

Trabajo Práctico de R

Denisse Blasco

Nicolás Castelao

21 de octubre de 2025

Índice

1. Ejercicio 1: Análisis de datos	2
1.1. ¿Cuál es el efecto de la localía en el resultado?	3
1.2. ¿Qué tan comunes son las remontadas?	4
1.3. ¿Cuáles son los árbitros con mayor severidad disciplinaria? ¿Hay relación con la localidad?	5
2. Ejercicio 2: Análisis econométrico	8
2.1. Ingreso per cápita de Argentina	8
2.2. Modelos de predicción y comparación	9
2.3. Co-movimiento regional del ingreso per cápita	11
2.4. Relación entre expectativa de vida total y femenina (2010)	12
2.5. Modelo simple: expectativa total vs femenina	13
2.6. Prueba t pareada: ¿la expectativa femenina es mayor?	14
2.7. Modelo múltiple con ingreso per cápita	14
2.8. Modelo alternativo: ingreso, religión y población	15
3. Ejercicio 3: Simulación y estática comparativa	16
3.1. Ingreso simulado y momentos teóricos	16
3.2. Demanda Cobb–Douglas y utilidad indirecta	18
3.3. Monte Carlo y estadísticos descriptivos	18
3.4. Probabilidad de bajo consumo	20
3.5. Shock de precio y comparación de distribuciones	21
3.6. Visualización: antes vs. después del shock	21
3.7. Heterogeneidad en preferencias	22

1. Ejercicio 1: Análisis de datos

Del sitio web <https://www.football-data.co.uk/englandm.php> extraimos datos sobre la Premier League de fútbol.

Nuestras preguntas de interés son:

- ¿Cuál es el efecto de la localía en el resultado?
- ¿Qué tan comunes son las remontadas?
- ¿Cuáles son los árbitros con mayor severidad disciplinaria? ¿Hay relación con la localidad?

Para esto, utilizamos datasets de partidos de fútbol de la Premier League. Restringimos la información desde la temporada 2010-2011 a 2023-2024, teniendo en cuenta cambios en el desarrollo del deporte y lo incompleto de los datasets futuros. Para la pregunta acerca de la localía, separadamente realizamos un análisis sin tener en cuenta los años de pandemia, donde no hubo público local.

Primero, realizamos las manipulaciones correspondientes a los datasets, explicadas con comentarios:

```
1 # Tenemos un dataset por cada temporada. Queremos juntarlos.
2 # El dataset t_i corresponde a la temporada del año i-i+1. Es decir, t_
  2023 tiene la temporada 2023-2024
3 t_2010 <- read.csv("C:/Users/gusta/Downloads/2010.csv", header = TRUE)
4 t_2011 <- read.csv("C:/Users/gusta/Downloads/2011.csv", header = TRUE)
5 t_2012 <- read.csv("C:/Users/gusta/Downloads/2012.csv", header = TRUE)
6 t_2013 <- read.csv("C:/Users/gusta/Downloads/2013.csv", header = TRUE)
7 t_2014 <- read.csv("C:/Users/gusta/Downloads/2014.csv", header = TRUE)
8 t_2015 <- read.csv("C:/Users/gusta/Downloads/2015.csv", header = TRUE)
9 t_2016 <- read.csv("C:/Users/gusta/Downloads/2016.csv", header = TRUE)
10 t_2017 <- read.csv("C:/Users/gusta/Downloads/2017.csv", header = TRUE)
11 t_2018 <- read.csv("C:/Users/gusta/Downloads/2018.csv", header = TRUE)
12 t_2019 <- read.csv("C:/Users/gusta/Downloads/2019.csv", header = TRUE)
13 t_2020 <- read.csv("C:/Users/gusta/Downloads/2020.csv", header = TRUE)
14 t_2021 <- read.csv("C:/Users/gusta/Downloads/2021.csv", header = TRUE)
15 t_2022 <- read.csv("C:/Users/gusta/Downloads/2022.csv", header = TRUE)
16 t_2023 <- read.csv("C:/Users/gusta/Downloads/2023.csv", header = TRUE)
17
18
19 # Nos quedaremos con las variables de interés.
20
21 vars_interes <- c(
22   "FTR", # resultado final
23   "HTR", # resultado al entretiempo.
24   "Referee", "HY", "AY", "HR", "AR" # disciplina. No ajustamos por cantidad de
    faltas pues suponemos que la elección de árbitros y la intensidad
    del juego no están correlacionadas.
25 )
26
27 df_2010 <- t_2010[, vars_interes]
28 df_2011 <- t_2011[, vars_interes]
29 df_2012 <- t_2012[, vars_interes]
30 df_2013 <- t_2013[, vars_interes]
31 df_2014 <- t_2014[, vars_interes]
32 df_2015 <- t_2015[, vars_interes]
```

```

33 df_2016 <- t_2016[, vars_interes]
34 df_2017 <- t_2017[, vars_interes]
35 df_2018 <- t_2018[, vars_interes]
36 df_2019 <- t_2019[, vars_interes]
37 df_2020 <- t_2020[, vars_interes]
38 df_2021 <- t_2021[, vars_interes]
39 df_2022 <- t_2022[, vars_interes]
40 df_2023 <- t_2023[, vars_interes]
41
42 # Los uno
43 data <- rbind(
44   df_2010,
45   df_2011,
46   df_2012,
47   df_2013,
48   df_2014,
49   df_2015,
50   df_2016,
51   df_2017,
52   df_2018,
53   df_2019,
54   df_2020,
55   df_2021,
56   df_2022,
57   df_2023
58 )
59
60 # Le saco la única fila con missing data
61 data <- data[-1901, ]
62 sum(is.na(data))

```

Luego, comenzamos a responder cada pregunta.

1.1. ¿Cuál es el efecto de la localía en el resultado?

Para comprender mejor el fenómeno, grafico la frecuencia de cada resultado a través del siguiente código:

```

1 barplot(table(data$FTR),
2         main = "Frecuencia de resultados ",
3         xlab = "Resultado final", ylab = "Frecuencia",
4         col = c(rgb(0.7, 0.85, 1), rgb(0.4, 0.6, 0.9), rgb(0.1, 0.3, 0.7))
5         ),
6         border = "black")

```

El gráfico resultante es

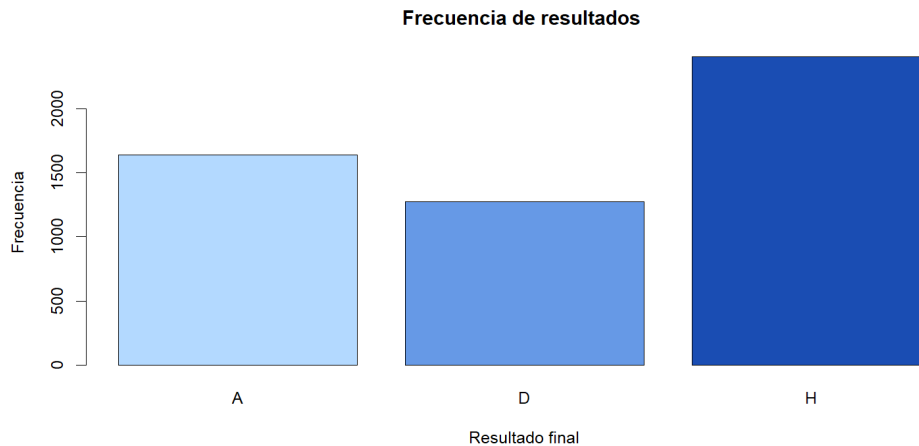


Figura 1: Gráfico de resultados de la Premier League.

Notemos que, dado que hay igual cantidad de partidos de visitante que de local (cada partido tiene un visitante y un local en la Premier), entonces las probabilidades pueden calcularse de la siguiente manera.

$$P(\text{ganar} \mid \text{soy local}) = \frac{\# \text{victorias como local}}{\# \text{partidos totales}}$$

$$P(\text{ganar} \mid \text{soy visitante}) = \frac{\# \text{victorias como visitante}}{\# \text{partidos totales}}$$

$$P(\text{empatar} \mid \text{soy local}) = \frac{\# \text{empates como local}}{\# \text{partidos totales}} = P(\text{empatar} \mid \text{siendo visitante}) = P(\text{empatar})$$

Calculemos todo esto.

```

1 p_ganar_local = sum(data$FTR == "H") / nrow(data)
2 p_ganar_visita = sum(data$FTR == "A") / nrow(data)
3 p_empate = sum(data$FTR == "D") / nrow(data)

```

Lo observado es que en el 45 % de los partidos gana el local, mientras que el visitante solo gana en el 30 % de ellos. El empate ocurre en el 24 % de los partidos.

1.2. ¿Qué tan comunes son las remontadas?

Para ver esto, agregamos una columna que valga 1 si hubo remontada y valga 0 si no. Se considera remontada un partido que se estaba perdiendo en el primer tiempo y se termina ganando hacia el final del partido. Además, grafico la frecuencia de las remontadas. Esto se hace a través del siguiente código:

```

1 data$remontada = as.integer((data$HTR == "H" & data$FTR == "A") | (
2     data$HTR == "A" & data$FTR == "H"))
3 barplot(table(data$remontada),
4     main = "Frecuencia de remontadas ",
5     xlab = "Remontada", ylab = "Frecuencia",

```

```

6   col = c(rgb(0.7, 0.85, 1), rgb(0.1, 0.3, 0.7) ),
7   border = "black")

```

Y el gráfico resultante es:

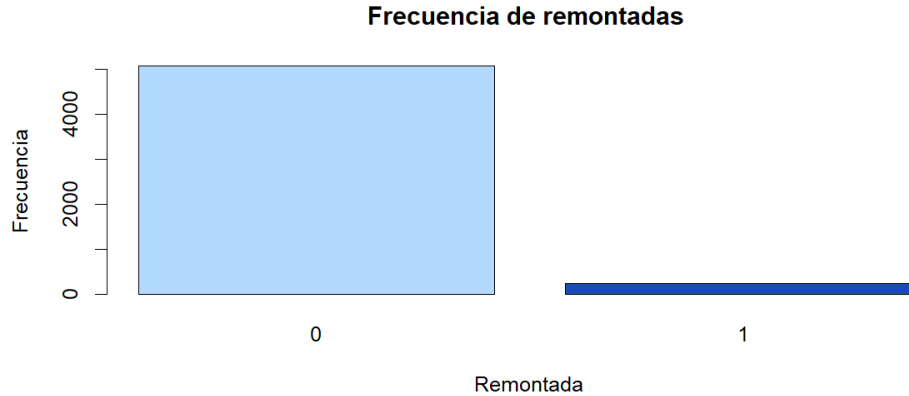


Figura 2: Gráfico de remontadas de la Premier League.

Es decir, se remontan aproximadamente 4.5 % de los partidos, ya que hay 238 remontadas y 5082 partidos que no fueron remontados.

Ahora calculamos la probabilidad de remontar dado que llegamos al entretiempo perdiendo, tanto para el local como para el visitante:

```

1   cat("P(remontar | perdía el local al ET) = ",round(100 * (sum(data$HTR
2   == "A" & data$FTR == "H") / sum(data$HTR == "A")), 4), "%\n")
3   cat("P(remontar | perdía el visitante al ET) = ",round(100 * (sum(data$HTR
4   == "H" & data$FTR == "A") / sum(data$HTR == "H")), 4), "%\n")

```

Y el resultado es que el local tiene el doble de probabilidades de remontar si llega perdiendo al entretiempo, siendo de 10.3139 % contra 5.4915 %.

1.3. ¿Cuáles son los árbitros con mayor severidad disciplinaria? ¿Hay relación con la localidad?

Hay 50 árbitros distintos, lo cual se obtiene usando:

```

1   arbitros <- unique(data$Referee)
2   n_arbitros <- length(arbitros)

```

Armo un índice de severidad disciplinaria, asignando 10 = tarjeta amarilla y 25 = tarjeta roja. También junto otras columnas.

```

1   data$disciplina <- 10*(data$HY + data$AY) + 25*(data$HR + data$AR)
2   data$tot_yellow <- data$HY + data$AY
3
4   arbitros_df <- data.frame(Referee = arbitros, disciplina = numeric(n_
5   arbitros))
6   for (i in 1:n_arbitros){
7     arbitros_df$disciplina[i] = mean(data$disciplina[data$Referee ==
8     arbitros_df$Referee[i]])

```

```
7 }
```

Grafico los 15 árbitros más severos:

```
1 top_severos <- order(arbitros_df$disciplina, decreasing = TRUE)[1:15]
2 barplot(
3   arbitros_df$disciplina[top_severos],
4   col = rgb(0.1,0.83,0.7,0.85),
5   main = "Árbitros con mayor severidad disciplinaria",
6   ylab = "Promedio por partido",
7   names.arg = arbitros_df$Referee[top_severos],
8   las = 2, cex.names = 0.7, border = "black"
9 )
```

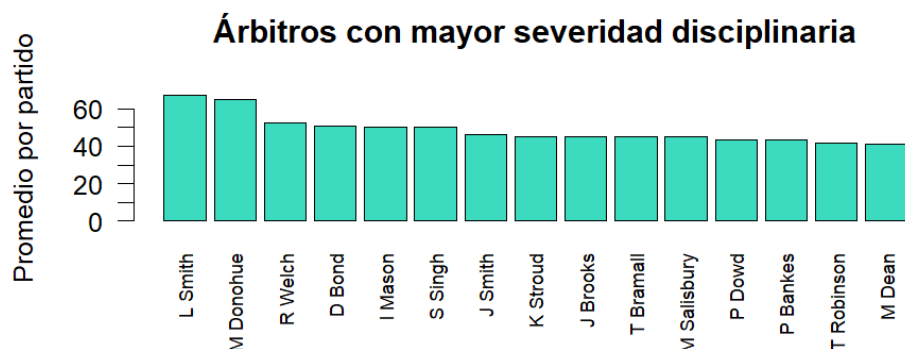


Figura 3: Árbitros más severos.

Veamos ahora información sobre la cantidad de amarillas según el resultado, usando el siguiente código:

```
1 y_H <- data$tot_yellow[data$FTR == "H"]
2 y_A <- data$tot_yellow[data$FTR == "A"]
3 y_D <- data$tot_yellow[data$FTR == "D"]
4
5 hist(y_H,
6   col = rgb(0.10, 0.55, 0.10, 0.35), border = rgb(0.10, 0.55, 0.10,
7   0.9),
8   main = "Distribución de amarillas (HY+AY) por resultado",
9   xlab = "Amarillas totales")
10 hist(y_D, add = TRUE,
11   col = rgb(0.35, 0.70, 0.35, 0.35), border = rgb(0.35, 0.70, 0.35,
12   0.9))
13 hist(y_A, add = TRUE,
14   col = rgb(0.40, 0.60, 0.90, 0.35), border = rgb(0.40,0.60,0.90,0.9))
15 legend("topright",
16   fill = c(rgb(0.10, 0.55, 0.10, 0.35), rgb(0.35, 0.70, 0.35, 0.35),
17   rgb(0.40, 0.60, 0.90, 0.35)),
18   border = c(rgb(0.10, 0.55, 0.10, 0.9), rgb(0.35, 0.70, 0.35, 0.9),
19   rgb(0.40,0.60,0.90,0.9)),
20   legend = c("Gana local", "Empate", "Gana visita"),
21   bty = "n")
```

Gráficamente se observa que hay una mayor cantidad de amarillas cuando gana el local, intermedia cuando gana la visita y menor cantidad en partidos que resultan en empate:

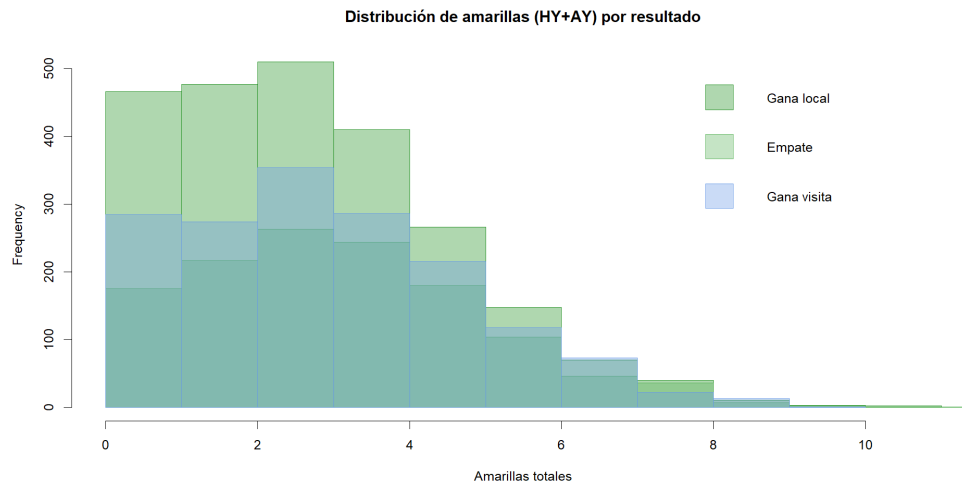


Figura 4: Amarillas según resultado.

Ahora analizamos el efecto de la localidad en las tarjetas amarillas, codificando según:

```

1 plot(jitter(data$HY, amount = 0.15), #dispersa los puntos
      aleatoriamente para que se puedan visualizar mejor
2 jitter(data$AY, amount = 0.15),
3 pch = 16, col = rgb(0.10,0.30,0.70,0.25),
4 xlab = "Amarillas del local",
5 ylab = "Amarillas del visitante",
6 main = "Amarillas: local vs visitante")
7 abline(0, 1, col = "red", lwd = 2, lty = 2) # La linea que refleja la
  imparcialidad.

```

El gráfico resultante no muestra un sesgo claro hacia la localidad o la visita, pero sí se ve que una mayor cantidad de amarillas para un equipo correlaciona con una mayor cantidad de amarillas para el otro:

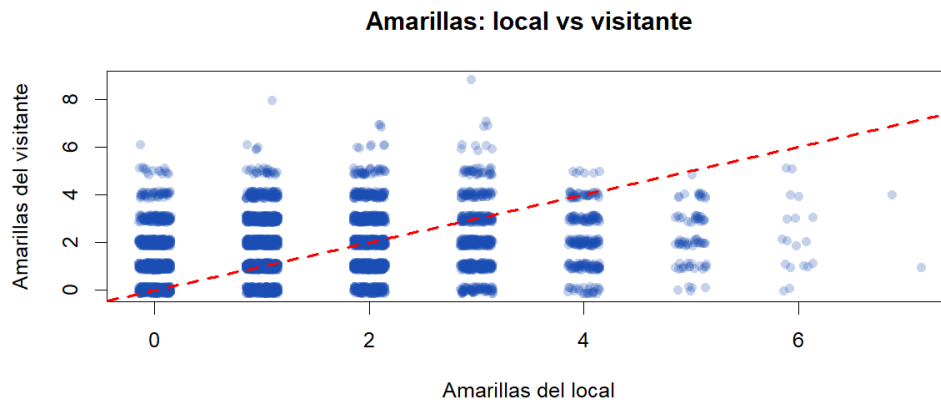


Figura 5: Localidad en tarjetas amarillas.

2. Ejercicio 2: Análisis econométrico

Para este ejercicio, trabajamos con el dataset `gapminder.csv`, que contiene información de distintos países sobre ingreso per cápita, expectativa de vida, religión y otros indicadores. El análisis se centra en Argentina y luego se amplía al estudio comparativo regional y global.

```
1 data <- read.csv("C:/Users/gusta/Downloads/gapminder.csv", header = TRUE)
2 library(ggplot2)
```

2.1. Ingreso per cápita de Argentina

Filtramos los datos correspondientes a Argentina y graficamos la evolución del ingreso per cápita a lo largo de los años:

```
1 arg <- subset(data, country=="Argentina")
2
3 ggplot(arg, aes(x = year, y = income_per_person)) +
4   geom_line() +
5   geom_point() +
6   labs(title = "Ingreso per cápita argentino",
7         x = "Año", y = "Ingreso per cápita") +
8   theme_minimal()
```

El gráfico obtenido se muestra a continuación:



Figura 6: Ingreso per cápita argentino a lo largo del tiempo.

Se observa que antes de 1875 la tendencia era creciente. Luego, entre 1875 y 1990, el ingreso per cápita disminuyó. A partir de 1990 volvió a crecer sostenidamente hasta 1998, cayendo nuevamente en 2002. Desde entonces, la tendencia vuelve a ser positiva, con una leve caída en 2009.

2.2. Modelos de predicción y comparación

Separamos los datos en conjuntos de entrenamiento y prueba, y ajustamos tres modelos: lineal, polinómico de grado 2 y polinómico de grado 10.

```
1 train <- arg[1:(nrow(arg) - 10), ]
2 test  <- arg[(nrow(arg)-9):nrow(arg),]
3
4 lineal <- lm(income_per_person ~ year, data = train)
5 pol2   <- lm(income_per_person ~ poly(year, 2, raw = TRUE), data = train)
6 pol10  <- lm(income_per_person ~ poly(year, 10, raw = TRUE), data = train)
```

Luego, comparamos los modelos utilizando el RMSE en el conjunto de testeo:

```
1 rmse <- function(obs, pred) sqrt(mean((obs - pred)^2))
2
3 pred_lin   <- predict(lineal, newdata = test)
4 pred_pol2  <- predict(pol2, newdata = test)
5 pred_pol10 <- predict(pol10, newdata = test)
6
7 rmse_lin   <- rmse(test$income_per_person, pred_lin)
8 rmse_pol2  <- rmse(test$income_per_person, pred_pol2)
9 rmse_pol10 <- rmse(test$income_per_person, pred_pol10)
```

El RMSE en el dataset para el testeo de los modelos es de 2012.47 para el modelo lineal, 3114.47 para el cuadrático y 11311.22 para el polinomio de grado 10.

Visualmente, los tres modelos ajustados sobre toda la serie lucen así:

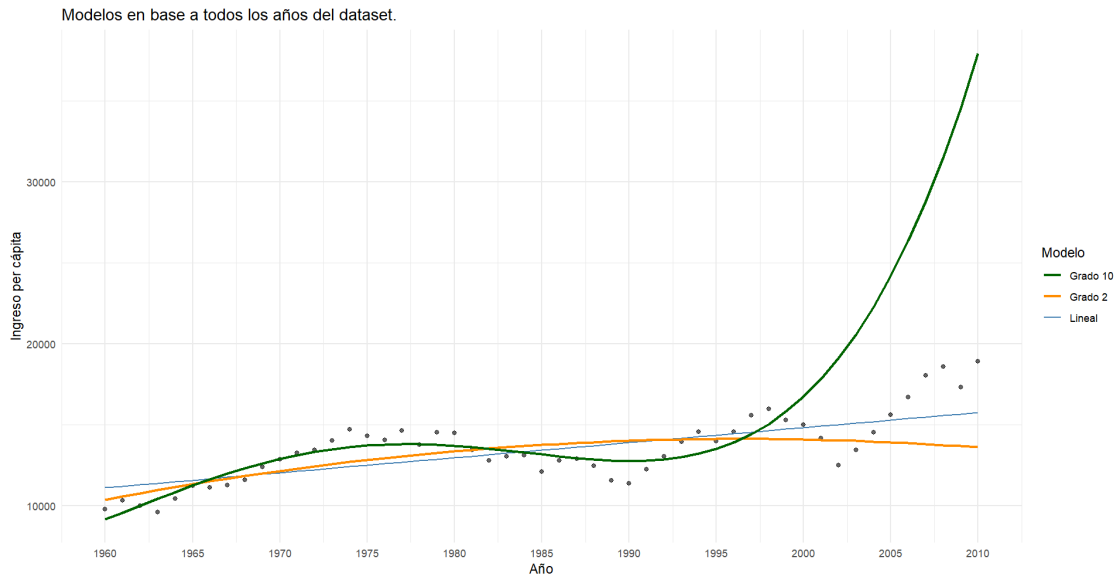


Figura 7: Ajuste de modelos sobre el ingreso per cápita argentino.

Y las predicciones sobre los últimos 10 años (conjunto de test) se muestran en la siguiente figura:

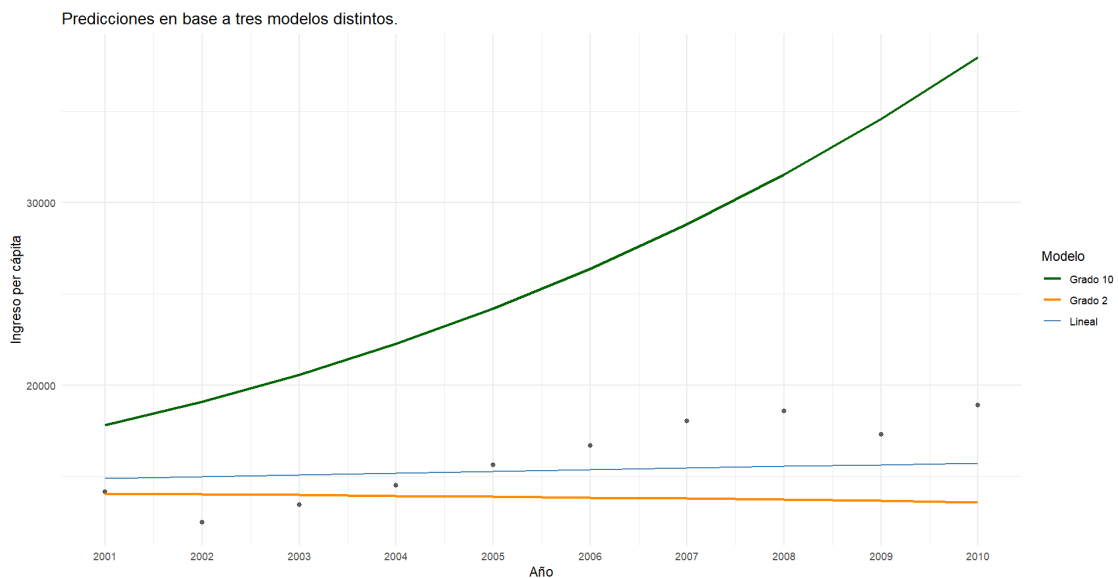


Figura 8: Predicciones sobre el conjunto de testeo (últimos 10 años).

Todos estos gráficos se obtienen mediante las siguientes líneas:

```

1  grid <- data.frame(year = arg$year)
2  grid$y_lin   <- predict(lineal,  newdata = grid)
3  grid$y_pol2  <- predict(pol2,   newdata = grid)
4  grid$y_pol10 <- predict(pol10,  newdata = grid)
5
6  ggplot() +
7    geom_point(data = arg, aes(year, income_per_person), alpha = 0.6) +

```

```

8   geom_line(data = grid, aes(year, y_lin,   color = "Lineal")) +
9   geom_line(data = grid, aes(year, y_pol2,   color = "Grado 2"), linewidth
    = 1) +
10  geom_line(data = grid, aes(year, y_pol10, color = "Grado 10"),
    linewidth = 1) +
11  labs(
12    title = "Modelos en base a todos los años del dataset.",
13    x = "Año", y = "Ingreso per cápita", color = "Modelo"
14  ) +
15  scale_x_continuous(breaks = seq(min(arg$year), max(arg$year), by = 5)) +
16  scale_color_manual(values = c("Lineal" = "steelblue", "Grado 2" = "
    darkorange", "Grado 10" = "darkgreen")) +
17  theme_minimal()
18
19
20 # Ahora veo las predicciones sobre el conjunto de testeo
21 grid2 <- data.frame(year = test$year)
22 grid2$y_lin   <- predict(lineal,   newdata = grid2)
23 grid2$y_pol2  <- predict(pol2,    newdata = grid2)
24 grid2$y_pol10 <- predict(pol10,   newdata = grid2)
25
26 ggplot() +
27   geom_point(data = test, aes(year, income_per_person), alpha = 0.6) +
28   geom_line(data = grid2, aes(year, y_lin,   color = "Lineal")) +
29   geom_line(data = grid2, aes(year, y_pol2,   color = "Grado 2"), linewidth
    = 1) +
30   geom_line(data = grid2, aes(year, y_pol10, color = "Grado 10"),
    linewidth = 1) +
31   labs(
32     title = "Predicciones en base a tres modelos distintos.",
33     x = "Año", y = "Ingreso per cápita", color = "Modelo"
34   ) +
35   scale_x_continuous(breaks = seq(min(test$year), max(test$year), by = 1))
    +
36   scale_color_manual(values = c("Lineal" = "steelblue", "Grado 2" = "
    darkorange", "Grado 10" = "darkgreen")) +
37   theme_minimal()

```

2.3. Co-movimiento regional del ingreso per cápita

Se analizan los países Argentina, Brasil, Chile, Perú y Uruguay, construyendo una matriz con sus ingresos per cápita para medir correlaciones tanto en niveles como en variaciones porcentuales.

```

1  library(tidyverse)
2  paises <- c("Argentina", "Brazil", "Chile", "Peru", "Uruguay")
3
4  paises_df <- data %>%
5    filter(country %in% paises) %>%
6    select(year, country, income_per_person) %>%
7    pivot_wider(names_from = country, values_from = income_per_person) %>%
8    arrange(year)
9
10 correlaciones <- cor(paises_df, use="pairwise.complete.obs")

```

```

11
12
13 paises_df_pct <- paises_df %>%
14   mutate(across(-year,
15     ~ . / dplyr::lag(.), #. para hablar de la columna actual,
16       lag para hacer referencia al año anterior, la división
17       para calcular la tasa.
18     .names = "{.col}_pct"))
19
20 paises_df_pct <- paises_df_pct[2:nrow(paises_df_pct),-c(2:6)]
21
22 correlaciones_pct <- cor(paises_df_pct)[2:nrow(cor(paises_df_pct)),2:ncol(
23   cor(paises_df_pct))]

```

Las tablas de correlaciones son las siguientes:

Tabla 1: Matriz de correlaciones entre países elegidos (niveles de ingreso per cápita).

	Argentina	Brazil	Chile	Peru	Uruguay
Argentina	1.0000	0.7951	0.7650	0.7817	0.8292
Brazil	0.7951	1.0000	0.7717	0.5649	0.8713
Chile	0.7650	0.7717	1.0000	0.5485	0.9408
Peru	0.7817	0.5649	0.5485	1.0000	0.5773
Uruguay	0.8292	0.8713	0.9408	0.5773	1.0000

Tabla 2: Matriz de correlaciones entre países elegidos (variaciones porcentuales del ingreso per cápita).

	Argentina_pct	Brazil_pct	Chile_pct	Peru_pct	Uruguay_pct
Argentina_pct	1.0000	0.2721	0.1692	0.3553	0.5128
Brazil_pct	0.2721	1.0000	0.0061	0.4105	0.2780
Chile_pct	0.1692	0.0061	1.0000	0.0432	0.3655
Peru_pct	0.3553	0.4105	0.0432	1.0000	0.4127
Uruguay_pct	0.5128	0.2780	0.3655	0.4127	1.0000

El resultado muestra correlaciones altas y positivas en niveles, indicando que las economías tienden a moverse conjuntamente en el largo plazo. En cambio, al analizar las variaciones porcentuales anuales, las correlaciones se reducen y se vuelven más heterogéneas, lo que refleja diferencias en los ciclos económicos de corto plazo.

2.4. Relación entre expectativa de vida total y femenina (2010)

Seleccionamos el año 2010 y graficamos la relación entre expectativa de vida total y femenina:

```

1 year_df <- subset(data, year == 2010)
2 year_df <- year_df %>% drop_na(life_expectancy)
3 year_df <- year_df[year_df$life_expectancy_female != "-",]
4 year_df$life_expectancy_female <- as.numeric(year_df$life_expectancy_
5   female)

```

```

5
6 ggplot(year_df, aes(x = life_expectancy_female, y = life_expectancy)) +
7   geom_point() +
8   geom_abline(color = "red") +
9   geom_smooth(method = lm, color = "blue") +
10  theme_minimal()

```

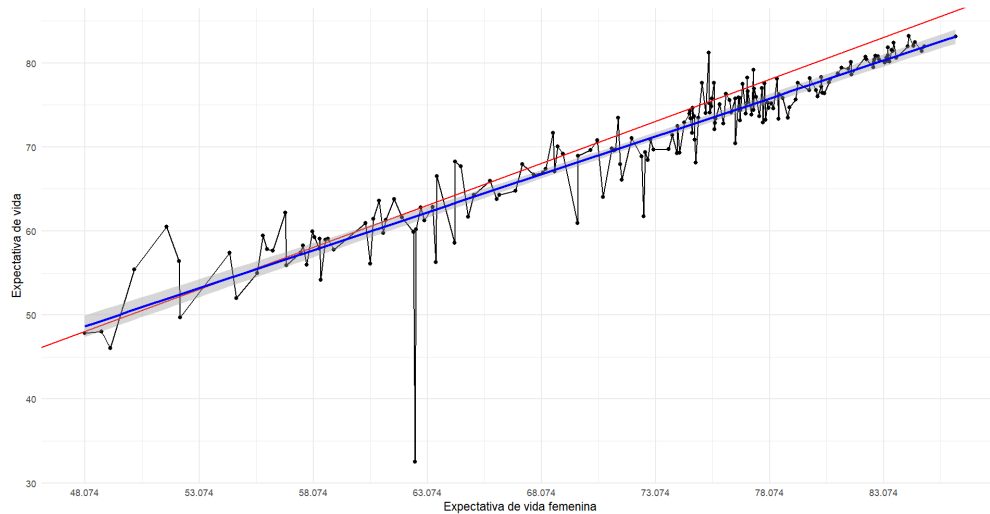


Figura 9: Expectativa de vida total vs. femenina (año 2010).

La recta azul es la regresión lineal, y la recta roja es la recta $y = x$. La relación es claramente positiva, con pendiente menor que 1, indicando que los hombres viven menos en promedio.

2.5. Modelo simple: expectativa total vs femenina

Creo el modelo a través del siguiente código:

```

1 modelo <- lm(life_expectancy ~ life_expectancy_female, data = year_df)
2 summary(modelo)
3
4 beta0_hat <- summary(modelo)$coefficients["(Intercept)", "Estimate"]
5 beta1_hat <- summary(modelo)$coefficients["life_expectancy_female", "
  Estimate"]
6 R2 <- summary(modelo)$r.squared

```

El modelo lineal simple arroja:

$$\hat{\beta}_0 = 5,22, \quad \hat{\beta}_1 = 0,90, \quad R^2 = 0,874$$

Que el coeficiente β_1 sea 0.9037 indica que cuando la expectativa de vida femenina aumenta 1 año, la expectativa de vida total aumenta 0.9 años aproximadamente. El intercepto no tiene interpretación causal, solo ajusta el modelo. El R^2 es muy alto, lo cual dice que la expectativa de vida femenina explica en gran parte cambios en el nivel de expectativa de vida total.

2.6. Prueba t pareada: ¿la expectativa femenina es mayor?

Se realiza el test con H_0 : `life_expectancy_female = life_expectancy` y H_1 : `life_expectancy_female > life_expectancy`. Por esa misma razón, uso un test pareado a la derecha unilateralmente.

```
1 t.test(year_df$life_expectancy_female,
2        year_df$life_expectancy,
3        paired = TRUE,
4        alternative = "greater")
5 print(tt)
```

El t-test pareado unilateral entre female y total dio $t = 7.105$, p muy cercano a cero. La diferencia media (female – total) de 1.72 años con un intervalo de confianza del 95 %. Como el intervalo queda estrictamente por encima de 0 y el p -valor es demasiado pequeño, rechazamos la hipótesis nula y concluimos que la esperanza de vida femenina es mayor que la total en estos datos.

2.7. Modelo múltiple con ingreso per cápita

El modelo es creado a través del siguiente código:

```
1 modelo_multiple <- lm(life_expectancy ~
2                        life_expectancy_female + income_per_person,
3                        data = year_df)
4 print(summary(modelo_multiple))
```

Y el resumen del modelo es como sigue:

Tabla 3: Resumen del modelo de regresión múltiple: expectativa de vida total vs. femenina e ingreso per cápita.

Variable	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	7.403	2.111	3.508	0.00057***
life_expectancy_female	0.8663	0.0310	27.980	<2e-16***
income_per_person	0.0000320	0.0000154	2.077	0.03919*
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1				

Residual standard error: 3.148 on 181 df

Multiple R-squared: 0.8772 *Adjusted R-squared:* 0.8759

F-statistic: 646.6 on 2 and 181 DF, *p-value:* < 2.2e-16

Vemos que el coeficiente asociado a la expectativa de vida femenina es 0.866 y es significativo. Esto significa que un año más de expectativa de vida femenina se asocia a 0.866 años más de expectativa de vida total.

El coeficiente asociado al ingreso per cápita es positivo pero casi cero, con cierta significatividad. Suponiendo que el ingreso per cápita está en miles de dolares, esto significa que un aumento en 1000 dolares al ingreso per cápita produce un aumento positivo pero casi cero en años de expectativa de vida total.

El R^2 de esta regresión es 0.8772, por lo que las variables independientes explican con fortaleza a la variable dependiente.

Al controlar por ingreso, el coeficiente de la expectativa de vida femenina disminuye, pasando de 0.904 a 0.866. Esto muestra que parte del efecto anteriormente atribuido a la expectativa de vida femenina en realidad era explicado por el ingreso per cápita. Sin embargo, el efecto de la expectativa de vida femenina permanece siendo grande y significativo.

Luego, vemos que incluir la variable de ingreso per cápita mejora la estimación central y aumenta el R^2 del modelo. Sin embargo, el aumento del R^2 es muy chico, por lo que a fines prácticos podría ser un problema incluir esta variable en la regresión.

2.8. Modelo alternativo: ingreso, religión y población

Las covariables a elegir son: logaritmo del ingreso per cápita, la población y la religión principal. Hago el modelo:

```

1 year_df$log_income <- log(year_df$income_per_person)
2
3 year_df$main_religion <- trimws(year_df$main_religion)
4 year_df$main_religion <- gsub("\\s+", " ", year_df$main_religion)
5 tmp <- tolower(year_df$main_religion)
6 tmp[tmp %in% c("", "na", "n/a")] <- NA
7 tmp[tmp %in% c("christian")] <- "Christian"
8 tmp[tmp %in% c("muslim")] <- "Muslim"
9 tmp[tmp %in% c("eastern religions")] <- "Eastern religions"
10 year_df$main_religion <- factor(tmp)
11
12 modelo_propio <- lm(life_expectancy ~ log_income + main_religion +
13                       population,
14                       data = year_df)
15 summary(modelo_propio)

```

El resumen del modelo se muestra en la siguiente tabla:

Tabla 4: Modelo alternativo: expectativa de vida total vs. log(ingreso per cápita), religión y población.

Variable	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	18.480	3.532	5.232	5.13e-07***
log_income	5.691	0.376	15.146	<2e-16***
main_religionEastern religions	3.312	1.872	1.770	0.0787.
main_religionMuslim	0.673	1.007	0.668	0.5048
population	$-4,274 \times 10^{-11}$	$3,264 \times 10^{-9}$	-0,013	0.9896
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1				

Residual standard error: 5.619 on 162 df (17 observaciones eliminadas por faltantes)

Multiple R-squared: 0.5901 Adjusted R-squared: 0.5800

F-statistic: 58.3 on 4 and 162 DF, p-value: < 2.2e-16

Vemos que el coeficiente asociado al logaritmo natural del ingreso per cápita es 5.691. Esto significa que un aumento en una unidad del log_income genera un aumento en 5.691 años de la expectativa de vida total. Este coeficiente es altamente significativo. Esto significa que un aumento

del 1 % en el ingreso per cápita esta asociado con 0.0569 años más de expectativa de vida, cercano a 21 días.

El coeficiente asociado a la población es $-4,274,10^{-11}$, negativo pero muy cercano a cero. Aunque esto no es significativo, por lo que podemos asumir que la estimación es cero. Esto significa que no hay relación entre la cantidad de población y la expectativa de vida total según nuestros datos.

Los coeficientes de las religiones usan al cristianismo como religión base. Luego, que el coeficiente de religiones del Este sea 3.312, con cierta significatividad, significa que hay un efecto positivo marginal con respecto al cristianismo en la expectativa de vida. El coeficiente de los musulmanes es 0.673 sin significatividad. Luego, no podemos decir mucho de este regresor.

El R^2 de la regresión es de 0.5901, lo cual muestra que los regresores elegidos solamente explican el 59 % de la variación en la expectativa de vida.

Además, hubo 17 observaciones eliminadas por faltantes en alguna variable.

3. Ejercicio 3: Simulación y estática comparativa

En este ejercicio simulamos ingresos a partir de una distribución $\chi^2(k)$, analizamos una demanda tipo Cobb–Douglas y estudiamos cómo cambian las distribuciones de consumo y la utilidad indirecta ante un shock de precios y con heterogeneidad de preferencias.

3.1. Ingreso simulado y momentos teóricos

Definimos una función para simular ingresos $Y \sim \chi^2(k)$. Teóricamente, $E[Y] = k$ y $\text{Var}(Y) = 2k$, por lo que k determina el nivel y la dispersión. Luego, comparamos histogramas empíricos con la densidad teórica para distintos k .

```
1  simular_ingreso <- function(k,n){
2    Y <- rchisq(n, df = k)
3    return(Y)
4  }
5
6  # Teoricamente, los valores son
7  media_teo <- k
8  varianza_teo <- 2*k
9  sd_teo <- sqrt(varianza_teo)
10
11 # Veamos cómo cambia para distintos valores de k. Se realiza un gráfico
    para cada k.
12 ks = c(5,10,20)
13
14 for (k in ks) {
15   Y <- simular_ingreso(k,1000)
16   hist(Y, freq = FALSE, col = "lightblue",
17        main = paste("Ingreso ~ Chi-cuadrado(k =", k, ")"),
18        xlab = "Y")
19   curve(dchisq(x, df = k), add = TRUE, col = "red", lwd = 2)
20 }
```

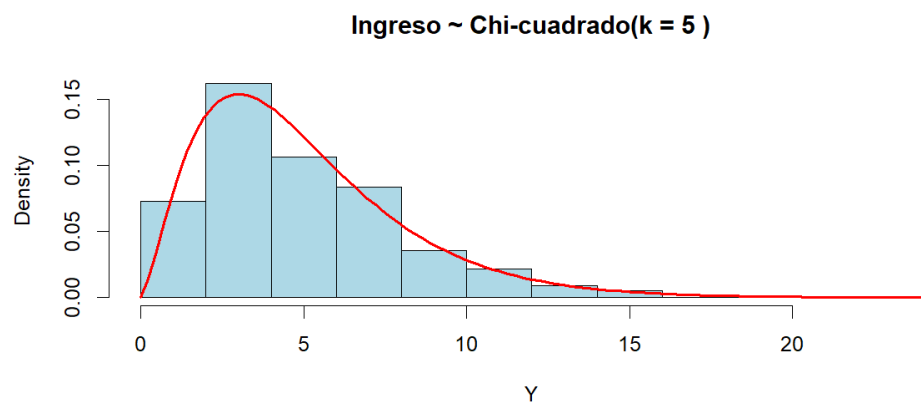



Figura 10: Ingreso simulado vs densidad teórica ($k = 5$).

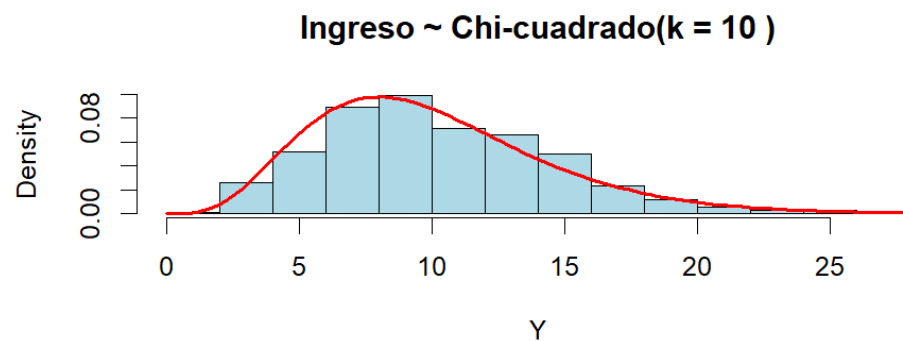


Figura 11: Ingreso simulado vs densidad teórica ($k = 10$).

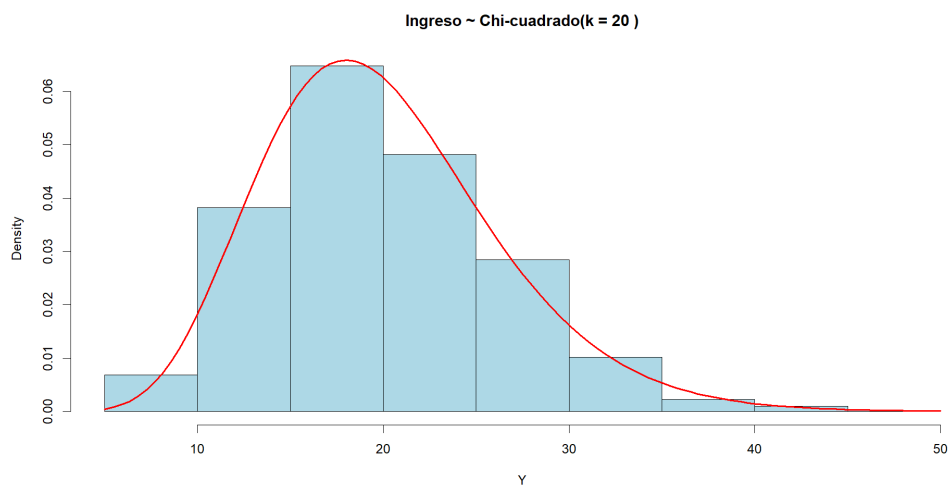


Figura 12: Ingreso simulado vs densidad teórica ($k = 20$).

3.2. Demanda Cobb–Douglas y utilidad indirecta

Bajo precios $p_1, p_2 > 0$, ingreso $Y \geq 0$ y preferencias $a_1, a_2 > 0$ con $a_1 + a_2 = 1$, la demanda óptima es:

$$x_1^* = \frac{a_1 Y}{p_1}, \quad x_2^* = \frac{a_2 Y}{p_2}, \quad V = (x_1^*)^{a_1} (x_2^*)^{a_2}.$$

Utilizando el código mostrado a continuación, creamos la función *demanda_cd* que permite estimar los valores óptimamente demandados dados los precios, los ingresos y las elasticidades.

```
1 Y <- simular_ingreso(5,1000)
2
3 demanda_cd <- function(Y, p1, p2, a1, a2){
4   if (p1 <= 0 || p2 <= 0 || a1 <= 0 || a2 <= 0 || Y < 0){
5     print("Los valores ingresados deben ser positivos o, en el caso del
6       ingreso, no negativo")
7   } else if (a1 + a2 - 1 > 1e-8){
8     print("la suma de los alphas debe dar 1")
9   } else {
10    x_1 <- (a1 * Y) / p1
11    x_2 <- (a2 * Y) / p2
12    V <- (x_1 ^ a1) * (x_2 ^ a2)
13    return(c(x_1,x_2,V))
14  }
15 }
16 demanda_cd(mean(Y), 10, 5, 0.5, 0.5)
```

Para la demanda estudiada en la última línea, donde $Y = \text{mean}(Y)$, $p_1 = 10$, $p_2 = 5$, $a_1 = 0,5$, $a_2 = 0,5$, los valores son:

$$x_1^* = 0,2495419, \quad x_2^* = 0,4990839, \quad V^* = 0,3529056$$

3.3. Monte Carlo y estadísticos descriptivos

Simulamos $n = 10,000$ consumidores con $k = 5$, $p_1 = 10$, $p_2 = 8$, $a_1 = 0,3$, $a_2 = 0,7$; y reportamos medias y cuartiles de x_1^*, x_2^*, V . Grafico estas situaciones

```
1 n = 10000
2 k = 5
3 p1 = 10
4 p2 = 8
5 a1 = 0.3
6 a2 = 0.7
7 Y <- simular_ingreso(k,n)
8 resultados <- matrix(NA_real_,n,3)
9
10 for (j in 1:n){
11   for (i in 1:3){
12     resultados[j,i] = demanda_cd(Y[j],p1,p2,a1,a2)[i]
13   }
14 }
15
16 # Extraigo los resultados como vectores para poder graficar
17 x1 <- resultados[,1]
```

```

18 x2 <- resultados[,2]
19 V <- resultados[,3]
20
21 hist(x1, freq = FALSE, col = "lightblue",
22      main = "x_1* (empírico)", xlab = "x_1*")
23 abline(v = mean(x1), col = "blue", lwd = 2)
24 abline(v = quantile(x1, c(0.25, 0.5, 0.75)), col = "red", lty = 2)
25
26 hist(x2, freq = FALSE, col = "lightblue",
27      main = "x_2* (empírico)", xlab = "x_2*")
28 abline(v = mean(x2), col = "blue", lwd = 2)
29 abline(v = quantile(x2, c(0.25, 0.5, 0.75)), col = "red", lty = 2)
30
31 hist(V, freq = FALSE, col = "lightblue",
32      main = "V* (empírico)", xlab = "V*")
33 abline(v = mean(V), col = "blue", lwd = 2)
34 abline(v = quantile(V, c(0.25, 0.5, 0.75)), col = "red", lty = 2)
35
36 # Reporto los valores.
37 qx1 <- quantile(x1, c(0.25, 0.5, 0.75))
38 qx2 <- quantile(x2, c(0.25, 0.5, 0.75))
39 qV <- quantile(V, c(0.25, 0.5, 0.75))

```

Los valores relevantes reportados son:

$$\overline{x_1^*} = 0,1511403, \quad \overline{x_2^*} = 0,4408259, \quad \overline{V^*} = 0,319743$$

Cuartiles de x_1^* : 0,0804 0,1316 0,2007

Cuartiles de x_2^* : 0,2344 0,3837 0,5854

Cuartiles de V^* : 0,1700 0,2783 0,4246

Toda esta información es resumida en los siguientes gráficos:

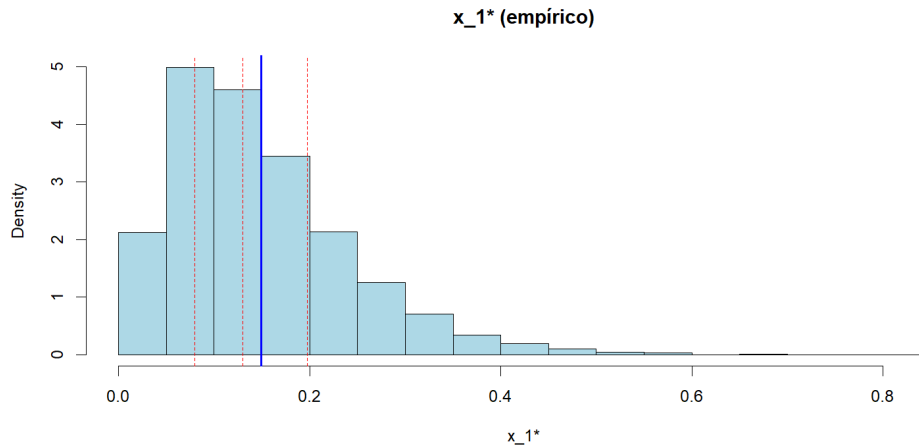


Figura 13: Distribución empírica de x_1^* con media y cuartiles.

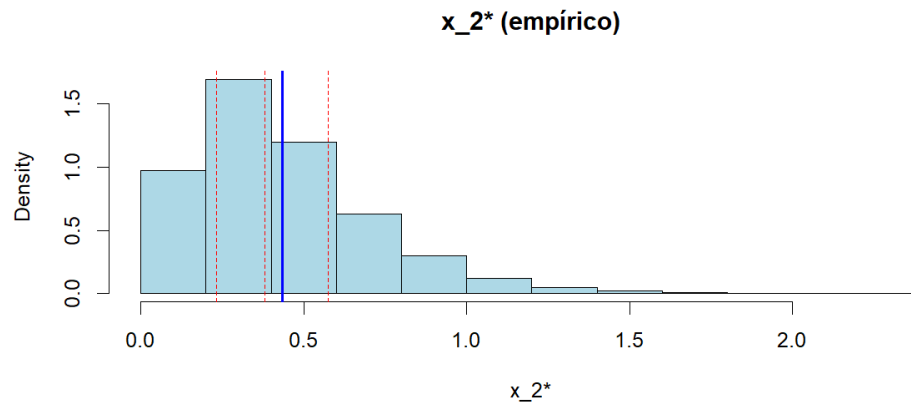


Figura 14: Distribución empírica de x_2^* con media y cuartiles.

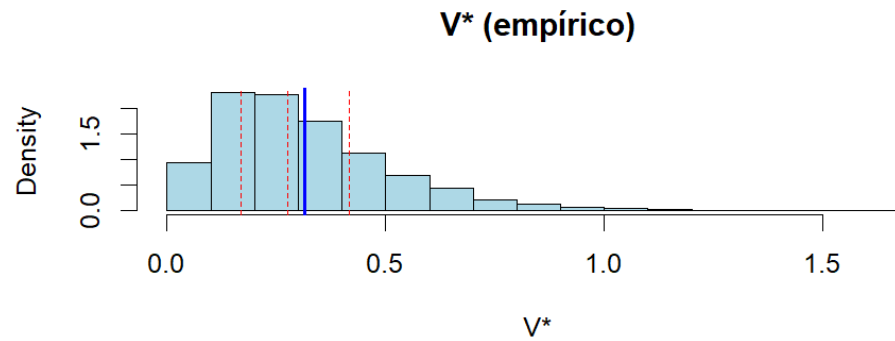


Figura 15: Distribución empírica de V con media y cuartiles.

3.4. Probabilidad de bajo consumo

Estimamos $\Pr(x_j < c)$ como frecuencia muestral, para $j \in \{1, 2\}$.

```

1 prob_bajo_consumo <- function(c, j){
2   if (j == 1){
3     x = x1
4   } else {x = x2}
5   aciertos <- 0
6   for (i in 1:n){
7     if (x[i] < c){
8       aciertos = aciertos + 1
9     }
10  }
11  return(aciertos/n)
12 }
```

3.5. Shock de precio y comparación de distribuciones

Aplicamos un shock $p_1' = 1,2 p_1$ y comparamos medias y cuartiles de x_1^*, x_2^*, V antes y después. El resultado es que x_1^* cae, x_2^* no cambia (no depende de p_1) y la utilidad indirecta V disminuye.

```

1 p1_nuevo = p1*1.2
2 resultados_nuevos <- matrix(NA_real_,n,3)
3
4 for (j in 1:n){
5   for (i in 1:3){
6     resultados_nuevos[j,i] = demanda_cd(Y[j],p1_nuevo,p2,a1,a2)[i]
7   }
8 }
9
10 x1_nuevos <- resultados_nuevos[,1]
11 x2_nuevos <- resultados_nuevos[,2]
12 V_nuevos <- resultados_nuevos[,3]
13
14 qx1_nuevos <- quantile(x1_nuevos, c(0.25, 0.5, 0.75))
15 qx2_nuevos <- quantile(x2_nuevos, c(0.25, 0.5, 0.75))
16 qV_nuevos <- quantile(V_nuevos, c(0.25, 0.5, 0.75))

```

Los efectos de este shock en las variables de interés son resumidos como sigue:

Diferencia de medias (empíricas):

$$\overline{x_1^*} - \overline{x_{1,nuevo}^*} = 0,02519$$

$$\overline{x_2^*} - \overline{x_{2,nuevo}^*} = 0$$

$$\overline{V^*} - \overline{V_{nuevo}^*} = 0,01702$$

Diferencia de cuartiles (empíricos):

$$x_1^* - x_{1,nuevo}^* : \quad 0,0134 \quad 0,0220 \quad 0,0334$$

$$x_2^* - x_{2,nuevo}^* : \quad 0,0000 \quad 0,0000 \quad 0,0000$$

$$V^* - V_{nuevo}^* : \quad 0,0090 \quad 0,0148 \quad 0,0226$$

3.6. Visualización: antes vs. después del shock

Graficamos la superposición de histogramas de x_1^* antes y después del aumento de p_1 .

```

1 hist(x1, freq = FALSE, col = "lightblue",
2     main = "x1*: antes vs. después (shock p1 +20%)", xlab = "x1*")
3 hist(x1_nuevos, freq = FALSE, col = NA, border = "red", add = TRUE, lty =
4     2)
5 abline(v = mean(x1), col = "blue", lwd = 2)
6 abline(v = mean(x1_nuevos), col = "darkred", lwd = 2, lty = 2)
7
8 legend("topright",
9     legend = c("Antes", "Después (borde)"),
10    fill = c("lightblue", "red"),
11    border = c("black", "lightblue"),

```

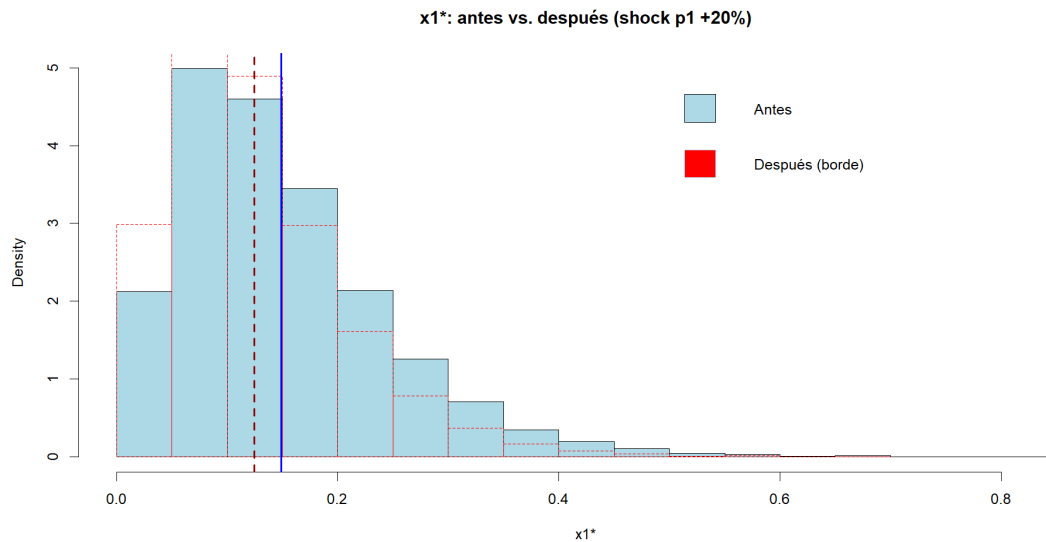


Figura 16: Distribución de x_1^* : antes vs. después del shock en p_1 .

Vemos que la distribución tiene mayor densidad en valores bajos y menor densidad en valores altos. Esto también se visualiza en un menor promedio, tal como es mostrado por las rectas verticales.

La intuición económica es que el precio del bien 1 aumentó, y todo lo demás quedó constante. Por lo tanto, el individuo decidirá consumir menos del bien 1. Dado que el consumo del bien 1 también depende del ingreso, el resultado es que toda la distribución se mueve hacia valores menores de consumo del bien 1.

Lógicamente, ante menor consumo, con todo lo demás constante, la utilidad indirecta V cae. Esto fue mostrado empíricamente en el inciso 5.

3.7. Heterogeneidad en preferencias

Introducimos heterogeneidad $a_1 \sim \text{Beta}(a, b)$ y repetimos el ejercicio base y el shock. Esperamos que mayor a_1 implique mayor x_1^* y menor x_2^* . Elijo $a = 1$ y $b = 2$ arbitrarios, sabiendo que la esperanza es $a/(a+b)$ y la varianza es $ab/((a+b)^2 * (a+b+1))$

```

1 alpha1 <- numeric(n)
2 a = 1
3 b = 2
4 for(i in 1:n){
5   alpha1[i] <- rbeta(n = 1, shape1 = a, shape2 = b)
6 }
7
8 resultados_het <- matrix(NA_real_, n, 3)
9 for (j in 1:n){
10   for (i in 1:3){
11     resultados_het[j, i] = demanda_cd(Y[j], p1, p2, alpha1[j], 1-alpha1[j])[i]
12   }
13 }

```

```

14
15 # Extraigo vectores
16 x1_het <- resultados_het[,1]
17 x2_het <- resultados_het[,2]
18 V_het  <- resultados_het[,3]
19
20 qx1_het <- quantile(x1_het, c(0.25, 0.5, 0.75))
21 qx2_het <- quantile(x2_het, c(0.25, 0.5, 0.75))
22 qV_het  <- quantile(V_het,  c(0.25, 0.5, 0.75))

```

Las diferencias de las distribuciones iniciales con las heterogéneas son:

Diferencia de medias (empíricas):

$$\overline{x_1^*} - \overline{x_{1,\text{het}}^*} = -0,01763$$

$$\overline{x_2^*} - \overline{x_{2,\text{het}}^*} = 0,02203$$

$$\overline{V^*} - \overline{V_{\text{het}}^*} = -0,04322$$

Diferencia de cuartiles (empíricos):

$$x_1^* - x_{1,\text{het}}^* : \quad 0,0353 \quad 0,0184 \quad -0,0345$$

$$x_2^* - x_{2,\text{het}}^* : \quad 0,0503 \quad 0,0432 \quad 0,0180$$

$$V^* - V_{\text{het}}^* : \quad -0,0117 \quad -0,0261 \quad -0,0520$$

Ahora obtengo el cambio de precio con las nuevas preferencias heterogéneas, corriendo las siguientes líneas:

```

1 # Cambio de precio con preferencias heterogéneas
2 resultados_nuevos_het <- matrix(NA_real_,n,3)
3 for (j in 1:n){
4   for (i in 1:3){
5     resultados_nuevos_het[j,i] = demanda_cd(Y[j],p1_nuevo,p2,alpha1[j],1-
6       alpha1[j])[i]
7   }
8 }
9
10 x1_nuevos_het <- resultados_nuevos_het[,1]
11 x2_nuevos_het <- resultados_nuevos_het[,2]
12 V_nuevos_het  <- resultados_nuevos_het[,3]
13
14 qx1_nuevos_het <- quantile(x1_nuevos_het, c(0.25, 0.5, 0.75))
15 qx2_nuevos_het <- quantile(x2_nuevos_het, c(0.25, 0.5, 0.75))
16 qV_nuevos_het  <- quantile(V_nuevos_het,  c(0.25, 0.5, 0.75))

```

Vemos que el resultado de este shock es, como sigue:

Diferencia de medias (empíricas):

$$\overline{x_1^*} - \overline{x_{1,\text{nuevos het}}^*} = 0,01050$$

$$\overline{x_2^*} - \overline{x_{2,\text{nuevos het}}^*} = 0,02203$$

$$\overline{V^*} - \overline{V_{\text{nuevos het}}^*} = -0,02471$$

Diferencia de cuartiles (empíricos):

$$x_1^* - x_{1,\text{nuevos het}}^* : 0,0428 \quad 0,0372 \quad 0,0047$$

$$x_2^* - x_{2,\text{nuevos het}}^* : 0,0503 \quad 0,0432 \quad 0,0180$$

$$V^* - V_{\text{nuevos het}}^* : 0,0004 \quad - 0,0081 \quad - 0,0257$$

Intuitivamente, un mayor α_1 simboliza mayor importancia del bien 1 en la canaste del consumidor. Racionalmente, el consumidor destinará mayor parte de su ingreso al bien 1 a costa de consumo del bien 2, el cual se volvió menos deseado.

Entonces, es esperable que un aumento de α_1 tenga como consecuencia un aumento del consumo de x_1 y una disminución del consumo del bien 2, a igual ingreso y precios.