp_mean_05, samp_mean_20)
income$type = factor(income$type, levels=c('data_dist', 'mean_of_5',
                                           'mean_of_20'), labels=c('Data', 'Mean of 5', 'Mean of 20'))

# Graficar
ggplot(income, aes(x=income))+geom_histogram(bins=40)+facet_grid(type ~ .)
```



8.1.2. *Bootstrap*

El código anterior utiliza tal cual cómo aparece el ejemplo en el libro, solo que aquí se carga el *dataset* a la variable `loans_income`. Después crea tres muestras diferentes; `samp_data` que toma 1000 valores aleatorios directamente del *dataset* y los guarda con la etiqueta “data_dist”; `samp_mean_05` toma una muestra de 5000 valores (la operación que se hace son 1000 grupos de 5), calcula el promedio de cada grupo de 5, y se obtienen 1000 promedios; `samp_mean_20`, es parecido al anterior pero con grupos de valores de 20, se obtienen 1000 promedios de 20 valores cada uno. Después de tomar las muestras se combinan los datos en un solo *dataframe* `income` y después convierte la columna `type` en un factor categórico con etiquetas que facilitan la lectura. Por último, se grafican los histogramas, cada uno con 40 barras y se utiliza `face_grid(type ~ .)` para mostrar tres gráficos separados.

Se va a utilizar el paquete `boot` que permite realizar las operaciones para un muestreo *bootstrap* replicando 1000 veces la función `stat_fun()` al *dataset* `stat_fun`

```
library(boot)
stat_fun <- function(x, idx) median(x[idx])
boot_obj <- boot(loans_income, R=1000, statistic=stat_fun)
boot_obj

##
## ORDINARY NONPARAMETRIC BOOTSTRAP
##
##
## Call:
## boot(data = loans_income, statistic = stat_fun, R = 1000)
##
##
## Bootstrap Statistics :
##      original      bias      std. error
## t1*      62000 -82.5945      222.1515
```

La función `stat_fun()` recibe un vector `x` (los datos originales) y otro vector `idx` (índices de muestra aleatoria con reemplazo generados por el *bootstrap*), y devuelve la mediana de los valores muestreados. Después se ejecuta el *bootstrap* mediante la función `boot()` que toma el *dataset* `loans_income`, genera 1000 muestras con reemplazo, para cada muestra calcula la mediana usando la función previamente definida. Por último, muestra el valor original de la mediana, las 1000 replicaciones (medianas *bootstrap*) e información de las muestras como el bias y el error estándar.

8.1.3. *Confidence Intervals*

Se utiliza el código que se encuentra en el repositorio del libro.

```

set.seed(5)
set.seed(7)
sample20 <- sample(loans_income, 20)
sampleMean <- mean(sample20)

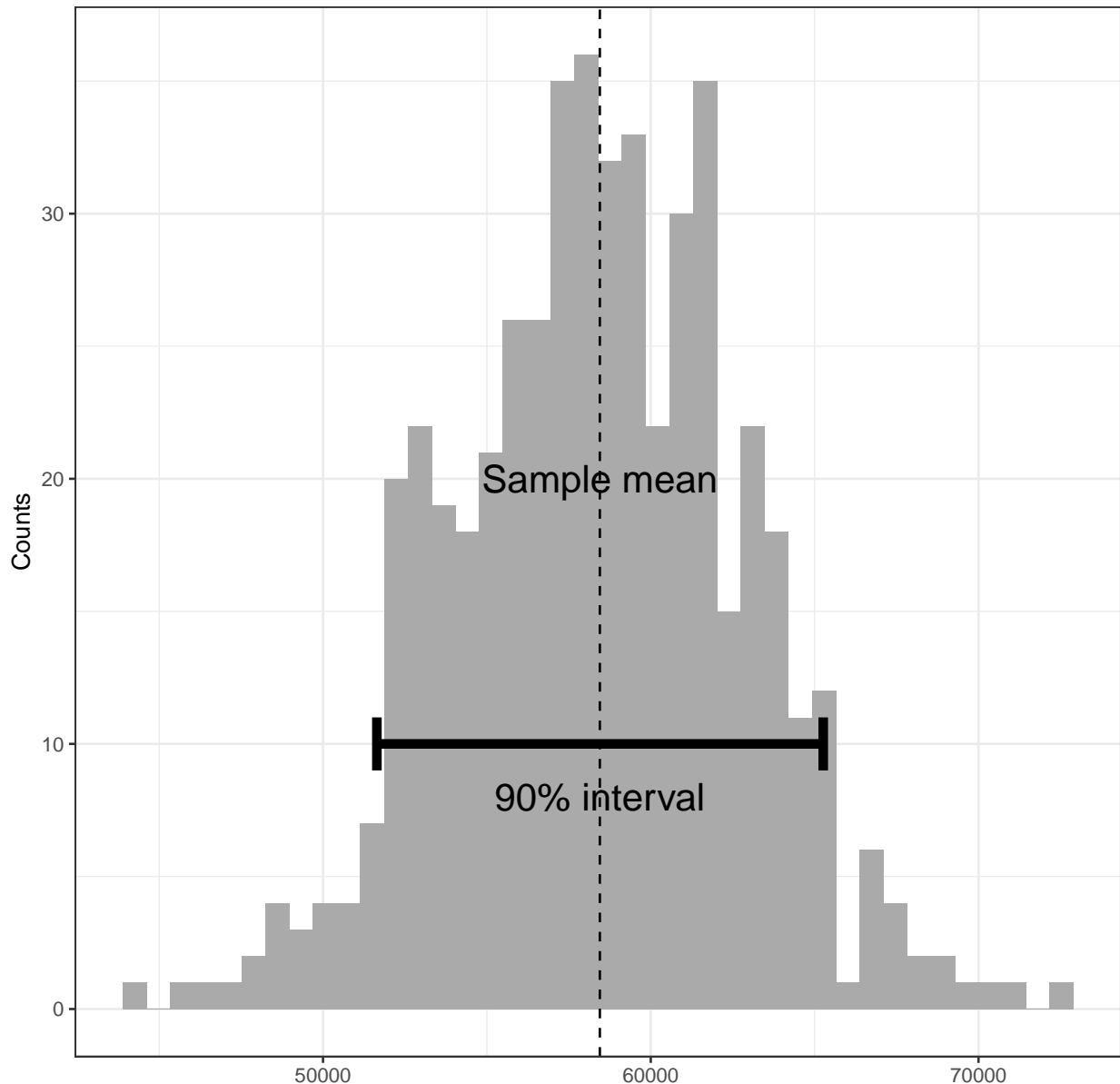
stat_fun <- function(x, idx) mean(x[idx])
boot_obj <- boot(sample20, R=500, statistic=stat_fun)
boot_ci <- boot.ci(boot_obj, conf=0.9, type='basic')
X <- data.frame(mean=boot_obj$t)
ci90 <- boot_ci$basic[4:5]
ci <- data.frame(ci=ci90, y=c(9, 11))
ci

##           ci  y
## 1 51643.09  9
## 2 65262.95 11

ggplot(X, aes(x=mean)) +
  geom_histogram(bins=40, fill='#AAAAAA') +
  geom_vline(xintercept=sampleMean, linetype=2) +
  geom_path(aes(x=ci, y=10), data=ci, size=2) +
  geom_path(aes(x=ci90[1], y=y), data=ci, size=2) +
  geom_path(aes(x=ci90[2], y=y), data=ci, size=2) +
  annotate('text', x=sampleMean, y=20, label='Sample mean', size=6) +
  annotate('text', x=sampleMean, y=8, label='90% interval', size=6) +
  theme_bw() +
  labs(x = '', y='Counts')

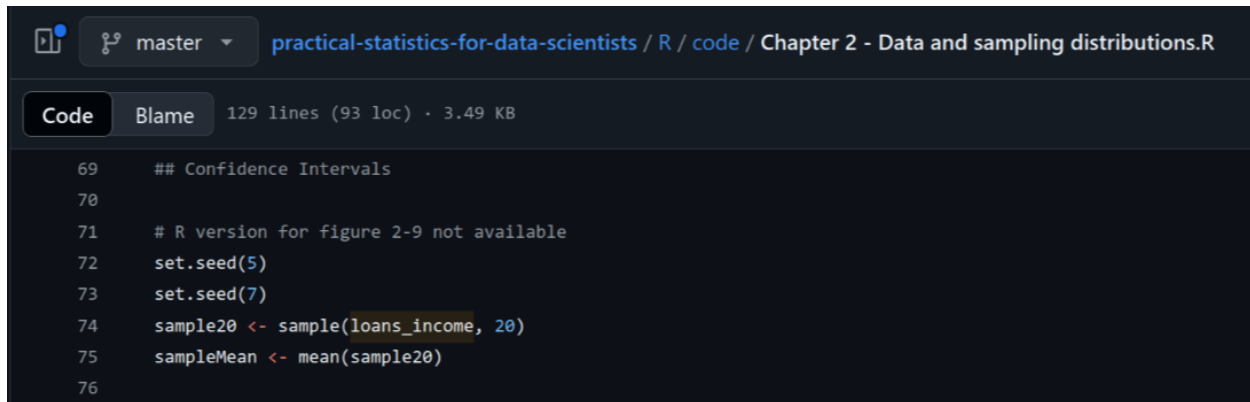
## Warning: Using 'size' aesthetic for lines was deprecated in ggplot2 3.4.0.
## i Please use 'linewidth' instead.
## This warning is displayed once every 8 hours.
## Call 'lifecycle::last_lifecycle_warnings()' to see where this warning was
## generated.

```



El código no genera la misma salida de la figura que se encuentra en el libro, esto debido a que se replicó el que está disponible en GitHub y dentro del archivo del *script* hay un comentario diciendo que la figura no está disponible, esto se puede mostrar como evidencia en la figura 1.

Por otra parte, se muestra una advertencia por el estilo de la figura, el parámetro `size` ha sido deprecado, esto quiere decir que esta edición del libro se realizó con una versión de `ggplot2` mucho más vieja que la que está actualmente.

A screenshot of a GitHub code viewer interface. At the top, there's a navigation bar with a repository icon, a dropdown menu showing 'master', and the file path 'practical-statistics-for-data-scientists / R / code / Chapter 2 - Data and sampling distributions.R'. Below this, there are tabs for 'Code' (selected) and 'Blame', followed by file statistics: '129 lines (93 loc) · 3.49 KB'. The main area displays R code with line numbers 69 to 76. The code includes comments about confidence intervals and R version compatibility, and sets two seeds (5 and 7) before sampling from 'loans_income' and calculating the mean.

```
69     ## Confidence Intervals
70
71     # R version for figure 2-9 not available
72     set.seed(5)
73     set.seed(7)
74     sample20 <- sample(loans_income, 20)
75     sampleMean <- mean(sample20)
76
```

Figura 1: *Script del libro Practical statistics for data scientists: 50+ essential concepts using R and Python*[1].

Referencias

- [1] Peter Bruce, Andrew Bruce y Peter Gedeck. *Practical statistics for data scientists: 50+ essential concepts using R and Python*. O'Reilly Media, 2020.
- [2] Sheldon M Ross. *Introduction to probability models*. Academic press, 2014.