

Una introducción a la epistemología formal

Diego Tajer

SADAF

Tajer, Diego

Una introducción a la epistemología formal / Diego Tajer. - 1ª ed - Ciudad de Buenos Aires : SADAF, 2025.

204 p. ; 24 x 17 cm.

ISBN 978-987-47781-7-8

1. Epistemología. I. Título.

CDD 120

© 2025, Diego Tajer

© 2024, por esta edición: SADAF

1ª edición: noviembre 2025

SADAF

www.sadaf.org.ar

Diseño de tapa: Iñaki Jankowski | www.jij.com.ar

Índice

INTRODUCCIÓN	V
CAPÍTULO 1: PROBABILIDADES	1
Parte A. Axiomas y principios	1
Parte B: Probabilidad condicional.....	6
Parte C: Regla de Bayes.....	10
Parte D: Probabilidades y mundos posibles	16
Parte E: Probabilidad y Lógica	21
Parte F: Filosofía de la probabilidad	23
Parte G: Coherencia probabilística y racionalidad	25
Parte H. Nociones de estadística	36
Soluciones para el capítulo 1.....	44
CAPÍTULO 2: TEORÍA DE LA DECISIÓN.....	51
Parte A: Matrices, actos y resultados	51
Parte B: Decisiones bajo ignorancia	53
Parte C: Escalas de utilidad.....	59
Parte D: Maximización de utilidad esperada.....	62
Parte E: El Teorema von Neumann-Morgenstern.....	68
Parte F: Utilidad conductual y utilidad sustantiva	72
Parte G: La Paradoja de Allais.....	75
Parte H: La Paradoja de Ellsberg	78
Parte I: Psicología de la decisión.....	82
Parte J: Buchak y el riesgo como factor	87
Parte K: Experiencia transformadora.....	90
Soluciones para el capítulo 2.....	92
CAPÍTULO 3: TEORÍA DE JUEGOS.....	97
Parte A: Juegos estratégicos	97
Parte B: Juegos dinámicos.....	111
Parte C: Cooperación	119
Parte D: Paradojas y experimentos	128
Parte E: Juegos de información incompleta	131
*Parte F: Estrategias mixtas y probabilidades	136
Soluciones para el capítulo 3.....	140

CAPÍTULO 4: ELECCIÓN SOCIAL.....	145
Parte A: Antecedentes y teoría del voto	145
*Parte B: Blindar el voto por mayoría	148
Parte C: Arrow y funciones de bienestar social.....	150
*Parte D: Prueba del Teorema de Arrow	153
Parte E: Lecturas del Teorema de Arrow	159
Parte F: Soluciones al Teorema de Arrow	161
Parte G: Paradoja Del Liberal Paretiano	166
Parte H: Resultados positivos sobre el voto	169
Parte I: Agregación de juicios.....	173
Soluciones para el capítulo 4.....	177
EPÍLOGO.....	183
BIBLIOGRAFÍA	185

INTRODUCCIÓN

La filosofía analítica se ha distinguido históricamente por su uso de herramientas lógicas. Por esta razón, encontramos cursos de lógica en cualquier carrera de Filosofía a lo largo del mundo. Incluso la epistemología analítica tradicional, que se ha preocupado fundamentalmente por preguntas conceptuales como la naturaleza del *conocimiento* y de la *justificación*, suele utilizar la Lógica como herramienta.

En contraste con la epistemología analítica tradicional, la *epistemología formal* nos propone responder preguntas de epistemología utilizando herramientas matemáticas más variadas, tales como probabilidades, teoría de la decisión, teoría de juegos, o teoría de la elección social. Muchas de estas áreas fueron discutidas por filósofos provenientes de la filosofía de las ciencias, la filosofía política o la ética. Por ejemplo, uno podría preguntarse por el valor de la estadística para la explicación científica; o por la importancia de la teoría de juegos para entender el origen del concepto de justicia; o por el valor del Teorema de Arrow para establecer límites a los sistemas políticos conocidos.

La epistemología formal, entonces, es un término muy general que engloba distintas áreas de conocimiento que han sido aplicadas en distintos debates, pero que (a diferencia de la Lógica) no suelen formar parte de la formación básica de un filósofo. El propósito de este libro es justamente ayudar a remediar esa situación, dentro del área hispanohablante. Este libro nos propone una introducción general a las distintas áreas centrales de la epistemología formal para filósofos, de forma accesible y a la vez sólida.

En primer lugar, este libro no presupone un nivel elevado de matemáticas. Esto no significa que su lectura sea fácil: algunas demostraciones son largas y complejas. El único conocimiento requerido es un buen manejo de operaciones algebraicas (saber multiplicar fracciones, dividir, trabajar ecuaciones e inecuaciones, etc.) y haber aprobado un curso básico de lógica.

Las secciones o los ejercicios que requieren más destreza matemática o lógica están marcados con un *.

En segundo lugar, este libro no está orientado necesariamente a estudiantes especializados, sino a estudiantes con suficiente interés en el tema. De este modo, puede usarse como libro de referencia para personas interesadas más indirectamente en estos temas que quieran utilizar ideas de teoría de juegos, teoría de la decisión o teoría de la elección social en sus áreas de trabajo, pero no sepan por dónde empezar.

Por último, este no es un libro “de divulgación”, que expone ideas complejas en un vocabulario sencillo pero sin entrar en grandes detalles. Por el contrario, este libro es un *manual*: expone ideas relativamente sencillas con suficiente detalle para entender los conceptos y métodos, y propone ejercicios a los lectores. Solucionar los ejercicios en cada sección dará a los lectores un conocimiento más sólido sobre estos temas.

El primer capítulo del libro trata sobre el concepto de *probabilidad*. En la primera parte se explican los principios básicos de las probabilidades y la regla de Bayes para la condicionalización. Luego entramos en detalles más filosóficos, como las interpretaciones del concepto de probabilidad, la relación entre probabilidad y validez lógica y los distintos resultados para establecer que cumplir las leyes probabilísticas es propiamente racional. Al final, se introducen algunos conceptos generales sobre estadística, en particular para su aplicación en el testeo de hipótesis científicas.

El segundo capítulo del libro trata sobre la *teoría de la decisión*. Empezamos por introducir el esquema conceptual de la teoría de la decisión, y la división entre acciones y estados del mundo. Luego discutimos brevemente la teoría de la decisión “bajo ignorancia”, que nos permite tomar decisiones racionales incluso cuando ignoramos las probabilidades de los estados del mundo. Más adelante se introduce la teoría de la decisión “bajo riesgo”, y el método de *maximización de la utilidad esperada*. La segunda parte de este capítulo se mete en las discusiones conceptuales, tales como la naturaleza de la *utilidad* según la teoría de la decisión, y las distintas objeciones o paradojas que se han

presentado contra el principio de maximización de utilidad. Hacia el final, introducimos algunas teorías contemporáneas (posteriores al 2010), tales como las experiencias transformadoras de L. A. Paul y la decisión basada en riesgos de L. Buchak.

El tercer capítulo trata sobre la *teoría de juegos*. Aquí, empezamos por distinguir entre juegos estratégicos y dinámicos. Primero se explican algunos métodos de solución para los juegos estratégicos, tales como el método de borrado iterado de estrategias dominadas, y la búsqueda de equilibrios de Nash. Luego se introduce el método de inducción hacia atrás para los juegos dinámicos, usando árboles de decisión. Más adelante, discutimos el concepto de Cooperación a partir del Dilema del Prisionero, y exponemos su importancia en la filosofía contemporánea. Hacia el final, exponemos brevemente algunos juegos de “información incompleta”, y mostramos cómo obtener equilibrios de Nash para estrategias *mixtas*, es decir, estrategias que utilizan algún tipo de procedimiento probabilístico.

El cuarto y último capítulo trata sobre la *teoría de la elección social*. Primero se introducen algunos métodos y conceptos en la clásica teoría del voto, y se explica la paradoja de los ciclos de Condorcet. Más adelante, se explican los conceptos fundamentales de la teoría contemporánea de la Elección Social, partiendo del concepto de “función de bienestar social”. En ese contexto, se introduce el teorema fundamental de esta área de investigación: el Teorema de Arrow. También agregamos una demostración sencilla del teorema, que los lectores más atentos pueden seguir. Luego discutimos la relevancia filosófica del Teorema de Arrow y las distintas lecturas que se han hecho de este resultado. Más adelante, introducimos resultados más “positivos” sobre el voto democrático, como el Teorema del Jurado de Condorcet. Por último, presentamos la *teoría de agregación de juicios*, un área muy similar a la Elección Social, pero más centrada en las relaciones lógicas entre las proposiciones.

El propósito de este libro, entonces, es introducir a los lectores en los métodos y los conceptos fundamentales de la epistemología formal. Incluso para aquellos que no se dediquen específicamente a estos temas, el libro les servirá para aprender nuevas herramientas técnicas y conceptuales que podrán aplicar

en sus futuras investigaciones. Aunque la matemática y la filosofía son cosas bien distintas, la filosofía puede sacar mucho provecho de los métodos formales para su propósito de clarificar y desarrollar conceptos.

Para terminar, quiero agradecer a Eleonora Cresto, con quien trabajo desde hace más de diez años, y que me acompañó en el aprendizaje de todos los temas que trato en el libro. Le agradezco especialmente por haber leído y comentado el manuscrito del libro. También agradezco a Javier Arróspide por sus comentarios sobre el capítulo 1, a mi editor Mariano Blatt, y a Juan Redmond por su trabajo sobre una versión anterior del texto.

Diego Tajer
Instituto de Investigaciones Filosóficas
CONICET, Argentina
diegotajer@gmail.com

CAPÍTULO 1: PROBABILIDADES

Parte A. Axiomas y principios

Muchas de las teorías que explicaremos en este libro presuponen, de algún u otro modo, la teoría clásica sobre la *probabilidad*. La probabilidad aparece en nuestro razonamiento cotidiano, a veces de forma explícita. Por ejemplo, mi celular indica que la probabilidad de lluvia hoy en Buenos Aires es del 65 %. Sé, al menos, que la lluvia es bastante probable.

¿Pero en qué consiste la probabilidad? En términos generales, la probabilidad es una *función* que asigna un número a los eventos posibles. Este número será mayor en tanto la probabilidad de que el evento ocurra sea más alta. En el contexto de las teorías de la racionalidad, entendemos usualmente a la probabilidad de forma subjetiva: lo que importa realmente no es la probabilidad del hecho en sí, sino la probabilidad que le asigna un agente, y qué es lo que ese agente hace a partir de ello. Sin embargo, en otros contextos, podemos entender las probabilidades como un hecho objetivo.

Empecemos por lo básico. Llamamos *espacio muestral* al conjunto de posibles resultados de un experimento. En fenómenos aleatorios (como tirar un dado o una moneda), estos eventos pueden ocurrir con la misma probabilidad. Por ejemplo, si tiramos un dado, el espacio muestral será $\{1, 2, 3, 4, 5, 6\}$. La probabilidad de un evento E en estos contextos aleatorios podría calcularse de este modo, conocido como *Regla de Laplace*:

$$P(E) = \frac{\text{Número de casos donde sucede E}}{\text{Número de casos posibles}}$$

Por ejemplo, la probabilidad de sacar un número par al tirar un dado será $\frac{1}{2}$, porque hay 6 casos posibles, y solo 3 casos de números pares (2, 4 y 6).

Los espacios muestrales podrían ser mucho más grandes. Por ejemplo, si tiro el dado dos veces, ya no tendré un espacio muestral de 6 elementos, sino uno de 36: $\{(1,1), (1,2), (1,3), \dots, (2,1), (2,2), (2,3), \dots\}$. Allí se indica lo que sale en la primera tirada y también lo que sale en la segunda. De este modo, por ejemplo, la probabilidad de sacar primero un número y luego su doble es $3/36 = 1/12$, porque entre los 36 posibles los resultados los únicos donde eso sucede son $(1,2)$, $(2,4)$ y $(3,6)$.

Las probabilidades cumplen con algunos **axiomas**, es decir, principios fundamentales. Estos principios suelen llamarse Axiomas de Kolmogórov, porque fueron establecidos por el matemático ruso Andréi Kolmogórov en 1933 (aunque la versión original es ligeramente distinta a la presentada aquí).¹

Vamos a suponer que A es una proposición y que $P(x)$ es la función de probabilidad. El primer axioma dice que la probabilidad de una proposición es un número real entre 0 y 1. Este axioma suele llamarse *Normalidad*.²

$$\text{Axioma 1 (Normalidad): } 0 \leq P(A) \leq 1$$

Es decir, lo menos probable tendrá probabilidad 0, y lo más probable tendrá probabilidad 1. Si entendemos la probabilidad como un cociente entre casos donde sucede el evento y todos los casos posibles, es natural pensar que una probabilidad estará entre 0 (no se da en ningún caso) y 1 (se da en todos los casos).

Los otros axiomas involucran, de algún u otro modo, cuestiones lógicas. En primer lugar, la probabilidad de las tautologías debe ser 1. Este axioma suele llamarse *Certeza*.

¹ En la presentación original de Kolmogórov, la función de probabilidad se aplica a eventos, no a proposiciones. Aquí usaremos ambos enfoques de forma equivalente. Después de todo, suele asumirse que una proposición es un conjunto de mundos posibles (y un evento también, como veremos en la Parte D).

² Los nombres para los axiomas siguen las convenciones de Hacking (2001, p. 61).

Axioma 2 (Certeza): Si A es una tautología³, entonces $P(A) = 1$.

Este axioma es muy importante para las aplicaciones de la teoría de la probabilidad. Desde el punto de vista objetivista, es completamente obvio que las verdades lógicas ocurrirán en todos los casos posibles. Desde el punto de vista más subjetivista, este axioma es más discutible, pues asume que los agentes asignan a las tautologías una probabilidad máxima. En otras palabras, se asume que los agentes saben todo sobre lógica (en ocasiones esto se llama “omnisciencia lógica”).

Es importante observar también que el axioma no es bicondicional. Es decir, si bien las tautologías tendrán siempre probabilidad 1, podría haber otras proposiciones con probabilidad 1. Esto incluye a aquellas proposiciones que son absolutamente seguras en un contexto determinado.

El tercer axioma nos dice que, si dos proposiciones son incompatibles entre sí, entonces la probabilidad de su disyunción será simplemente la suma de las probabilidades. Llamamos a este axioma *Aditividad*.

Axioma 3 (Aditividad): Si A y B son incompatibles entre sí, es decir $\{A, B\} \vdash \perp$, entonces $P(A \vee B) = P(A) + P(B)$.

Tal como sucede con el axioma anterior, aquí la incompatibilidad no será siempre lógica; basta que dos eventos no puedan ocurrir al mismo tiempo para que sean incompatibles.

Un caso obvio que puede inferirse a partir de estos axiomas es lo que sucede con las proposiciones y sus negaciones. Si la probabilidad de A es p , la probabilidad de $\neg A$ será $(1 - p)$. Es decir, si la probabilidad de que llueva es 0.65, la probabilidad de que no llueva es 0.35. Podemos ver ahora la prueba de este teorema:

³ Este libro presupone ciertos conocimientos de Lógica, en particular de lógica proposicional clásica. A lo largo del libro, esencialmente utilizaré tres conectivos: $\neg A$ es la *negación* de A ; por otro lado, $A \vee B$ es la *disyunción* entre A y B ; finalmente, $A \& B$ es la *conjunción* entre A y B . Una *tautología* es una oración necesariamente verdadera en virtud de su forma lógica, como $p \vee \neg p$.

Teorema (Negación): $P(\neg A) = 1 - P(A)$.

Prueba. $P(A \vee \neg A) = 1$, por Certeza.

$P(A \vee \neg A) = P(A) + P(\neg A)$, por Aditividad.

Entonces $P(A) + P(\neg A) = 1$.

Por lo tanto, $P(\neg A) = 1 - P(A)$. QED

De este teorema también se infiere, casi automáticamente, que la probabilidad de las contradicciones siempre será 0 (véase Ejercicios).

También podemos probar un principio sobre la equivalencia lógica, que usaremos repetidamente más adelante⁴:

Teorema (Equivalencia):

Si A y B son lógicamente equivalentes, $P(A) = P(B)$.

Prueba. Supongamos que A y B son lógicamente equivalentes.

Dado que A y $\neg A$ son incompatibles, y que A y B son lógicamente equivalentes, entonces A y $\neg B$ son incompatibles (por sustitución de equivalentes). Entonces Aditividad implica que $P(A \vee \neg B) = P(A) + P(\neg B)$.

Por la equivalencia entre A y B, $A \vee \neg B$ es una tautología, entonces $P(A \vee \neg B) = 1$, por Certeza.

Entonces $P(A) + P(\neg B) = 1$.

También sabemos por el Teorema de Negación que $P(B) + P(\neg B) = 1$. Entonces $P(B) + P(\neg B) = P(A) + P(\neg B)$, y simplificando idénticos obtenemos $P(B) = P(A)$. QED

Un importante teorema que usaremos luego (especialmente en la parte C de este capítulo) nos permite inferir la probabilidad de A a partir de la probabilidad de $A \& B$ y de $A \& \neg B$.

Teorema (Probabilidad Total):

$P(A) = P(A \& B) + P(A \& \neg B)$

⁴ Las pruebas de los siguientes teoremas son relativamente complejas. El lector sin mucha habilidad en lógica puede saltarlas y quedarse con el enunciado de los teoremas.

Prueba. En lógica clásica, $A \equiv ((A \& B) \vee (A \& \neg B))$, por la ley de Tercero Excluido. Entonces $P(A) = P((A \& B) \vee (A \& \neg B))$, por Equivalencia. Dado que $A \& B$ y $A \& \neg B$ son incompatibles, $P((A \& B) \vee (A \& \neg B)) = P(A \& B) + P(A \& \neg B)$, por Aditividad. Entonces $P(A) = P(A \& B) + P(A \& \neg B)$. QED

A partir de estos principios, podemos probar dos teoremas importantes. El primero relaciona Validez y Lógica (un tema que profundizaremos en la parte E):

Teorema (Validez):

Si A implica lógicamente B, entonces $P(A) \leq P(B)$.

Prueba. Si A implica lógicamente B, entonces $A \equiv A \& B$.

$P(B) = P(B \& A) + P(B \& \neg A)$ [Probabilidad Total]

$= P(A) + P(B \& \neg A)$ [Equivalencia]

Dado que $P(B \& \neg A)$ es mayor o igual a 0 (por Normalidad), se infiere que $P(A) \leq P(B)$. QED

En segundo lugar, podemos inferir también un importante teorema sobre la probabilidad de las disyunciones:

Teorema (Probabilidad de Disyunciones):

$P(A \vee B) = P(A) + P(B) - P(A \& B)$

Prueba.

$P(A) + P(B) =$

$= P(A \& B) + P(A \& \neg B) + P(A \& B) + P(\neg A \& B)$ Pr. Total

$= P(A \& B) + P[(A \& \neg B) \vee (\neg A \& B)]$ Aditiv.⁵

$= P(A \& B) + P(A \vee B)$ Equiv.

Por ende, $P(A) + P(B) = P(A \& B) + P(A \vee B)$.

Si cambiamos los términos, obtenemos:

$P(A \vee B) = P(A) + P(B) - P(A \& B)$. QED

La probabilidad de las disyunciones es una típica fuente de error. Por ejemplo, uno pensaría que, si la probabilidad de sacar un 3

⁵ Aquí usamos una Aditividad de tres proposiciones, que se sigue del axioma de Aditividad para dos proposiciones (omitimos la prueba).

en un dado es $1/6$, la probabilidad de sacar un 3 si uno tira el dado dos veces es el doble, es decir, $2/6$. Sin embargo, la probabilidad es un poco menos que eso. Luego veremos cómo calcularlo.

Ejercicios

1. Tiro dos dados.

- ¿Cuál es la probabilidad de sacar el mismo número en ambos dados?
- ¿Cuál es la probabilidad de sacar 4 en un dado y 3 en el otro?
- ¿Qué es más probable, que ambos dados sumen 7 o que sumen 6? ¿O es igualmente probable?

2. Bingo con letras: Un “bingo” nos da una letra entre la A y la E, y un número del 1 al 9, por ejemplo, B8. Antes de sacar la bolilla, el jugador trata de adivinar si va a salir *par/impar* y *vocal/consonante*. Por ejemplo, el jugador podría apostar a que saldrá una consonante con un número par. Si adivinas que sea vocal/consonante y que sea par/impar, ganas este juego.

- ¿Cuál es la probabilidad de ganar jugando *vocal impar*?
- ¿Cuál es la probabilidad de ganar jugando *consonante par*?
- ¿Cuál es la probabilidad de ganar en *alguno* de los dos anteriores?
- ¿Qué te conviene jugar (si puedo elegir solo una apuesta)? ¿Y qué probabilidad hay de ganar jugando eso?

3. Pruebe (usando los teoremas o axiomas vistos en este capítulo) que, si A es una contradicción, entonces $P(A) = 0$.

Parte B: Probabilidad condicional

Otro concepto usual en probabilidad es el de *probabilidad condicional*. A veces no queremos saber cuál es la probabilidad de

determinado evento, sino la probabilidad de un evento asumiendo que otro evento también ocurre.⁶ Por ejemplo, podríamos querer determinar la probabilidad de que la temperatura en Buenos Aires exceda los 30 grados. Pero también podría interesarnos la probabilidad de que la ciudad supere esa temperatura asumiendo que es verano. Obviamente, asumiendo que es verano, la probabilidad será más alta. La probabilidad de que ocurra un evento A, asumiendo que ocurre un evento B, la escribimos $P(A|B)$. En este caso, por ejemplo, $P(\text{Más de 30 grados} | \text{Verano}) > P(\text{Más de 30 grados})$.

Una formulación sencilla de la probabilidad condicional nos dice lo siguiente, siempre tomando en cuenta que $P(B) \neq 0$:

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

Por ejemplo, en el mazo de cartas españolas, hay 4 “palos” (espada, oro, bastos y copas) y 10 números con cada uno de esos palos. Podríamos preguntarnos cuál es la probabilidad de sacar el 1 de bastos. Claramente es $1/40$. Pero podríamos también preguntarnos la probabilidad de sacar el 1 de bastos (1B) asumiendo que sacamos una carta de bastos (B). En ese caso, intuitivamente sabemos que es $1/10$. Pero también podríamos calcularlo así:

$$\begin{aligned} P(1B | B) &= P(1B \& B) / P(B) \\ &= P(1B) / P(B) && \text{Equivalencia}^7 \\ &= 1/40 / 1/4 \\ &= 1/40 \times 4 = 1/10 \end{aligned}$$

De este modo vemos que el cálculo de la probabilidad condicional coincide con nuestras intuiciones al respecto.

⁶ En sentido estricto, las probabilidades descritas aquí se aplican a proposiciones, no a eventos. Por simplicidad, uso “evento” o “proposición” de forma equivalente a lo largo del libro.

⁷ Por las propiedades de la conjunción, $(1\&B)\&B$ es equivalente a 1B.

Es importante observar que, si movemos los términos en la definición de probabilidad condicional, nos quedará lo siguiente:

$$P(A \& B) = P(B) \times P(A|B)$$

Dado que, por equivalencia lógica, $P(A \& B) = P(B \& A)$, podemos también inferir que:

$$P(A \& B) = P(B \& A) = P(A) \times P(B|A).$$

Esta es la forma más general de calcular la probabilidad de las conjunciones.

Independencia

Usando el concepto de probabilidad condicional, podemos definir el concepto de *independencia probabilística*. Intuitivamente, decimos que dos eventos son independientes cuando no se influyen entre sí, es decir, cuando el hecho de que uno sucede no cambia la probabilidad del otro. En términos formales, necesitamos el concepto de probabilidad condicional:

(Independencia) Dos eventos A y B son *independientes* si y sólo si $P(A) = P(A|B)$.

De aquí se sigue el teorema de Probabilidad de Conjunciones:

Teorema (Probabilidad de Conjunciones): Si A y B son dos eventos independientes, entonces $P(A \& B) = P(A) \times P(B)$.

Prueba: Sabemos que $P(A \& B) = P(A) \times P(B|A)$, por la definición de Probabilidad Condicional. Pero dado que A y B son independientes, $P(B|A) = P(B)$. Entonces $P(A \& B) = P(A) \times P(B)$. QED

En otras palabras, la probabilidad de eventos independientes se multiplica. La probabilidad de que una moneda salga “cara” dos veces seguidas será $0.5 \times 0.5 = 0.25$. Esta regla nos permite calcular probabilidades un poco más complejas.

Por ejemplo, supongamos que tiro dos veces un dado. Podría querer calcular la probabilidad de sacar primero un número par, y luego un número mayor que 2. La probabilidad será $P(\text{Par en primera tirada} \& \{3,4,5,6\} \text{ en segunda tirada}) = P(\text{Par en primera tirada}) \times P(\{3,4,5,6\} \text{ en segunda tirada}) = \frac{1}{2} \times \frac{2}{3} = \frac{1}{3}$.

Con el teorema sobre la probabilidad de conjunciones, podemos calcular la probabilidad disyuntiva mencionada en la sección anterior. Decíamos que la probabilidad de sacar un 3 en una tirada de dado es $1/6$. Pero la probabilidad de sacar un 3 en alguna de *dos* tiradas no es el doble. Es un poco menos que el doble. Podemos calcular la probabilidad exacta, usando el principio de probabilidad de la disyunción:

$$\begin{aligned} P(3 \text{ en primera tirada} \vee 3 \text{ en segunda tirada}) &= \\ P(3 \text{ en primera tirada}) + P(3 \text{ en segunda tirada}) - P(3 \text{ en ambas tiradas}) &= \\ = 1/6 + 1/6 - (1/6 \times 1/6) = 6/36 + 6/36 - 1/36 = 11/36 \end{aligned}$$

Hay instancias mucho más intuitivas de esta idea. Por ejemplo, la probabilidad de sacar al menos una “cara” tirando dos veces una moneda va a ser 0.75 (tenemos tres probabilidades sobre cuatro). Esto puede inferirse con la misma fórmula, porque $P(\text{Cara en 1} \vee \text{Cara en 2}) = P(\text{Cara en 1}) + P(\text{Cara en 2}) - P(\text{Dos veces cara}) = 0.5 + 0.5 - 0.25 = 0.75$.

La independencia probabilística no siempre es percibida como tal. Por ejemplo, si vamos a un casino veremos personas usando el siguiente razonamiento: “En las últimas tres rondas de la ruleta salió Rojo, entonces en la próxima ya debería salir Negro”. Llamamos a este tipo de falacia “Falacia del Apostador”. Este es un tipo de razonamiento muy usual, pero falaz, donde percibimos probabilidades independientes como si no lo fueran. Para la teoría de la probabilidad, el hecho de que haya salido tres veces Rojo no indica que la próxima vaya a salir Negro, o que sea más probable que eso suceda. Son hechos probabilísticamente independientes.

Ejercicios

1. Tengo un mazo de cartas españolas. Es decir, hay 40 cartas y 4 “palos” (espada, copa, oro y bastos), y cada uno de ellos tiene diez números (1, 2, 3, 4, 5, 6, 7, 10, 11 y 12).

a. ¿Cuál es la probabilidad de sacar un 7 dos veces seguidas, con reemplazo? (es decir, sacar un 7, volver a poner la carta, mezclar y volver a sacar un 7)

b. ¿Cuál es la probabilidad de sacar un 7 dos veces seguidas, sin reemplazo?

c. ¿Cuál es la probabilidad de sacar dos cartas (con reemplazo) de las cuales ninguna es de espada? ¿Y sin reemplazo?

d. *¿Cuál es la probabilidad de sacar un 7 en alguna de dos tiradas (con reemplazo)?

2. 10% de los habitantes de Buenos Aires viven en el barrio de Palermo. De los habitantes de Palermo, 40% da positivo a una prueba de toxicidad en la piel (por contaminantes de desechos químicos). De los habitantes de Buenos Aires que no están en Palermo, solo 10% da positivo la prueba de toxicidad en la piel.

a. ¿Cuál es la probabilidad de que una persona (de Buenos Aires), elegida al azar, viva en Palermo y dé positivo a la prueba?

b. ¿Cuál es la probabilidad de que una persona (de Buenos Aires), elegida al azar, de positivo la prueba? [Usar probabilidad total]

c. ¿Cuál es la probabilidad de que una persona (de Buenos Aires), elegida al azar, que da positivo la prueba, viva en Palermo? [Usar probabilidad condicional]

d. *Viene al hospital el hermano de esa persona, que *vive en el mismo barrio*. Le hacen la prueba y también da positivo. ¿Cuál es ahora la probabilidad de que ambos sean de Palermo?

Parte C: Regla de Bayes

Usando la definición probabilística de la conjunción, podemos reformular la probabilidad condicional, y expresar lo que usualmente se conoce como *regla de Bayes*:

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

La regla de Bayes, llamada así por su descubridor, Thomas Bayes (1702-1761), se usa en disciplinas científicas para evaluar la probabilidad de una hipótesis a partir de cierta evidencia. En general, llamamos “probabilidad previa” o “prior” a $P(A)$, y “likelihood” a $P(B|A)$.⁸

La idea es la siguiente. Supongamos que a la mañana de un día de invierno me desperté con fiebre. Entonces, empiezo a pensar qué podría haberme causado la fiebre. Ciertamente, si tuviera una enfermedad grave, seguramente me daría fiebre; por ejemplo, la *likelihood* de tener fiebre si tuviera un cáncer avanzado sería muy alta. Pero no tengo razones para pensar que tengo una enfermedad grave. Es decir, tomando en cuenta mi edad y mi estado general de salud, la probabilidad previa de tener cáncer avanzado es muy baja. Pienso en otra hipótesis: me intoxicqué con comida. De nuevo, la *likelihood* de tener fiebre al intoxicarme con comida es alta. Sin embargo, recuerdo que anoche comí una ensalada de tomate, lechuga y pollo bien cocinado. Entonces, es muy improbable haberme intoxicado con esa cena. Es decir, la probabilidad previa también es bajísima. Se me ocurre una última hipótesis: quizás tengo una gripe. La probabilidad previa de tener gripe en invierno es más alta, y si tuviera gripe también tendría fiebre (es decir, la *likelihood* de tener fiebre dado que tengo gripe también es alta). Por eso, la probabilidad de tener gripe en invierno, asumiendo que tengo fiebre, es bastante alta.

Usando la regla de Bayes, podríamos representar la situación de este modo. Tengo fiebre, y quiero calcular la probabilidad de tener gripe:

$$P(\text{Gripe} | \text{Fiebre}) = \frac{P(\text{Fiebre} | \text{Gripe}) \times P(\text{Gripe})}{P(\text{Fiebre})}$$

⁸ No hay una buena traducción de *likelihood* al español, por lo cual dejamos esa palabra en inglés.

El problema ahora es que no tengo cómo establecer la probabilidad de la evidencia por sí sola (es decir, en este caso, la probabilidad de tener fiebre). Ahora supongamos (para simplificar) que, al tener fiebre, solo pienso en dos hipótesis: o gripe, o cáncer. En este contexto, son hipótesis exclusivas y excluyentes.

Con esto, puedo volver a los principios probabilísticos antes mencionados (especialmente, la Probabilidad Total), y recordar que la probabilidad de tener fiebre es la suma entre la probabilidad de fiebre&cáncer y la probabilidad de fiebre&gripe (asumiendo que tener gripe y tener cáncer son hipótesis exhaustivas y excluyentes). Y el resto lo calculo usando la regla de probabilidad de conjunciones.

Es decir:

$$\begin{aligned} P(\text{Fiebre}) &= P(\text{Fiebre} \ \& \ \text{Cáncer}) + P(\text{Fiebre} \ \& \ \text{Gripe}) \\ &= P(\text{Cáncer}) \times P(\text{Fiebre} \mid \text{Cáncer}) + P(\text{Gripe}) \times P(\text{Fiebre} \mid \text{Gripe}) \end{aligned}$$

En otras palabras: la probabilidad de la evidencia (tener fiebre) es la suma de los priors por likelihoods de cada hipótesis, siempre y cuando estas hipótesis sean exhaustivas y excluyentes.

Si asumimos que H_1 y H_2 son hipótesis exhaustivas y excluyentes, llegamos a otra versión muy conocida de la regla de Bayes, que establece la probabilidad de una hipótesis H_1 sobre la base de la evidencia e :

$$P(H_1|e) = \frac{P(e|H_1) \times P(H_1)}{P(e|H_1) \times P(H_1) + P(e|H_2) \times P(H_2)}$$

En esta versión, el denominador tiene la probabilidad previa y el likelihood de todas las hipótesis alternativas, asumiendo que son exhaustivas y excluyentes.

En el ejemplo en cuestión, la ecuación quedaría de este modo:

$$\begin{aligned} P(\text{Gripe} \mid \text{Fiebre}) = & \\ & \frac{P(\text{Fiebre} \mid \text{Gripe}) \times P(\text{Gripe})}{P(\text{Fiebre} \mid \text{Gripe}) \times P(\text{Gripe}) + P(\text{Fiebre} \mid \text{Cáncer}) \times P(\text{Cáncer})} \end{aligned}$$

Asumiendo que los likelihoods son parecidos, (Gripe | Fiebre) será mucho mayor a P(Cáncer | Fiebre) porque la probabilidad previa de tener gripe en invierno, P(Gripe), es mucho mayor a la de tener cáncer, P(Cáncer). Justamente, las personas que se atribuyen la peor enfermedad ante el mínimo síntoma, los hipocondríacos, cometen esta falacia: ignoran la probabilidad previa de las hipótesis. Esta falacia se conoce como *falacia de tasa base* (*base rate fallacy*).

La regla de Bayes puede generalizarse para más de dos hipótesis, asumiendo también que son exhaustivas y excluyentes. En ese caso, extendemos esta definición del siguiente modo:

$$P(H_i|e) = \frac{P(e|H_i) \times P(H_i)}{\sum_{i=1}^n P(e|H_i) \times P(H_i)}$$

Es decir, el denominador será la suma de probabilidad previa por likelihood de *todas* las hipótesis.

Ejemplo. En mayo de 2020 en Buenos Aires, me despierto con una fuerte tos. Interpreto que puede ser COVID o una gripe estacional; seguro tengo alguna de las dos cosas. En este momento, la probabilidad de tener COVID es 0.001, mientras que la probabilidad de tener una gripe estacional es 0.002. Por otro lado, la gripe estacional provoca tos en 50% de los casos, mientras que el COVID provoca tos en 80% de los casos. Entonces, ¿cuál es la probabilidad de tener COVID?

$$\begin{aligned} P(\text{Covid} | \text{Tos}) &= \frac{P(\text{Tos} | \text{Covid}) \times P(\text{Covid})}{P(\text{Tos} | \text{Covid}) \times P(\text{Covid}) + P(\text{Tos} | \text{Gripe}) \times P(\text{Gripe})} \\ &= \frac{0.8 \times 0.001}{0.8 \times 0.001 + 0.5 \times 0.002} = \frac{0.0008}{(0.0008 + 0.001)} = \frac{0.0008}{0.0018} = 0.44 \end{aligned}$$

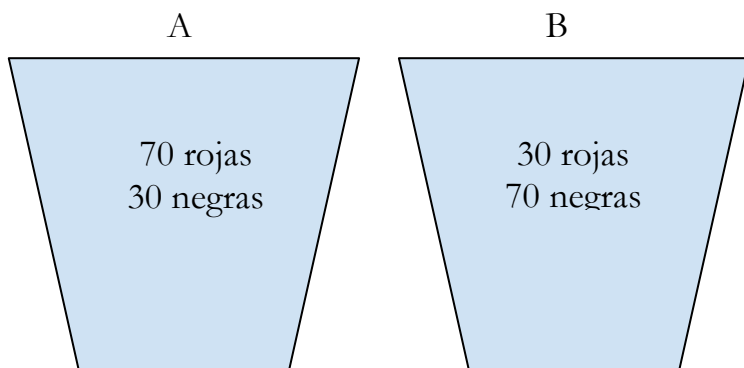
Es decir, la probabilidad de tener COVID será 0.44. Como son probabilidades excluyentes y exhaustivas, la probabilidad de tener gripe estacional será 0.56.

Usando la regla de Bayes, podemos aprender a actualizar nuestras creencias a partir de nueva evidencia. Lo interesante también es que la probabilidad actualizada pasa a ser probabilidad “vieja” cuando aparece nueva evidencia. La actualización puede plantearse como un proceso diacrónico, donde mientras va apareciendo evidencia, vamos actualizando nuestras creencias.

Un problema típico del bayesianismo como teoría del conocimiento es la imposibilidad de establecer “probabilidades previas” en muchos casos. Esto puede ser importante, ya que una probabilidad previa muy alta será más resistente a evidencia contraria, y una muy baja será más resistente a la evidencia a favor. Incluso, si el *prior* es 1, la hipótesis no puede revisarse. En casos de desconocimiento, algunos proponen asignar una probabilidad previa matemáticamente neutra, como 0.5. Esta estrategia se origina en el “Principio de Indiferencia” de Keynes, según el cual, en ausencia de evidencia, uno debería asignar probabilidades idénticas a todas las posibilidades en juego.

Veamos ahora un clásico ejemplo de razonamiento bayesiano.

Ejemplo. Supongamos que tengo una urna frente a mí, y sé que hay dos posibilidades: o bien la urna tiene 70 bolas rojas y 30 negras (urna A), o bien tiene 70 negras y 30 rojas (urna B).



Naturalmente, si saco una bola roja, voy a pensar que seguramente sea la urna A. Y si saco *dos* bolas rojas, voy a tener una creencia más convencida de que se trata de la urna A.

Esto lo podemos modelar perfectamente con la regla de Bayes. Supongamos que hasta ahora no sé cuál es la urna que tengo frente a mí, y no tengo ninguna pista al respecto. Usando el ya mencionado “Principio de Indiferencia”, le daremos una probabilidad de 0.5 a la urna A, y una probabilidad de 0.5 a la urna B. Es decir, $P(A) = 0.5$ y $P(B) = 0.5$. Ahora supongamos que saco una bola, y esa bola es roja (R). ¿Cómo actualizo mis creencias? Usando la regla de Bayes, lo puedo hacer del siguiente modo:

$$\begin{aligned} P(A | \text{Roja}) &= \frac{P(\text{Roja} | A) \times P(A)}{P(\text{Roja} | A) \times P(A) + P(\text{Roja} | B) \times P(B)} = \\ &= \frac{0.7 \times 0.5}{0.7 \times 0.5 + 0.3 \times 0.5} = \frac{0.35}{0.35 + 0.15} = \frac{0.35}{0.5} = 0.7 \end{aligned}$$

Es decir, usando la regla de Bayes, le terminaré atribuyendo a la hipótesis A la probabilidad 0.7.

Ahora, ¿qué sucede si encuentro una *nueva* bola roja? Supongamos, por simplicidad, que se trata de una urna “con reemplazo”, es decir, lo que saco lo vuelvo a poner. Ahora la probabilidad previa de A es 0.7, y actualizo del siguiente modo:

$$\begin{aligned} P(A | \text{Roja}) &= \frac{P(\text{Roja} | A) \times P(A)}{P(\text{Roja} | A) \times P(A) + P(\text{Roja} | B) \times P(B)} = \\ &= \frac{0.7 \times 0.7}{0.7 \times 0.7 + 0.3 \times 0.3} = \frac{0.49}{0.49 + 0.09} = \frac{0.49}{0.58} = 0.84 \end{aligned}$$

Así podemos ver que mi creencia en que tengo frente a mí la urna A crece cuando saco una segunda bola roja. Esto nos permite modelar procesos diacrónicos de aprendizaje, donde voy actualizando probabilidades, pero al mismo tiempo voy obteniendo evidencia nueva.

Ejercicios

1.

a. Pedro tiene problemas de insomnio, y se pregunta si podría estar causado por una depresión. Calcula la probabilidad de tener depresión, asumiendo que la prevalencia de depresión es 10%, el 60% de los depresivos tienen insomnio, y entre personas sin depresión la prevalencia de insomnio es 20%.

b. Pedro tiene ahora un ataque de pánico, y se preocupa más. Usando las probabilidades nuevas, calcula la probabilidad de que tenga depresión, si la prevalencia de ataques de pánico en depresivos es 50%, y en no-depresivos es de 10%. (Por simplicidad, asumimos que al tener depresión, la probabilidad de ansiedad e insomnio son independientes)

2.

a. El test de HIV a los 28 días es bastante confiable: si tienes el virus, es 95% probable que el resultado sea Positivo; si no lo tienes, es 99% probable que el resultado sea Negativo.

Calcule la probabilidad de tener HIV luego de un testeo negativo a los 28 días, en caso de que uno sea un/a trabajador/a sexual, donde la prevalencia de HIV es 3%.

b. El test de HIV a los 3 meses es mucho más confiable: si tienes el virus, es 98.5% probable que el resultado sea Positivo. Si no lo tienes, es 99.7% probable que el resultado sea Negativo.

Calcule la probabilidad de tener HIV luego de un testeo positivo a los 3 meses, para el promedio de la población, donde la prevalencia de HIV es 0.4%.

Parte D: Probabilidades y mundos posibles

Hay distintas formas de representar nociones probabilísticas. La forma más usual es la que mencionamos en las secciones anteriores: usando una función sobre proposiciones. Sin embargo, y especialmente en los textos de filosofía, también se suelen utilizar mundos posibles y álgebras. La idea contemporánea de *mundo posible*, que tiene su inspiración en la obra de Leibniz del

siglo XVIII, se origina en la obra de Kripke (1963) y Hintikka (1963). De acuerdo con estos autores, un mundo posible es una forma en que el mundo podría ser (o podría haber sido).

En el contexto de la epistemología, y principalmente a partir de la obra de Hintikka (1963), usamos los mundos posibles para hablar sobre *modos en que el mundo podría ser, de acuerdo con nuestro conocimiento*. Esto suele leerse de un modo subjetivo: por ejemplo, si jugamos dominó, sabemos que las fichas del otro pueden ser cualquiera salvo las nuestras y las que ya se tiraron al tablero. Todas esas posibilidades aún no descartadas son “mundos posibles” en este sentido puramente epistémico.

En términos técnicos, un *mundo posible* es una asignación de valor de verdad a todas las oraciones. En ese sentido, un mundo posible es equivalente a la noción de valuación en la lógica proposicional. En un mundo posible, todo está determinado: los jugadores de dominó tienen determinadas fichas, en Buenos Aires hace determinada temperatura, etc.

Ahora podemos explicar la noción de *álgebra*, que es particularmente relevante para el cálculo de probabilidades. Un *álgebra* F sobre un conjunto W de mundos posibles es un conjunto de subconjuntos de W que cumple con algunas condiciones:

1. El conjunto W está en F .
2. Si $A \in F$, entonces el $\underline{A} \in F$, donde \underline{A} es el complemento de A (relativo a W).
3. Si A y $B \in F$, entonces $A \cup B \in F$.

En otras palabras, un álgebra sobre un conjunto de mundos posibles W incluye el conjunto W , y está cerrada bajo unión y bajo complemento. El lector puede probar que también está cerrada bajo intersección.

Ejemplos

Sea $W = \{a, b\}$

Conjunto F : $\{\emptyset, W, \{a\}\}$

Conjunto G : $\{\emptyset, W\}$

Conjunto H : $\{\emptyset, W, \{a\}, \{b\}\}$

El conjunto H es un álgebra sobre W , porque contiene los complementos y las uniones de todos los conjuntos. De hecho, H es idéntico al conjunto potencia de W (es decir, el conjunto de subconjuntos de W). Veamos qué sucede con G . Aquí, el conjunto no contiene $\{a\}$ ni $\{b\}$, pero cumple con los axiomas, porque está cerrado bajo complemento y unión. Entonces también es un álgebra. Ahora bien, no sucede lo mismo con el conjunto F . Porque el conjunto F está cerrado bajo unión, pero no bajo complemento. El complemento de $\{a\}$ es $\{b\}$, pero $\{b\}$ no está en F . De modo que F no es un álgebra sobre W .

Ahora que tenemos el concepto de “álgebra”, podemos definir una función de probabilidad.⁹ Siendo F un álgebra sobre W , y μ una función de F en $[0,1]$, μ es una función de probabilidad si y sólo si cumple con los siguientes axiomas:

1. $\mu(W) = 1$
2. Si $A \cap B = \emptyset$, entonces $\mu(A \cup B) = \mu(A) + \mu(B)$

Podemos ver entonces la analogía entre los axiomas presentados en la Parte A y la descripción de una función de probabilidad. La Normalidad se infiere de la definición de la función μ , cuya imagen es el conjunto $[0,1]$. El primer axioma nos dice que W va a valer 1, y esto equivale a darle 1 a las tautologías (las oraciones verdaderas en todo mundo); se trata del axioma de Certeza. Mientras que el segundo axioma nos dice que la probabilidad de la unión entre dos conjuntos “incompatibles” es la suma de sus probabilidades por separado, y esto equivale al Axioma de Aditividad.

Por último, es importante señalar que de acuerdo con la tradición filosófica, podemos entender a las *proposiciones* como conjuntos de mundos posibles. Es decir, una proposición A puede entenderse como el conjunto de mundos posibles donde A es verdadera. En el contexto de las álgebras, no usamos *todos* los mundos sino un conjunto muy acotado, que elegimos según nuestros propósitos.

⁹ Véase Halpern (2003), cap. 2.

Por ejemplo, supongamos que salgo de la Facultad de Filosofía y quiero ir al centro en bus. Mientras espero, veo un bus viniendo a lo lejos. Tomando en cuenta que estoy en la Avenida Rivadavia (una avenida muy importante de Buenos Aires), infiero que puede ser la línea 132 (Flores - Retiro), la 145 (Mercado Central - Plaza Italia) o la 8 (Liniers - Retiro). Puedo representar el conjunto de mundos posibles de este modo:

$$W = \{8, 145, 132\}$$

Supongamos que tengo un álgebra muy fina, de modo tal que cualquiera de esos mundos posibles pertenece a ella como conjunto de único elemento¹⁰:

$$M = \{W, \emptyset, \{8\}, \{145\}, \{132\}, \{145, 8\}, \{8, 132\}, \{145, 132\}\}$$

Ahora supongamos que por un simple asunto de frecuencias (es decir, cuáles buses suelen pasar más seguido), esta es mi atribución de probabilidades:

$$\mu(\{145\}) = 0.3, \mu(\{8\}) = 0.2, \mu(\{132\}) = 0.5$$

Podemos notar dos cosas. Primero: al atribuir probabilidades a los conjuntos más pequeños, ya podemos inferir las probabilidades de todos los otros conjuntos (son simplemente las sumas). Segundo, que estas probabilidades deben sumar 1.

Ahora bien, podríamos definir muchas proposiciones a partir de los mundos donde son verdaderas. Aquí lo leeremos en el contexto epistémico particular en el que nos encontramos. Por ejemplo:

$$\text{“El bus me lleva a la Estación de Retiro”} = \{8, 132\}$$

$$\text{“El nombre del bus tiene tres cifras”} = \{145, 132\}$$

¹⁰ Cuando esto sucede, el álgebra será el conjunto potencia de W .

Como sabemos, dado que el conjunto de mundos es muy acotado, un mismo conjunto puede expresarse de distintos modos. Por ejemplo, en este contexto, afirmar que el bus tiene una cifra equivale a decir que sale de Liniers.

Una cuestión interesante es cómo representar en este enfoque las probabilidades condicionales. Afortunadamente, esto es bastante sencillo. En general, podemos decir (asumiendo que $P(B) \neq 0$):

$$P(A | B) = P(A \cap B) / P(B)$$

Por cuestiones puramente matemáticas, el resultado será siempre un número entre 0 y 1.

Ejemplo. Volvamos al ejemplo anterior. Tenemos el mismo modelo con las mismas atribuciones de probabilidad. Ahora veo al bus un poco más cerca, y noto que el cartel tiene tres números. ¿Cuál es ahora la probabilidad de que sea el bus 145? Primero, debo entender cuál fue la evidencia. La evidencia en este caso, en este contexto específico, es $\{145, 132\}$.

Entonces la nueva probabilidad de que sea el 145 es:

$$\begin{aligned} P(\{145\} | \{145, 132\}) &= P(\{145\} \cap \{145, 132\}) / P(\{145, 132\}) = \\ &= P(\{145\}) / P(\{145, 132\}) = 0.3 / 0.8 = 0.375 \end{aligned}$$

Ahora la probabilidad de que el bus sea el 145 pasa a ser 0.375.

Recientemente, Leitgeb (2014) utilizó este aparato teórico para dar respuesta a la Paradoja de la Lotería (Kyburg 1961). Esta paradoja intenta reducir al absurdo la *Tesis de Locke*, según la cual podemos creer racionalmente una proposición si su probabilidad es mayor a un parámetro r (supongamos, 0.8). El argumento dice lo siguiente: si existiera una lotería de 10 tickets, la chance de que *no* salga cada uno es 0.9. Entonces, por la Tesis de Locke, debo creer “No va a salir el ticket 1”, “No va a salir el ticket 2”, etc. Pero también debo creer que “Va a salir algún ticket” (cuya probabilidad es 1). Esto me lleva a una inconsistencia.

La respuesta de Leitgeb es que las proposiciones creíbles racionalmente no requieren solo probabilidad alta sino también *estabilidad*. Esto significa que resisten bien a la información nueva. Formalmente, una proposición A es *estable* si y sólo si $P(A) > 0.5$, y además $P(A|B) > 0.5$ para toda proposición B del álgebra compatible con A , es decir, $A \cap B \neq \emptyset$. En el contexto de la lotería, no es estable creer que “No va a salir el ticket 1”, es decir $\{2, 3, 4, 5, 6, 7, 8, 9, 10\}$, porque condicionalizado bajo “salieron el 1 o el 2”, es decir $\{1, 2\}$, nos arroja $P(\{2\}) / P(\{1, 2\}) = 0.5$. La única proposición estable es que “Va a salir algún ticket”.

Ejercicio

1. En el escenario anterior, veo de lejos un bus que dice que va a Retiro. Calcule la probabilidad de que sea la línea 8.
2. Pruebe que $\{132, 145\}$ es una proposición estable.

Parte E: Probabilidad y Lógica

¿Cuál es la relación entre la lógica y la probabilidad? Como vimos, los axiomas de Kolmogórov establecen algunas relaciones: las tautologías tendrán probabilidad 1, y si A implica B , entonces $P(A) \leq P(B)$. Sin embargo, las inferencias válidas suelen tener más de una premisa. Una buena pregunta es qué podemos inferir sobre la relación entre lógica y probabilidad en esos casos.

Uno de los resultados más conocidos se lo debemos a Ernest Adams. En su libro *A primer on Probability Logic* (1998), prueba una importante inecuación:

Ley de Adams: Si A_1, \dots, A_n implican lógicamente B , entonces
$$P(A_1) + \dots + P(A_n) - (n - 1) \leq P(B)$$

La idea de Adams es que, si sumamos la ‘incerteza’ de las premisas, esta suma debe ser mayor a la ‘incerteza’ de la conclusión. Por ejemplo, si A, B implican C , y $P(A) = 0.8$, y $P(B) = 0.9$, la suma de incertezas es 0.3, de modo tal que la incerteza de C debe ser a lo sumo 0.3, es decir, la probabilidad de C debe ser como

mínimo 0.7. Si la probabilidad de C fuera 0.6, la incerteza sería 0.4, y no se cumpliría esa inecuación.

Es fácil ver que el Teorema de Validez para argumentos con una sola premisa es una instancia de este principio, donde $n = 1$. Por otro lado, también podemos ver que, si el argumento es válido y la probabilidad de las premisas es 1, la conclusión debe también tener probabilidad 1. La Ley de Adams se infiere de los axiomas de Kolmogórov (Adams 1998, p. 32).

Resultados como el de Adams son importantes para la filosofía porque establecen los paralelos posibles entre la probabilidad y la lógica. Además, nos sirven para comprender determinados escenarios epistémicos.

Un famoso escenario es conocido como *Paradoja del Prefacio*. En este escenario, una profesora escribe un libro de historia luego de años de investigación. En el prefacio, sin embargo, dice: “como la historia es una ciencia empírica, este libro contiene errores, de los que me hago totalmente responsable”. La profesora admite que el libro contiene algunos errores. Al mismo tiempo, si le preguntan, ella cree cada oración del libro individualmente. Es decir, la profesora cree A_1, \dots, A_n , para todas las oraciones A_i del libro, y a la vez cree $\neg(A_1 \& \dots \& A_n)$. Según Makinson (1965), este caso muestra que podemos ser inconsistentes y racionales al mismo tiempo.

La paradoja del prefacio tuvo decenas de respuestas. Desde un punto de vista meramente probabilístico, apelando a la ley de Adams, la solución es obvia. La probabilidad permite que las incertezas se vayan sumando, de modo tal que la creencia en una conjunción gigante de oraciones no del todo seguras, podría tener una probabilidad muy baja. Muchos autores consideran que esto no termina de resolver el problema, porque incluso si aceptamos que todas las oraciones dentro de un conjunto inconsistente podrían tener probabilidad alta, es difícil aceptar que un conjunto lógicamente inconsistente sea aceptable como conjunto racional de creencias.

Ejercicios

1. ¿Puede suceder que $P(A)$ y $P(\neg(A \vee B))$ sean ambas mayores a 0.5?
2. Supongamos que A y B implican C. Además, $P(A) = 0.9$ y $P(B) = 0.9$. ¿Puede suceder que $P(C) = 0.6$? ¿Cuál es la probabilidad permitida para C según el teorema de Adams?
3. Piense un ejemplo de A y B donde $P(A)$, $P(B)$ y $P(\neg(A \& B))$ son al mismo tiempo mayores a 0.5 (por ejemplo, tirando un dado).

Parte F: Filosofía de la probabilidad

¿Qué significa que la probabilidad de un evento sea 0.5? Debería ser una pregunta fácil de responder.

Hemos visto que, para la perspectiva bayesiana, la probabilidad es entendida como una estimación que hace un agente sobre las posibilidades disponibles. Llamamos a esa noción *probabilidad subjetiva*. La idea de la probabilidad subjetiva es que las probabilidades son principalmente representaciones de las actitudes subjetivas de los agentes. Podemos ver a las probabilidades como simples hipótesis. Dentro de la epistemología formal suele usarse esta lectura subjetiva de la probabilidad. Uno de los primeros filósofos en proponer una lectura subjetivista de la probabilidad fue Frank Ramsey, en “Truth and probability” (1926).

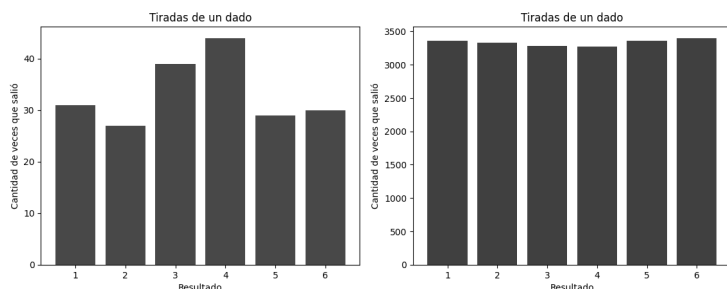
Como vimos anteriormente, la lectura bayesiana de la probabilidad se enfrenta a algunos problemas conceptuales, como la imposibilidad de tener criterios justificados para otorgarle a los eventos su “probabilidad inicial” o *prior*.

Existen también lecturas más *objetivas* de la probabilidad. Por ejemplo, si tiro una moneda que no está sesgada, la probabilidad de que salga cara y de que salga cruz es la misma. Esto no depende de lo que yo crea, sino de cómo funciona el *azar*.

El concepto de *azar* es muy problemático en física, en matemática y en filosofía. Pero en casos sencillos, siempre se asume que ciertos fenómenos tienen ciertas tendencias. Por ejemplo, en el caso anterior, una moneda “justa” tiende a caer cara en un 50% de las veces. Obviamente, esta probabilidad no proviene de una simple inducción: si tiro la moneda cinco veces, *puede pasar* que

salga cara cuatro de esas cinco veces. Ahora bien, ¿qué pasa si tiramos la moneda un número *infinito* de veces? Siguiendo la *Ley de los Grandes Números* (que veremos más adelante), mientras más grande sea el número de intentos, más cerca estará la proporción de su tendencia natural.

Podemos observarlo con la siguiente simulación (hecha en Python), donde tiré un dado muchas veces y anoté cuántas veces salió cada número. En el gráfico de la izquierda, tiré el dado 200 veces, y en el de la derecha, 20.000 veces. Podemos ver que, al tirar el dado más veces, la frecuencia termina acercándose a su probabilidad original.



El *frecuentismo* es la interpretación según la cual una probabilidad expresa una frecuencia relativa. Es decir, la probabilidad de que la moneda salga cara es $\frac{1}{2}$ porque si tiráramos la moneda infinitas veces, la mitad de ellas saldría cara. La lectura frecuentista es la más común en ciencias empíricas como la física. Otras disciplinas como la medicina usan visiones frecuentistas en la mayoría de los casos, y otras veces usan lecturas subjetivistas bayesianas. Un defensor del frecuentismo en su versión contemporánea fue Richard Von Mises (1928).

Un problema de las teorías frecuentistas es la apelación a una serie infinita de repeticiones de un experimento, que por obvias razones no puede llevarse a cabo en la realidad. Sin embargo, un frecuentista podría observar que, si se repiten los experimentos suficientes veces, obtendremos una buena estimación de lo que pasaría en una serie infinita.

Otro problema para el frecuentismo es la incapacidad de explicar la probabilidad de hechos únicos. Por ejemplo, ¿cuál es la probabilidad de que Argentina derrotara a Inglaterra en la Guerra de Malvinas? No podemos repetir la guerra infinitamente, ni un número suficientemente grande de veces como para saberlo.

Para estos problemas existen otras teorías objetivas, pero no frecuentistas, de la probabilidad. Por ejemplo, la lectura objetivista de la probabilidad como *propensión*, propuesta por autores como Karl Popper (1983). Para Popper, la probabilidad de un hecho proviene de la configuración física de la realidad. Por ejemplo, la probabilidad de que la moneda salga cara es simplemente una consecuencia de la forma y la composición material de la moneda.

Otra lectura objetivista es la interpretación *lógica* de la probabilidad, de autores como Keynes (1921). Keynes propone que, en ausencia de información, la probabilidad de eventos indistinguibles lógicamente será idéntica: lo llamamos Principio de Indiferencia. Por eso, la probabilidad de que la moneda salga cara es $\frac{1}{2}$ y la probabilidad de que un dado salga 3 es $\frac{1}{6}$.

Algunos filósofos han discutido las posibles relaciones entre la probabilidad subjetiva y la probabilidad objetiva. La lectura subjetiva parece dejar demasiada libertad a los agentes, porque solo se preocupa por la consistencia de las creencias. Uno de los principios que suelen postularse para unir ambas concepciones de la probabilidad es el *Principio Principal* de David Lewis (1986). Este principio dice (a grandes rasgos) que, cuando la probabilidad objetiva es conocida, la probabilidad subjetiva debe ser igual a la probabilidad objetiva. Por ejemplo, si estamos tirando una moneda justa, la probabilidad objetiva de que salga “cara” es $\frac{1}{2}$, entonces la probabilidad subjetiva que deberías asignarle a ese evento es también $\frac{1}{2}$.

Parte G: Coherencia probabilística y racionalidad

Una pregunta que se han hecho muchos filósofos es la siguiente: ¿Qué nos fuerza a nosotros, los seres humanos, a cumplir con los axiomas de la probabilidad? Por ejemplo: ¿Por qué, si creemos

que la probabilidad de A es $\frac{1}{2}$, deberíamos creer que la probabilidad de $\neg A$ también es $\frac{1}{2}$?

Argumentos pragmáticos: Dutch Books y explotabilidad

Una idea fundamental dentro del enfoque bayesiano es que la incoherencia provoca *explotabilidad*. Es decir, si no cumplimos con los axiomas probabilísticos en nuestras creencias, un agente suficientemente perspicaz podrá estafarnos y sacarnos todo nuestro dinero.¹¹ La forma en que suele explicarse la explotabilidad es mediante apuestas en donde estamos condenados a perder.

Por ejemplo, supongamos que yo estoy segurísimo de que Boca Juniors (el equipo de fútbol más popular de Argentina) va a ganar la final de la Copa Libertadores. Estoy 99% seguro. Entonces un amigo me ofrece una apuesta:

Apuesta 1: Me das un dólar ahora. Si gana Boca, te doy 100 dólares. Si pierde, me quedo con el dólar.

Obviamente, si soy suficientemente racional, voy a entrar en esa apuesta. Se trata de una apuesta muy conveniente: con sólo poner 1 dólar, sé que muy posiblemente voy a ganar 100.

Pero una buena pregunta es: ¿hasta qué punto voy a seguir apostando? Asumiendo que en estos casos la utilidad del dinero es lineal (esto lo discutiremos en detalle en el próximo capítulo), si estoy 99% seguro de que ganará Boca, podría entrar en esta apuesta:

Apuesta 98: Me das 98 dólares ahora. Si gana Boca, te doy 100 dólares. Si pierde, me quedo con los 98 dólares.

Incluso debería ser *indiferente* respecto a esta apuesta:

Apuesta 99: Me das 99 dólares ahora. Si gana Boca, te doy 100 dólares. Si pierde, me quedo con los 99 dólares.

¹¹ La idea de que la irracionalidad genera explotabilidad es más general, no sólo aplica a los axiomas de la probabilidad sino presuntamente también a cualquier otro principio de la racionalidad.

Decimos que una apuesta es *justa* cuando, tomando en cuenta nuestras probabilidades subjetivas, no tendríamos por qué rechazarla. La apuesta 99 es justa.

En términos más generales, si lo que está en juego en la apuesta es S (es decir, S es el “pozo” de la apuesta), y yo creo que el evento A va a pasar con probabilidad p , debería estar dispuesto a apostar $\$pS$ a que el evento A va a suceder, para ganar $\$(1 - p)S$. Por ejemplo, si yo creo en 0.99 que Boca va a ganar, y el pozo son 100.000 dólares, debería estar dispuesto a apostar 99.000 dólares para ganar 1.000. Esto no significa que voy a estar *feliz* de apostar de ese modo, sino simplemente que no tengo ningún motivo para no apostar así.

Todos estos conceptos nos sirven para entender por qué debemos cumplir con los axiomas probabilísticos. Pues supongamos que creo que Boca va a ganar con 0.6 y también creo que Boca *no* va a ganar con 0.6. Aquí, obviamente, incumplo los axiomas de Kolmogórov. Una forma de mostrar que se trata de un tipo de irracionalidad es que un agente podría ofrecirme este *conjunto de apuestas*:

Apuesta A: Me das 51 dólares. Si gana Boca, te doy 100 dólares. Si no gana, me quedo con tus 51 dólares.

Apuesta B: Me das 51 dólares. Si no gana Boca, te doy 100 dólares. Si gana, me quedo con tus 51 dólares.

Por creer que Boca va a ganar con 0.6 de probabilidades, voy a entrar en la apuesta A. Y por creer que Boca *no* va a ganar con 0.6 de probabilidades, voy a entrar en la apuesta B. El problema es que ahora, pase lo que pase, voy a terminar perdiendo dinero. Porque le di al agente 102 dólares. Si Boca gana, obtengo 100 dólares, y si pierde también. Pero en cualquier caso voy a tener una pérdida neta de 2 dólares.

El teorema de *Dutch Book* nos muestra que *cualquier* violación al cálculo de probabilidades nos vuelve pasibles de ser explotados de este modo: alguien nos puede ofrecer un conjunto de apuestas que estamos racionalmente obligados a aceptar, pero que nos dará finalmente pérdida, en cualquier caso.

***Una prueba del Teorema de Dutch Book**

De aquí en más supondremos que el pozo de la apuesta, es decir S , es \$1, así que podremos obviarlo. También por legibilidad vamos a ignorar los signos de dólar “\$” cuando no sean estrictamente necesarios.

Teorema de Dutch Book. Si nuestras probabilidades personales no satisfacen los axiomas de la probabilidad, nos enfrentamos a una apuesta de pérdida segura.

Prueba. Recordemos que en una apuesta justa, cuando mi creencia en A es p , voy a apostar p a favor de A (o $1-p$ contra A), y obtendré los siguientes resultados, dado que el pozo es 1:

	Resultado de apuesta por A	Resultado de apuesta contra A
Sucede A	$\$(1 - p)$	$\$-(1 - p)$
Sucede $\neg A$	$\$-p$	$\$p$

Ahora veremos que, frente a cada violación de los axiomas, un corredor nos puede ofrecer una apuesta “justa” o un conjunto de apuestas “justas” con pérdida segura.

Axioma 1: Normalidad

Si $P(A) = p$, requerimos que $0 \leq p \leq 1$.

Ahora supongamos que un agente incumple este principio. Pueden pasar dos cosas: o bien le asigna menos de 0, o bien le asigna más de 1.

A. Supongamos $p < 0$

[Creo demasiado poco en A]

Un corredor nos puede ofrecer una apuesta justa en contra de A , por lo que usaremos esta parte de la tabla:

	Resultado de apuesta contra A
Sucede A	$\$-(1 - p)$
Sucede $\neg A$	$\$p$

Resultado: Asumimos que $p = -r$, donde r es un número positivo (p es un número negativo). Si sucede $\neg A$, obtengo $-r$. Si sucede A, obtengo $-(1 - (-r)) = -(1 + r)$. En ambos casos pierdo dinero.

B. Supongamos $p > 1$

[Creo demasiado en A]

Un corredor nos puede ofrecer una apuesta justa a favor de A, por lo que usaremos la otra parte de la tabla:

	Resultado de apuesta por A
Sucede A	$\$(1 - p)$
Sucede $\neg A$	$\$-p$

Resultado: Supongamos que $p = 1 + r$. Si sucede $\neg A$, el resultado es $-(1 + r)$. Si sucede A, obtengo $1 - (1 + r) = 1 - 1 - r = -r$. En ambos casos pierdo dinero.

Axioma 2: Certeza

Supongamos que A sucede seguro, pero $p < 1$.

[Creo demasiado poco en A].

Un corredor nos puede ofrecer una apuesta justa contra A:

	Resultado de apuesta contra A
Sucede A	$\$-(1 - p)$
Sucede $\neg A$	$\$p$

Seguro pierdo $(1 - p)$, porque no es posible que suceda $\neg A$.

Axioma 3: Aditividad

El axioma 3 de Kolmogórov nos dice que si A y B son mutuamente excluyentes, entonces $P(A \vee B) = P(A) + P(B)$.

Usaremos esta nomenclatura para las creencias en cada proposición: $P(A) = p$, $P(B) = q$ y $P(A \vee B) = r$. Dados los axiomas debería suceder que $r = p + q$. Ahora supongamos que esto no sucediera. Esto podría pasar de dos formas distintas.

A. Supongamos $r < p + q$

[creo demasiado en A y B, y poco en $A \vee B$]

Entonces el corredor nos ofrece las siguientes apuestas:

1. Apostar p a favor de A
2. Apostar q a favor de B
3. Apostar $(1 - r)$ contra $A \vee B$

Debería aceptar estas apuestas, porque son “justas”. Ahora observemos que hay tres posibles estados del mundo: dado que A y B son mutuamente excluyentes, es imposible que suceda $A \& B$. Y en los tres posibles estados, terminaré perdiendo dinero.

Resultado:

	Pago por 1	Pago por 2	Pago por 3	Pago final
Sucede $A \& \neg B$	$\$(1-p)$	$\$-q$	$\$-(1-r)$	$\$(1-p-q-1+r) = \$(r-p-q)$
Sucede $\neg A \& B$	$\$-p$	$\$(1-q)$	$\$-(1-r)$	$\$(-p+1-q-1+r) = \$(r-p-q)$
Sucede $\neg A \& \neg B$	$\$-p$	$\$-q$	$\$r$	$\$(r-p-q)$

Siempre voy a obtener una misma cantidad $(r-p-q) = r - (p + q)$, que será negativa, porque asumimos que $r < p + q$.

Para una prueba completa también debo probar la parte [B], es decir, que hay un contrato de pérdida seguro si $r > p + q$. Esto es simétrico y queda al lector.

Con esto hemos probado que, si nuestras creencias no satisfacen los axiomas probabilísticos, nos enfrentamos a apuestas de pérdida segura.

QED

***Argumentos epistémicos: adecuación al mundo**

Algunos autores desarrollaron justificaciones puramente epistémicas, no pragmáticas, para la coherencia probabilística. Así como los argumentos pragmáticos nos muestran que las asignaciones no-probabilísticas nos vuelven posibles víctimas de apuestas de pérdida segura, los argumentos epistémicos muestran que las asignaciones no-probabilísticas nos alejan de la verdad y nos acercan al error.

Con más detalle, el argumento epistémico que presentaré en esta sección nos dice que, si hubiera distintos escenarios posibles, un estado epistémico probabilísticamente incoherente siempre (en todos los escenarios) va a ser menos adecuado a la realidad que uno probabilísticamente coherente. Para probar este resultado, necesitamos un método para *medir* la adecuación de una probabilidad respecto a un hecho. Por suerte, existe un método de medición para estos propósitos: el *puntaje de Brier*, llamado así por su creador, el meteorólogo Glenn Brier (1913-1998).

El puntaje de Brier mide la *distancia* entre una probabilidad y la realidad. Decimos que la realidad r tiene dos resultados posibles: verdad (1) o falsedad (0). El puntaje de Brier para una asignación de probabilidad $P(e)$ para un evento e se define del siguiente modo:

$$\text{El puntaje de Brier de } P(e) \text{ es } (P(e) - r)^2$$

La aplicación natural para el puntaje de Brier es la meteorología. Por ejemplo, si llueve (1) y la probabilidad que le asigno a que llueva es 0.6, mi grado de adecuación será $(0.6 - 1)^2 = (-0.4)^2 = 0.16$. Si la probabilidad que le asigno a que llueva es 0.1, mi grado de adecuación será $(0.1 - 1)^2 = (-0.9)^2 = 0.81$. Mientras más

alto sea el número de Brier, más *inadecuado* es mi estado epistémico: un agente omnisciente tendría un puntaje de Brier de 0. Aquí voy a desarrollar una prueba muy escueta de la coherencia por adecuación, para un par de eventos, p y $\neg p$, siguiendo a Fittelson (inédito). Dado que no asumimos la coherencia de los agentes, la probabilidad asignada a p será independiente de la asignada a $\neg p$. Nos interesan dos mundos, el mundo w_1 donde p es verdadera, y el mundo w_2 donde p es falsa. Usando la fórmula anterior, los puntajes en cada mundo pueden calcularse así:

- El puntaje en w_1 (donde p es verdadera) es $(P(p) - 1)^2 + P(\neg p)^2$
- El puntaje en w_2 (donde p es falsa) es $P(p)^2 + (P(\neg p) - 1)^2$

De Finetti mostró el siguiente teorema:

Teorema (De Finetti 1970): La asignación de probabilidad P es no probabilística si y sólo si hay otra asignación de probabilidad P' cuyo puntaje de Brier es más bajo (es decir, está más cerca de la realidad) en todo mundo posible.

Antes de hacer un esbozo de prueba, pensemos un ejemplo. Supongamos que $P(p) = 0.7$ y $P(\neg p) = 0.7$, siendo nuestra asignación obviamente irracional (dado que el grado de creencia en p y el de $\neg p$ deberían sumar 1). Si esto sucede, entonces el puntaje de Brier en caso de que p sea falso es $0.7^2 + 0.3^2 = 0.49 + 0.09 = 0.58$, mientras que cuando p es verdadero, el puntaje de Brier es también .58. Ahora pensemos en la asignación $P'(p) = P'(\neg p) = 0.5$. En caso de que p sea falso, el puntaje de Brier es $0.5^2 + 0.5^2 = 0.25 + 0.25 = 0.5$, mientras que cuando p es verdadero el puntaje será también $0.25 + 0.25 = 0.5$. Aquí vemos que P' , la asignación coherente, tiene un menor puntaje de Brier que P (la asignación incoherente) en ambos mundos posibles.

Ahora necesitamos probar que *siempre* que haya una asignación incoherente, hay otra asignación coherente con mejor puntaje de Brier en todo mundo posible.

Para visualizar informalmente la prueba, lo más sencillo es pensar en estos gráficos donde aparece tanto el valor de verdad de p como nuestro grado de creencia. El eje horizontal representa la

asignación de probabilidad a p , y el vertical representa la asignación de probabilidad a $\neg p$:

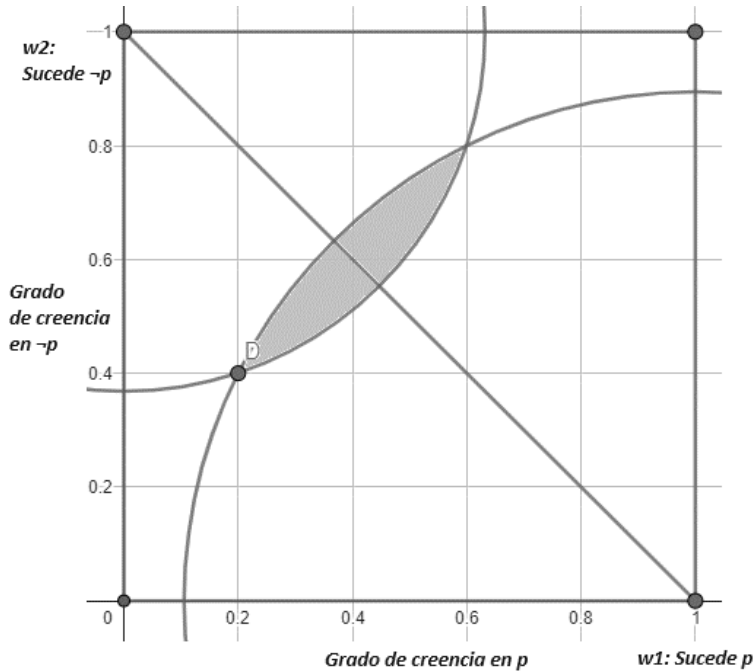


Gráfico 1

El plano también sirve para representar el mundo w_1 , donde p es verdadera, y el mundo w_2 , donde p es falsa. Cada asignación de probabilidad a p y $\neg p$ es representable como un punto en el plano. La línea diagonal representa las asignaciones coherentes de probabilidad (por ejemplo, 0.6 a p y 0.4 a $\neg p$).

En el Gráfico 1 podemos ver un punto marcado D, probabilísticamente incoherente (creencia de 0.2 en p y de 0.4 en $\neg p$) y dos circunferencias que pasan por él. Los puntos en una circunferencia son los que están a la misma distancia de w_1 en un caso, y en la otra circunferencia son los que están a la misma distancia de w_2 . De este modo, cada circunferencia también delimita un conjunto de puntos que son *más* cercanos a la realidad, ya sea w_1 o w_2 . La sección sombreada representa el conjunto de puntos que

son más cercanos a la realidad *en ambos mundos*, w_1 y w_2 , respecto al punto D que habíamos marcado.

Es sencillo observar que, dentro de ese conjunto, hay puntos que constituyen una asignación probabilística (i.e. aquellos que pertenecen a la diagonal). El argumento puede generalizarse para cualquier punto incoherente en ese plano. Así probamos un lado del enunciado de DeFinetti: siempre habrá un punto en la diagonal que *domina* al punto que no está en la diagonal, es decir, que está más cerca de la realidad cualquiera sea el hecho que suceda (ya sea w_1 o w_2).

Veamos ahora el otro lado del teorema.

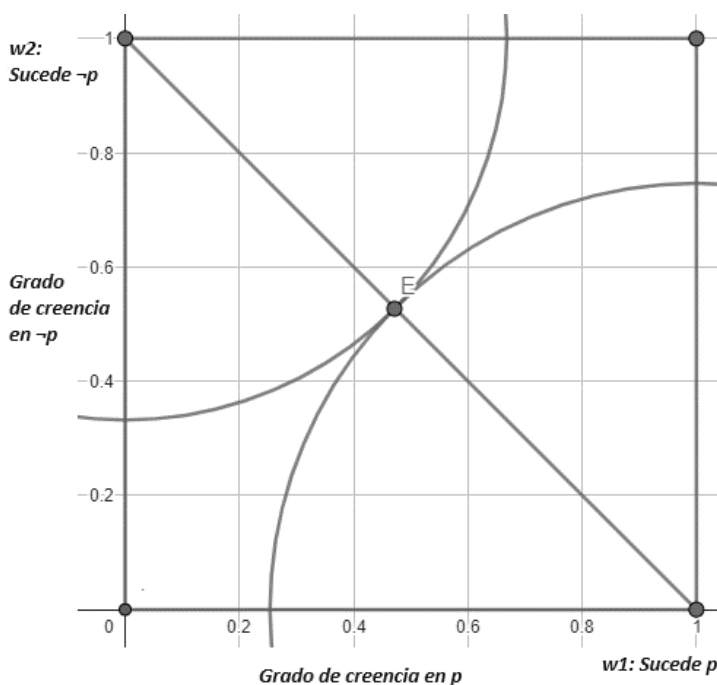


Gráfico 2

El Gráfico 2 (donde la asignación de probabilidades es aproximadamente 0.45 a p y 0.55 a $\neg p$) nos muestra que *sólo* una asig-

nación probabilística puede evitar que haya “espacios sombreados”, es decir, puntos más cercanos a la realidad en ambos mundos. Por razones geométricas, sólo los puntos que pertenecen a la diagonal pueden cumplir con esa propiedad.

Un aspecto frecuentemente observado de los resultados de adecuación es que no nos muestran que *cualquier* distribución probabilística sea de hecho mejor que *cualquier* distribución incoherente (Kolodny 2007); lo que muestran es que, si tenemos una distribución incoherente, podemos estar seguros de que habrá una distribución coherente que la supere en todos los escenarios posibles. Por eso, la aplicación práctica de este resultado dependerá de que seamos capaces de encontrar esa distribución coherente que supere a la distribución incoherente que tenemos.

En un artículo reciente, De Bona y Staffel (2018) desarrollan una manera de encontrar una asignación óptima a partir de una asignación incoherente, buscando el camino más corto entre la asignación incoherente y el conjunto de asignaciones coherentes (es decir, la diagonal en el plano). En el Gráfico 1, por ejemplo, el punto más cercano dentro de la diagonal al punto D (0.2, 0.4) es (0.4, 0.6).

Ejercicios

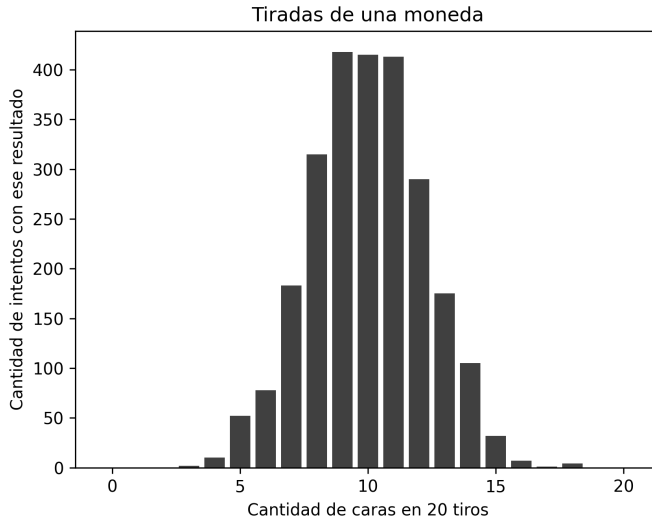
1. Probar usando Dutch Book que la probabilidad de A y la de $\neg A$ sumadas no pueden ser mayores a 1.
2. Respecto a la prueba por adecuación:
 - a. Creo que va a llover (p) con 0.6 y que no va a llover ($\neg p$) también con 0.6. Calcular el puntaje de Brier en los mundos posibles donde llueve y donde no llueve.
 - b. Comparar con el puntaje de Brier coherente cuando creo que va a llover con 0.6 y que no va a llover con 0.4. ¿Domina este al caso del ejercicio 2?
 - c. ¿Qué sucede si $P(\text{lloverá})$ es 0.5 y $P(\text{no lloverá})$ es 0.5? ¿Domina este al caso del ejercicio a?
 - d. Explique lo sucedido a partir del teorema de adecuación de De Finetti.

Parte H. Nociones de estadística

Las reglas básicas de la probabilidad pueden servir tanto para probabilidades subjetivas como para probabilidades objetivas. Para las probabilidades subjetivas, la regla fundamental es la de Bayes. Para las probabilidades objetivas, el resultado fundamental es el de Jakob Bernoulli (1654-1705); se considera que Bernoulli descubrió la “ley de grandes números”.

La ley de los grandes números nos dice que, si la repetición de un evento con probabilidad p tiende a infinito, la frecuencia relativa será más cercana a p . Es decir, si tiro 5 veces una moneda, casi cualquier cosa podría pasar; pero si tiro la moneda un millón de veces, aproximadamente saldrá la mitad “cara” y la mitad “cruz”. Sin embargo, decimos “aproximadamente”. El resultado raramente será la mitad exacta (es decir, 2.5 millones de veces “cara”, y 2.5 “cruz”). Pero tampoco podría desviarse demasiado hasta los costados. Estos son los temas que explora la estadística. Supongamos que tiro 20 veces una moneda. ¿Cuántas veces sale “cara”? Podría salir 10 veces, pero también 9 u 8 veces. Mucho más raro es que saliera 20 veces “cara”, o 20 veces “cruz”. Ahora bien, podríamos repetir el experimento muchas veces y anotar los resultados.

El siguiente gráfico es el resultado de una simulación (realizada en Python). Una simulación es un proceso computacional que intenta replicar otro fenómeno, que generalmente nos llevaría demasiado tiempo calcular manualmente. En este caso, tiramos una moneda 20 veces y anotamos cuántas veces sale “cara”. Eso podríamos hacerlo manualmente una vez (solo es cuestión de tirar una moneda 20 veces). Pero usamos la simulación para repetir este proceso 2500 veces, para darnos una idea general del fenómeno. Los lectores más escépticos podrían tirar una moneda 20 veces y calcular cuántas veces sale “cara”; y repetir este proceso 2500 veces manualmente (no debería llevar más que un par de días).



El resultado, como podemos ver, es una distribución en forma de “campana”, con barras más altas en el medio, y barras más bajas a los costados. ¿Cuál fue el promedio? A esto lo llamamos *promedio muestral*. En mi simulación me dio 9.96. Es decir, de cada 20 tiros, en promedio 9.96 salen cara. Aquí se cumple la predicción de Bernoulli, según el cual el promedio debería ser aproximadamente 10, es decir, $p \times n$.

También podemos observar que, en el gráfico, casi todos los resultados están entre 5 y 15, concentrándose la mayor parte alrededor del promedio. Para entender este fenómeno podemos introducir el concepto de *desviación estándar*. Se trata simplemente de una especie de promedio, el promedio de la desviación entre el promedio de la muestra y los resultados. Es decir, medimos la *dispersión* a partir del promedio X . Queremos obtener un número d tal que la gran mayoría de los resultados estén entre $(X + d)$ y $(X - d)$.

Para calcular la desviación a partir de una muestra, lo primero que podríamos hacer es notar las diferencias entre el promedio X

y cada resultado observado X_i . En nuestro caso, hay 2500 desviaciones. El promedio de las desviaciones (usualmente llamado “desviación media”) me dio 1.81.

La forma típica de medir la dispersión, usando el concepto de *desviación estándar*, no busca directamente el promedio de las desviaciones, sino algo un poco distinto. La idea es:

1. Calcular el *cuadrado* de cada desviación (número positivo). Así evitamos las desviaciones negativas.
2. Sumar todos esos cuadrados.
3. Dividir esa suma por la cantidad de resultados.
4. Sacar la raíz cuadrada de esa suma. Así compensamos el cuadrado calculado en el punto 1.

Es decir:

Siendo n la cantidad de tiros, X_i el resultado de cada tiro, y X el promedio:

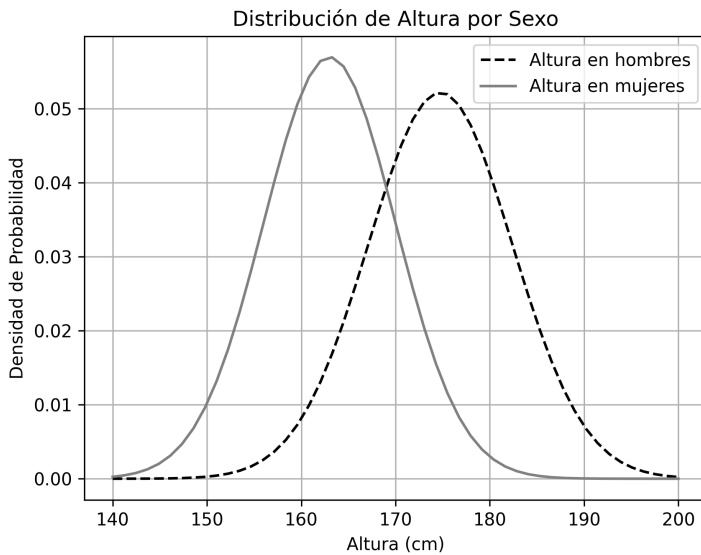
$$\text{Desviación estándar} = \sqrt{\frac{\sum_i (X_i - X)^2}{n}}$$

El resultado para nuestra simulación fue 2.27. Si vemos el gráfico, resultará que la gran mayoría de los resultados (algo así como el 68%) están entre el promedio y la primera desviación (es decir, aproximadamente entre 8 y 12). Si el experimento incluyera más tiros, la desviación sería menor.

Podemos observar que en este experimento estamos tirando una moneda, que podría salir cara o cruz. Llamamos a estos ejemplos *intentos binomiales*. En estos casos, hay dos eventos (por ejemplo, cara o cruz, verdad o falsedad, etc.), la probabilidad de que salga el primero (“éxito”) en vez del segundo (“fracaso”) es p , y yo pruebo n veces. Como ya señalamos, si n es suficientemente grande, la cantidad de “éxitos” será aproximadamente $p \times n$, por la ley de los grandes números. Más adelante veremos cómo podemos predecir la desviación estándar en estos casos.

Aproximaciones normales

Observemos ahora el gráfico de antes sobre cuántas veces sale cara si tiro una moneda 20 veces. Como señalamos, el gráfico tiene la forma (aproximada) de una campana. Muchas distribuciones tienen esta forma (mucho más, cuando se trata de fenómenos físicos o naturales). Por ejemplo, la altura en hombres y la altura en mujeres en las Islas Canarias tienen esta distribución (normalizada):¹²



Estas curvas suelen llamarse *normales* o *gaussianas*. Todas estas curvas tienen dos elementos:

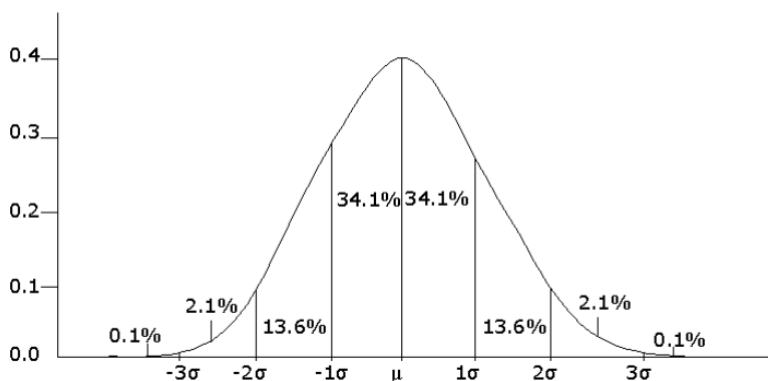
- El promedio o pico, que suele ser representado con μ .
- La desviación estándar, que suele ser representada con σ .

¹² Datos del Instituto Canario de Estadística de 2021.

En el caso de las curvas normales, se dan algunas regularidades que son muy relevantes para lo que veremos después. En una curva normal, pasará lo siguiente:

- 68% de los casos estarán entre el pico y la primera desviación.
- 95% de los casos estarán entre el pico y la segunda desviación.
- 99% de los casos estarán entre el pico y la tercera desviación.

Podemos diagramarlo así:



En general, cuando se trata de un fenómeno binomial (como tirar una moneda), y necesito calcular la cantidad de “éxitos” en una serie de intentos, es fácil calcular el promedio y la desviación estándar, siguiendo el método propuesto por Bernoulli. Si n es la cantidad de veces que hago el intento, y p es la probabilidad de “éxito” en cada caso, el promedio de veces que sucederá el evento podemos calcularlo con la fórmula $\mu = p \times n$. Para la desviación estándar, podemos usar la fórmula $\sigma = \sqrt{(1 - p)pn}$.

Por ejemplo, para el experimento de antes (tirar una moneda 20 veces), $n = 20$, y $p = .5$. Entonces usando la fórmula, el promedio será $p \times n = 20 \times .5 = 10$. Muy similar a lo que dio en la simulación (9.96). Y la desviación estándar debería ser la raíz cuadrada de $(0.5 \times 0.5 \times 20)$, es decir, la raíz cuadrada de 5, que es 2.23. Muy similar a lo que dio en la simulación (2.27). Es decir, la

aproximación de Bernoulli va a ser un excelente predictor de la distribución que tendremos.

Ejemplo: el prisionero aburrido

Un prisionero está muy aburrido y tira una moneda 10.000 veces. Supongamos que la moneda es justa. Obviamente el promedio será 5000. También podemos calcular la desviación estándar. Recordemos que $\sigma = \sqrt{(1-p)pn}$. Calculamos $(1-p)pn = \frac{1}{2} \times \frac{1}{2} \times 10.000 = \frac{1}{4} \times 10.000 = 2.500$. La raíz cuadrada de esto es 50. Es decir, si tiro la moneda 10.000 veces, hay aproximadamente 68% de probabilidad de que salga “cara” entre 4950 y 5050 veces. Mientras que hay aproximadamente 95% de probabilidad de que salga “cara” entre 4900 y 5100 veces.

Significancia

A veces se asume que la naturaleza se comporta de forma “normal”, es decir, distribuida en forma de campana. Una actividad donde las distribuciones normales son importantes es el testeo de hipótesis.

Supongamos que quiero testear que una pastilla es buena para el insomnio, ¿Qué hago? Lo que se suele hacer es lo siguiente:

- Tomo una muestra suficientemente grande.
- A la mitad de la muestra les doy la pastilla.
- A la otra mitad les doy un placebo (o no le doy nada).
- Comparo la efectividad de ambos tratamientos.
- Si hay suficiente diferencia positiva, entonces la pastilla es efectiva.

Pero todo esto suena muy impreciso (“suficientemente grande”, “suficiente diferencia”, etc.). La estadística nos puede ayudar a determinar parámetros para esta actividad. Para determinar estos parámetros, hace falta introducir una nueva terminología.

La *hipótesis nula* es la hipótesis de que “no hay un efecto significativo”. Si el efecto detectado es *grande*, entonces existe un efecto significativo, y la hipótesis nula queda “refutada”. Para determinar la significancia del resultado, necesito saber la probabilidad de obtener los datos que obtuve bajo la hipótesis nula

(es decir, si *no* hubiese efecto significativo). Esa probabilidad p es lo que suele llamarse el *nivel de significancia*. A veces se llama *p-valor* (o “p-value”, en inglés). Idealmente quisiera que ese valor sea lo más bajo posible. Por ejemplo, si quiero defender que un tratamiento médico es efectivo, necesito que los resultados obtenidos al usar ese tratamiento sean *muy improbables* de acuerdo con la hipótesis nula.

Para diseñar estas pruebas usamos ciertos criterios: en general, decimos que la evidencia E es *significativa* respecto a la hipótesis nula H cuando la probabilidad de E dado H es menor a 0.05 (o a veces, en el mejor de los casos, menor a 0.01). Por ejemplo, la pastilla es útil cuando suponiendo que el tratamiento no tiene efecto, la evidencia que encontramos “sorprende”. En otras palabras, el resultado *muy probablemente* no podría haberse dado si no se hubieran tomado las pastillas. El resultado puede ser significativo “al nivel del 5%” o incluso “al nivel del 1%”.

Filosóficamente, lo que podemos confirmar en un experimento (donde $p = 0.05$) es esta afirmación condicional: “si la hipótesis nula es cierta, lo que sucedió tiene una probabilidad de menos de 5%”. Esto se interpreta, en un contexto científico, como un rechazo de la hipótesis nula, o una afirmación de que el efecto sí es significativo.

Ejemplo: el prisionero aburrido (otra vez)

El prisionero tira la moneda 10.000 veces, suponiendo que es justa. Como vimos en la sección anterior, calculando la desviación estándar, obtenemos que hay:

- Aproximadamente 68% de probabilidad de que salga “cara” entre 4950 y 5050 veces.
- Aproximadamente 95% de probabilidad de que salga “cara” entre 4900 y 5100 veces.

Ahora supongamos que el prisionero hace la prueba y le sale “cara” 4500 veces. ¿Puede confirmar que la moneda de hecho está viciada?

Filosóficamente, si nos ponemos muy estrictos, solo podemos decir que “si la hipótesis nula fuera cierta, es decir, si la moneda

fuera justa, algo extremadamente raro ($< 1\%$) sucedió”. Científicamente, diríamos que la hipótesis nula queda rechazada (porque $p < 0.01$). Es decir, podemos afirmar que la moneda está viciada.

Los testeos de significancia se usan en muchas, o todas, las áreas de la ciencia empírica. En *medicina*, es común (casi universal) hacer este tipo de testeos. Por ejemplo, para evaluar la efectividad de un tratamiento, se mide la distancia entre el grupo tratado y el grupo control. En *psicología* es también común el uso de estos testeos (dentro de los estudios cuantitativos). Por ejemplo, si queremos saber si el sexo tiene un efecto en la autoestima, debemos comparar las medias en una escala de autoestima en hombres y mujeres, y determinar si la diferencia entre las medias es estadísticamente significativa. Para evaluar esas afirmaciones, en cualquier revista de psicología encontramos tablas y análisis estadísticos de resultados. Otras ciencias sociales, como la sociología, la economía y las ciencias políticas utilizan también esta metodología. Cada disciplina tiene sus propios métodos, que se utilizan para distintos tipos de experimentos. Explorar este tema en detalle requeriría un capítulo aparte.

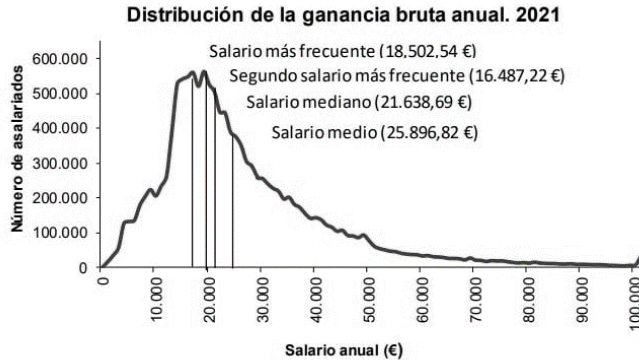
Sin embargo, esto no significa que los científicos conozcan los detalles matemáticos de la estadística. Muchos investigadores (¿la mayoría?) usan *paquetes estadísticos* como SPSS o STATA que hacen las cuentas solas, sin necesariamente comprender qué significa cada resultado. Por ejemplo, en un archivo de SPSS, podría anotar datos de salarios de una gran cantidad de hombres y mujeres. Luego, apretando un botón, podría comparar sus salarios y notar si su diferencia es significativa (es decir, si hay una “brecha de género”). Parte del entrenamiento de un científico, especialmente en medicina y distintas ciencias sociales, es aprender a utilizar esas herramientas computacionales.

Ejercicios

1. En general, cada 51 hombres, nacen 49 mujeres. Supongamos que en un país nacen 800.000 bebés al año.
¿Cuántas mujeres van a nacer el próximo año? Encuentre el promedio, la desviación estándar (redondear en un número entero),

y una estimación usando probabilidades, similar a la del ejemplo del prisionero.

2. La curva de salarios en España se ve así. ¿Es una curva de Gauss?



Soluciones para el capítulo 1

PARTE A

1.
 - a. Tengo $1/6$ de probabilidades de sacar el mismo número en ambos dados.
 - b. Aquí tengo $1/18$ chances de ganar: de entre los 36 posibles resultados, solo gano con (3, 4) y (4, 3).
 - c. Es más probable que sumen 7. Hay 6 posibles escenarios: (1,6), (2,5), (3,4), (4,3), (5,2), (6,1). En cambio, hay solo 5 posibilidades de que sume 6: (1,5), (2,4), (3,3), (4,2), (5,1).
2.
 - a. Con *vocal impar* tengo $2/5 \times 5/9 = 10/45 = 2/9$.
 - b. Con *consonante par* tengo $3/5 \times 4/9 = 12/45$.
 - c. Con alguna de las dos anteriores: $12/45 + 10/45 = 22/45$.
 - d. Me conviene jugar *consonante impar*. De este modo tengo $3/5 \times 5/9 = 15/45 = 1/3$.

3.

Una forma de probarlo es la siguiente:

Supongamos que A es una contradicción.

Entonces $\neg A$ es una tautología. Por Certeza, $P(\neg A) = 1$.

Pero $P(\neg A) = 1 - P(A)$, por Teorema de Negación.

Entonces $1 = 1 - P(A)$. Por lo tanto, $P(A) = 0$.

PARTE B

1.

a. $1/10 \times 1/10 = 1/100$.

b. $1/10 \times 3/39 = 3/390 = 1/130$.

c. Sin reemplazo: $\frac{3}{4} \times \frac{3}{4} = 9/16$.

Con reemplazo: $\frac{3}{4} \times \frac{29}{39} = 87/156 = 29/52$.

d. $P(\text{sacar 7 en la primera mano} \vee \text{sacar 7 en la segunda mano}) = P(\text{sacar 7 en la primera mano}) + P(\text{sacar 7 en la segunda mano}) - P(\text{sacar 7 en ambas manos}) = 1/10 + 1/10 - 1/100 = 0.19$.

Otra forma de calcular sin usar esa fórmula: la probabilidad de que salga algún número no-7 es 0.9, la probabilidad de que pase dos veces es 0.81. Entonces $1 - 0.81 = 0.19$.

2.

a. $P(\text{Palermo y Positivo}) = 1/10 \times 4/10 = 4/100$

b. $P(\text{No Palermo y Positivo}) = 9/10 \times 1/10 = 9/100$

$P(\text{Positivo}) = P(\text{Palermo y Positivo}) + P(\text{No Palermo y Positivo}) = (4 + 9) / 100 = 13/100$

c. $P(\text{Palermo} \mid \text{Positivo}) = P(\text{Palermo y Positivo}) / P(\text{Positivo}) = 4/13$

d. Usando el ejercicio anterior, una vez que le hago la prueba a uno y le da positivo, la probabilidad “previa” de que sean de Palermo es 4/13. Es solo cuestión de repetir los cálculos con nuevos priors. Ahora si le hago la prueba al otro y también da positivo, la probabilidad de que *ambos* sean de Palermo es:

$P(\text{Palermo y Positivo}) = 4/13 \times 4/10 = 16/130$

$P(\text{No Palermo y Positivo}) = 9/13 \times 1/10 = 9/130$

$P(\text{Positivo}) = 25/130$

$$P(\text{Palermo} \mid \text{Positivo}) = P(\text{Palermo} \ \& \ \text{Positivo})/P(\text{Positivo}) = 16/25 = 0.64$$

PARTE C

1.

a.

$$P(D) = 0.1$$

$$P(I|D) = 0.6$$

$$P(I|\neg D) = 0.2$$

$$P(D|I) =$$

$$\frac{P(I \mid D) \times P(D)}{P(I \mid D) \times P(D) + P(I \mid \neg D) \times P(\neg D)}$$

$$= (0.6 \times 0.1) / [(0.6 \times 0.1) + (0.2 \times 0.9)] = 0.06 / 0.24 = 1/4$$

b.

$$P(D) = 0.25 \text{ (nueva probabilidad)}$$

$$P(A|D) = 0.5$$

$$P(A|\neg D) = 0.1$$

$$P(D|A) =$$

$$\frac{P(A \mid D) \times P(D)}{P(A \mid D) \times P(D) + P(A \mid \neg D) \times P(\neg D)}$$

$$= (0.5 \times 0.25) / [(0.5 \times 0.25) + (0.1 \times 0.75)] = 0.125 / [0.125 + 0.075] = 0.125 / 0.2 = 0.625$$

2.

a.

$$P(\text{Negativo} \mid \text{Tengo}) = 0.05$$

$$P(\text{Negativo} \mid \text{No Tengo}) = 0.99$$

$$P(\text{Tengo} \mid \text{Negativo}) =$$

$$\frac{P(\text{Neg} \mid \text{Tengo}) \times P(\text{Tengo})}{P(\text{Neg} \mid \text{Tengo}) \times P(\text{Tengo}) + P(\text{Neg} \mid \text{No Tengo}) \times P(\text{No Tengo})}$$

$$= 0.05 \times 0.03 / (0.05 \times 0.03 + 0.99 \times 0.97)$$

$$= 0.0015 / 0.96015 = 0.001$$

Por esta razón, se considera que un resultado negativo de HIV a los 28 días es definitivo.

b.

$$P(\text{Positivo} \mid \text{Tengo}) = 0.985$$

$$P(\text{Positivo} \mid \text{No Tengo}) = 0.003$$

$$P(\text{Tengo}) = 0.004$$

$$P(\text{Tengo} \mid \text{Positivo}) =$$

$$\frac{P(\text{Pos} \mid \text{Tengo}) \times P(\text{Tengo})}{P(\text{Pos} \mid \text{Tengo}) \times P(\text{Tengo}) + P(\text{Pos} \mid \text{No Tengo}) \times P(\text{No Tengo})}$$

$$= 0.985 \times 0.004 / [0.985 \times 0.004 + 0.003 \times 0.996]$$

$$= 0.0039 / 0.0069 = 39/69 \cong 40/70 \cong 4/7$$

Por esta razón, un resultado positivo de HIV requiere de pruebas adicionales.

PARTE D

1.

Sabemos que $\mu(\{145\}) = 0.3$, $\mu(\{8\}) = 0.2$, $\mu(\{132\}) = 0.5$.

Entonces la nueva probabilidad de que sea el 8 es:

$$P(\{8\} \mid \{8, 132\}) = P(\{8\} \cap \{8, 132\}) / P(\{8, 132\}) = \\ = P(\{8\}) / P(\{8, 132\}) = 0.2 / 0.7 = 2/7$$

2.

La proposición “va a venir un bus de 3 cifras” o $\{145, 132\}$ es estable. En primer lugar, porque $P(\{145, 132\}) = 0.8$. Además su probabilidad sigue siendo > 0.5 bajo condicionalización con cualquier proposición compatible. Condicionada bajo W , su probabilidad sigue siendo 0.8. Condicionada bajo $\{132\}$ o $\{145\}$, su probabilidad será 1. Condicionada bajo $\{132, 8\}$ su probabilidad será $P(\{132\}) / P(\{132, 8\}) = 5/7$. Condicionada bajo $\{145, 8\}$, su probabilidad será $P(\{145\}) / P(\{145, 8\}) = 3/5$.

PARTE E

1. No, no puede suceder. Si $P(A)$ es mayor a 0.5, $P(A \vee B)$ también lo es, por el Teorema de Validez. Entonces $P(\neg(A \vee B))$ es menor a 0.5

2. Uso la fórmula de Adams: $P(A) + P(B) - 1 = 1.8 - 1 = 0.8$. Entonces debo creer C con probabilidad mayor o igual a 0.8. No puedo creer racionalmente C en probabilidad 0.6.

3. A es “va a salir 1, 2, 3 o 4”, B es “va a salir 3, 4, 5 o 6”. Ambos son probables en 0.66. Pero $A \& B$ es “va a salir 3 o 4” y esto solo es 0.33 probable. Entonces $\neg(A \& B)$ es probable en 0.66 también.

PARTE G

1.

Supongamos que $P(A) + P(\neg A) > 1$

[Creo demasiado en A y $\neg A$]

Digamos que $P(A) = p$ y $P(\neg A) = q$, y que $p + q > 1$.

Te hago apostar esto: p sobre A , y q sobre $\neg A$.

Entonces te ofrezco lo siguiente:

	Resultado por apuesta sobre A	Resultado por apuesta sobre $\neg A$
Sucede A	$1 - p$	$-q$
Sucede $\neg A$	$-p$	$1 - q$

Resultado: si sucede A , obtengo $(1 - p - q) = 1 - (p + q)$, que debe ser negativo. Lo mismo obtengo si sucede $\neg A$.

2.

a. Si llueve: obtengo $0.16 + 0.36 = 0.52$.

Si no llueve: obtengo $0.36 + 0.16 = 0.52$.

b. Si llueve: obtengo $0.16 + 0.16 = 0.32$.

Si no llueve: obtengo $0.36 + 0.36 = 0.72$.

Esta distribución es probabilística pero no domina a la anterior.

c. El puntaje es 0.5 en cada mundo. Entonces domina a la distribución del punto *a*.

d. El teorema solo nos dice que una asignación de probabilidad incoherente tiene una asignación probabilística que lo domina. En este caso, la tercera es coherente y domina a la primera. No estamos diciendo que *todas* las asignaciones coherentes dominan a todas las incoherentes; de hecho, la segunda no domina a la primera.

PARTE H

1.

$$\sigma = \sqrt{0.51 \times 0.49 \times 800000} \cong 447$$

$$\mu = 392.000$$

Es decir, hay 68% de probabilidades de que sean entre 391.553 y 392.447.

Hay 95% de probabilidades de que sean entre 391.106 y 392.894.

2. No es una curva normal porque no es simétrica.

CAPÍTULO 2: TEORÍA DE LA DECISIÓN

En este capítulo, abordaremos algunos conceptos de la *teoría de la decisión*. El propósito de esta teoría es comprender la racionalidad de las acciones humanas. Hay al menos dos enfoques sobre la teoría de la decisión: en el enfoque *deliberativo*, el objetivo de la teoría es ayudarnos a tomar decisiones, mientras que en el enfoque *explicativo*, la teoría explica por qué decidimos lo que decidimos. Más adelante veremos que ambos enfoques tienen sus ventajas y desventajas.

El punto de partida de la teoría de la decisión, sin embargo, es normativo. Es decir, la teoría de la decisión no es (al menos en principio) una teoría sobre cómo las personas deciden, o de los procesos psicológicos o neurológicos involucrados en la toma de decisiones. La teoría de la decisión nos muestra cómo las personas deberían decidir racionalmente, a partir de determinada reconstrucción de los contextos de decisión.

Parte A: Matrices, actos y resultados

En la teoría de la decisión, siguiendo el enfoque clásico de Savage (1954), un contexto de decisión incluye *actos* y *estados del mundo*. Por ejemplo, yo puedo decidir si salir con paraguas o sin paraguas. Los posibles estados del mundo relevantes son si llueve o no llueve. Usualmente escribimos estas situaciones en una tabla, donde las filas son los actos y las columnas son los estados del mundo.

	Llueve	No llueve
Salgo con paraguas		
Salgo sin paraguas		

Como puede verse, esa tabla está sin completar. Lo que falta es describir los posibles *resultados*, es decir, qué pasaría si tomamos determinada decisión en determinado estado del mundo. Por ejemplo, la tabla anterior podría completarse así:

	Llueve	No llueve
Salgo con paraguas	No me mojo	Debo cargar el paraguas innecesariamente
Salgo sin paraguas	Me termino mojando	No me mojo y no tengo que cargar un paraguas

Este procedimiento se ve muy sencillo, y de hecho no tiene por qué ser distinto a lo que nos representamos mentalmente al tomar una decisión. Típicamente consideramos los posibles escenarios antes de actuar. Las condiciones respecto a los actos y los estados del mundo no son muchas. Entre ellas:

1. Los estados del mundo deben ser excluyentes y exhaustivos. Es decir, siempre va a suceder uno de ellos, y no pueden suceder varios a la vez.
2. Los actos deben ser independientes del estado del mundo. Es decir, tomar un curso de acción no debería afectar la probabilidad de que el estado del mundo ocurra.

Esta última condición fue discutida en distintos contextos, y de hecho hay teorías (como la Teoría de la Decisión Causal) donde finalmente no se la adopta. Sin embargo, la condición es muy útil para representar decisiones cotidianas. En el contexto antes descrito, se trata de una condición natural: llevar el paraguas no hace más o menos probable que llueva.¹³

¹³ En el habla popular, existen las llamadas “Leyes de Murphy”. Por ejemplo: “el carril de al lado siempre es más rápido”, “cuando abras el paraguas dejará de llover”, “apenas enciendas el cigarro llegará el bus”, etc. Estas “leyes” solo tienen un valor humorístico, no científico.

Ejercicio

Julio César tiene que ir a la guerra. Debe decidir cuántos soldados enviar (en particular, si enviar pocos o muchos). No sabe si se enfrentará a un ejército grande o pequeño. Pero realiza el siguiente razonamiento:

	Gano la guerra	Pierdo la guerra
Envío muchos soldados	Victoria	Pierdo muchos soldados
Envío pocos soldados	Victoria épica	Pierdo pocos soldados

Al comprender que, en cualquier caso, es mejor mandar pocos soldados (un razonamiento que luego llamaremos *dominancia*), Julio César decide enviar pocos soldados.

¿Por qué es incorrecto su razonamiento?

¿Cómo debería formular esta matriz de decisión para evitar esta falacia?

Parte B: Decisiones bajo ignorancia

Teniendo tablas como la de antes, nos alcanza para tomar *algunas* decisiones. Pero necesitamos un poco más de información. Como mínimo, necesitamos saber cuáles son nuestros resultados preferidos. Caso contrario, no podríamos decidir qué rumbo de acción tomar.

La teoría de la decisión es fuertemente subjetivista respecto a las preferencias. Es decir, la teoría de la decisión no puede decirnos cuáles preferencias son “mejores” o “peores”. Sobre eso, cada uno tendrá su opinión. Lo único que impone la teoría de la decisión es cierta estructura racional de las preferencias, que nos permitirá luego tomar una decisión racional. En otras palabras, la teoría de la decisión nunca nos dirá “deberías hacer esto”, sino “dadas tus preferencias, deberías hacer esto”. Algunos autores llaman a esto “racionalidad instrumental”.

El elemento básico de esa estructura racional es una relación de *preferencia* entre resultados. Es decir, necesitamos un tipo de orden que nos indica qué resultados preferimos sobre otros. Usaremos el símbolo $x \succsim y$ para representar que el agente encuentra al resultado x mejor o igual que el resultado y . Cuando $x \succsim y$ ocurre junto con $y \succsim x$, decimos que $x \sim y$ (Indiferencia). Y cuando $x \succsim y$ pero $\neg(y \succsim x)$, decimos que $x \succ y$ (Preferencia estricta). Estos órdenes de preferencia tienen las siguientes propiedades:

Complejitud: Para todos los eventos x e y , se da que $x \succ y$, o $y \succ x$, o $x \sim y$. Es decir, los agentes siempre prefieren un resultado sobre otro, o son indiferentes.

Transitividad: Si $x \succsim y$, y también $y \succsim z$, entonces $x \succsim z$.

La *complejitud* establece que todos los agentes tienen preferencias respecto a todos los posibles resultados. Si voy al supermercado a buscar una fruta y hay solo una manzana y una banana, puedo preferir la banana, o la manzana, o puede darme igual entre ambas cosas. No hay resultados “incomparables”.

A la vez, esas preferencias deben ser *transitivas*. Si de hecho prefiero una banana antes que una manzana, y una manzana antes que un durazno, entonces obviamente prefiero una banana antes que un durazno. A modo de metáfora, podríamos pensar un orden como un edificio con pisos, donde algunas opciones (las preferidas) están arriba, otras más abajo, y las menos preferidas en la planta baja. Aquí un “piso” puede contener muchas opciones, es decir que permitimos la indiferencia entre opciones.

Es común defender la transitividad a partir de razones de explotabilidad, similares al argumento del Dutch Book. El argumento de la bomba de dinero (*money-pump argument*) nos dice que, si nuestras preferencias son cíclicas, un corredor nos puede dejar sin dinero (Gustaffson 2022). Porque si preferimos A sobre B, podríamos gastar una mínima suma de dinero en obtener A, cuando tenemos B. Y si las preferencias tuvieran un ciclo, no importa en qué parte del ciclo comenzamos, podríamos gastar sumas de dinero para obtener el evento preferible, hasta quedarnos sin nada.

En un clásico texto, Amartya Sen (1971) muestra que la transitividad es equivalente a dos condiciones sobre las preferencias: la condición α y la condición β . La condición α nos dice que, si entre un conjunto A nuestra opción preferida es x , entonces si le quitamos elementos a ese conjunto A , no puede cambiar nuestra opción preferida. Es decir, si entre {banana, pera, manzana} yo prefiero la pera, no puede suceder que entre {banana, pera} yo prefiera la banana. La condición β nos dice que, si tanto x como z son mis opciones preferidas dentro de A , y agregamos elementos a ese conjunto, no puede ser que x siga siendo una opción preferida, pero z no lo sea (o viceversa). Es decir, si entre {chocolate, vainilla, frutilla} yo prefiero vainilla y chocolate, no puede ser que entre {chocolate, vainilla, frutilla, limón} yo prefiera solamente vainilla, o solamente chocolate.

¿De qué forma podría combinar estos órdenes completos y transitivos de preferencias para tomar decisiones? Supongamos que estoy decidiendo si ir a almorzar al comedor universitario. No tengo muy claro qué comida van a servir hoy, pero sé que será barato. Mi única alternativa es almorzar en el restaurante tailandés de la esquina, que es delicioso pero caro.

Podría representar la situación del siguiente modo:

	El comedor sirve comida rica	El comedor sirve comida mediocre
Almuerzo en el comedor universitario	Gasto poco dinero y como algo rico	Gasto poco dinero y como algo mediocre
Almuerzo en restaurante tailandés	Gasto mucho dinero y como algo rico	Gasto mucho dinero y como algo rico

El resultado de almorzar en el restaurante de la esquina no depende de lo que sirva el comedor universitario, por eso será el mismo en ambos estados del mundo. Ahora supongamos que lo que más me preocupa es el dinero. En ese caso, mis preferencias se verán de este modo (donde arriba pongo lo que más prefiero):

Gasto poco dinero y como algo rico
Gasto poco dinero y como algo mediocre
Gasto mucho dinero y como algo rico

Vale recordar que en este contexto las preferencias son *subjetivas*. Nadie está diciendo que gastar mucho dinero para comer algo rico sea malo, simplemente estamos representando las preferencias de un sujeto determinado.

Desde la perspectiva de ese sujeto, si voy al comedor universitario, me garantizo gastar poco dinero, que es lo que más me importa. En otras palabras, sea cual sea el estado del mundo (sea cual sea la comida del comedor), voy a preferir comer allí antes que ir al restaurante de la esquina. Según un conocido principio de la racionalidad, en estos casos lo que corresponde hacer es ir al comedor universitario. Se trata de una decisión *dominante*.

(Dominancia estricta) Un acto X *domina estrictamente* a un acto Y si y sólo si en cualquier estado del mundo, el acto X trae un resultado estrictamente preferible al que trae el acto Y.

También diremos, en este caso, que la decisión Y está *estrictamente dominada*. Un acto que domina estrictamente a todos los otros actos se llamará *dominante*.

En ocasiones también se usa el concepto de *dominancia débil*, que es muy similar al anterior (aunque menos exigente):¹⁴

(Dominancia débil) Un acto X *domina débilmente* a un acto Y si y sólo si (a) en cualquier estado del mundo, el acto X trae un resultado igual o mejor que el acto Y, y además (b) en algún estado del mundo, el acto X trae un resultado estrictamente mejor que el acto Y.

¹⁴ En algunas presentaciones, lo que aquí llamamos “dominancia débil” es presentado como “dominancia estricta” (Peterson 2009).

En teoría de la decisión, suele asumirse que los agentes no deberían tomar decisiones dominadas.

(Regla de dominancia) La *regla de dominancia* establece que los actos dominados son irracionales, es decir, debemos elegir acciones que no estén dominadas.

En el ejemplo en cuestión, dadas nuestras preferencias puramente económicas, el acto de ir al comedor universitario es *dominante*, porque domina a la alternativa, que es ir al restaurante tailandés. Entonces, según la teoría de la decisión, deberíamos ir al comedor universitario.

Tomando en cuenta que la teoría de la decisión dialoga con la filosofía, es natural que todas las reglas de la racionalidad estén bajo discusión. Sin embargo, la regla de dominancia es una de las más sólidas en cualquier teoría de la racionalidad. Si quiero actuar racionalmente, lo mínimo que debo hacer es descartar las opciones dominadas.

Decisiones sin dominancia

En muchas ocasiones, la dominancia no nos ayuda a tomar una decisión. Supongamos que tengo que decidir si ir a la playa o al cine, pero las consecuencias dependen mayormente del clima. Puedo modelar la matriz de decisión de este modo:

	Llueve	No llueve
Ir a la playa	Día perdido	Mucha felicidad
Ir al cine	Felicidad moderada	Arrepentimiento leve

Y supongamos que mi orden de preferencias es el siguiente:

Mucha felicidad
Felicidad moderada
Arrepentimiento leve
Día perdido

Aquí la regla de dominancia no sirve, porque ir a la playa solo es preferible a ir al cine cuando no llueve; análogamente, ir al cine solo es preferible a ir a la playa cuando llueve. Ningún acto domina al otro.

Una regla posible para decidir es *evitar el peor resultado*. En ese caso, el peor resultado es ir a la playa y que llueva (“Día perdido”). Usando esa regla, lo que debemos hacer es ir al cine. Llamamos a esta regla *maximin*.

(Maximin) Entre varios actos, debo elegir aquel cuyo peor resultado sea el menos malo (“maximizar el mínimo”).

La regla *maximin* es la más conocida para estos contextos de poca información. Como veremos más adelante, esta regla caracteriza una forma extrema de la aversión al riesgo, y aparece en distintas versiones a lo largo de la literatura de la teoría de la decisión. De todos modos, vale repetir que no se trata de una regla canónica de la racionalidad, como la regla de dominancia. *Maximin* es simplemente una regla posible para tomar decisiones en contextos de escasez informativa.

Junto con *maximin*, existen otras reglas posibles para estos contextos sin mucha información. Por ejemplo, un agente más “optimista” podría elegir el acto donde el mejor resultado sea el mejor de todos. En el ejemplo, el agente iría a la playa, esperando que salga el sol. A esta regla la llamamos *maximax*. Existen muchas otras reglas posibles en estos escenarios donde solo tenemos un orden de preferencias, aunque aquí no entraremos en detalle sobre este tema.

Ejercicio

Tengo una cita mañana y planeo adónde llevar a mi pareja. Se me ocurre ir al Parque Rivadavia, quedarnos en mi casa viendo Netflix o ir a pasear por Recoleta (un barrio con parques y museos). Pero no sé si va a llover o no. Puedo representar mis preferencias de este modo:

	Soleado	Llueve
Parque Rivadavia	a	b
Netflix en casa	c	d
Recoleta	e	f

El orden de preferencias es $a > d > e \sim c > f > b$.

- ¿Hay alguna acción estrictamente (o débilmente) dominada?
- ¿Qué acción nos indicaría la regla Maximax?
- ¿Qué acción nos indicaría la regla Maximin?

Parte C: Escalas de utilidad

En contextos específicos, podemos tener escalas más informativas entre los posibles resultados, que vayan más allá de un simple orden. Por ejemplo, podríamos querer informar no sólo que preferimos tener una casa en las afueras a tener un departamento en la ciudad, sino en qué medida nos parece mejor (podría ser doblemente mejor, marginalmente mejor, etc.). Para hacer esto, lo usual es atribuir números a los posibles resultados, y generar una escala donde podamos calcular distancias.

Usando las nociones de la sección anterior: habrá una función de utilidad U que le asigna un número a cada resultado A/E, donde A es un acto, y E es un estado del mundo. Por ejemplo, si A es ir al cine, y E es que llueve, A/E es la situación donde llueve y vamos al cine, y $U(A/E)$ será la utilidad asignada a ese escenario. La idea de asignar utilidades a eventos proviene del *utilitarismo*. Para el utilitarismo hedonista de Bentham (1780), que luego fue

parcialmente adoptado por Mill (1861)¹⁵, podemos caracterizar un evento a partir de cuánto placer o dolor nos da. Éticamente, debemos maximizar el placer de la mayor cantidad posible de personas. Para eso, es necesario usar alguna escala que nos sirva de “hedonómetro”, es decir, para medir el placer y el dolor de determinado evento.

La teoría de la decisión reformula estas ideas utilitaristas, pero las vuelve más manejables técnicamente. La utilidad, según la perspectiva usual en teoría de la decisión, ya no indica nuestros “puntos de placer”, sino que da información sobre la intensidad de las preferencias. Por ejemplo, podríamos ordenar los resultados del ejemplo anterior del siguiente modo:

Ir a la playa un día soleado	100
Ir al cine un día lluvioso	40
Ir al cine un día soleado	20
Ir a la playa un día lluvioso	0

Aquí se mantiene el orden, pero damos más información sobre cuánto preferimos la plaza un día soleado antes que ir al cine. Otro agente podría ordenar sus preferencias así:

Ir a la playa un día soleado	100
Ir al cine un día lluvioso	2
Ir al cine un día soleado	1
Ir a la playa un día lluvioso	0

Este agente tiene el mismo *orden* de preferencias que el anterior, pero no valora mucho ir al cine. Solo lo considera muy levemente mejor que ir a la playa un día lluvioso. Lo verdaderamente deseable para este agente es ir a la playa en un día soleado.

Como vemos, en estas escalas “cardinales” no hay solo un orden, sino también distancias. Ahora bien, si lo importante es la *dis-*

¹⁵ El utilitarismo de Mill, a diferencia del de Bentham, suponía una jerarquía de placeres, donde el placer de una noche de cervezas está por debajo de placeres superiores, como una noche de ópera.

tancia relativa entre los eventos, y no el número preciso que usamos, muchas escalas van a darnos lo mismo. Por ejemplo, tener preferencias del tipo $(100, 10, 0)$ podría ser equivalente a tener las preferencias $(10, 1, 0)$. Estas equivalencias entre escalas de preferencias pueden explicarse con el concepto de *transformación*. En general, dos escalas son iguales en términos prácticos cuando una es simplemente una transformación de la otra.

Una *transformación de cociente* toma los elementos del orden y los multiplica por un número positivo a , es decir, $f(x) = ax$. Por ejemplo, podríamos tomar el orden $(100, 1, 0)$ y multiplicarlo por 2. Nos quedaría $(200, 2, 0)$. Aquí, no solo se mantienen las distancias relativas, sino también el cociente entre los valores: 100 es 100 veces 1, 200 es 100 veces 2. En términos coloquiales, esto significa que “achicamos” o “agrandamos” la escala.¹⁶

Algo más compleja, pero fundamental en la teoría de la decisión, es la noción de *transformación lineal positiva*. Una transformación lineal positiva de un orden (x_1, \dots, x_n) es una multiplicación de cada elemento de ese orden por una función $f(x) = ax + c$. En términos coloquiales, esto significa que además de “achicar” o “agrandar” la escala (con el elemento a), también podemos “correr” la escala (con el elemento c). Una consecuencia de esto es que las escalas transformables linealmente no tienen inicios, finales o puntos medios “naturales”.

Las transformaciones lineales mantienen las distancias relativas, aunque no los cocientes. Por ejemplo, $(101, 2, 0)$ es una transformación lineal positiva de $(100, 1, 0)$, donde corremos el eje solamente un punto (es decir, $f(x) = x + 1$). Pero el cociente no se mantiene, porque 101 no es 100 veces 2. Un caso conocido de transformación lineal es la conversión de grados Fahrenheit (que se usan en Estados Unidos) a grados Celsius (que se usan en casi todo el resto del mundo). Si tengo un número de grados Celsius x , puedo hacer la conversión usando la fórmula $f(x) = 9/5 x + 32$. Por ejemplo, 0 grados Celsius son 32 grados Fahrenheit. Se trata

¹⁶ Además de mantener el cociente, las escalas de cociente (a diferencia de las de intervalo) pueden tener un 0 único. En un libro reciente, Narens y Skyrms (2020) utilizan una escala de cociente, donde el 0 separa las sensaciones de placer de las sensaciones de dolor.

de dos formas de medir la temperatura, que no modifican aquello que miden: obviamente en Estados Unidos la temperatura no es un fenómeno natural distinto a la temperatura en Argentina.

Ejercicio

Tengo cuatro escalas distintas, que reflejan preferencias de cuatro agentes:

	Agente 1	Agente 2	Agente 3	Agente 4
Frutillas	6	30	112	5
Naranjas	3	15	106	10
Almendras	2	10	104	2
Chocolate	1	5	102	1

- Encuentre dos escalas que *no* son equivalentes en términos ordinales.
- Encuentre dos escalas equivalentes en términos de cociente (y la función de transformación correspondiente).
- Encuentre dos escalas equivalentes en términos de intervalo que no son equivalentes en términos de cociente (y la función de transformación correspondiente).

Parte D: Maximización de utilidad esperada

Para tomar decisiones informadas no alcanza con tener escalas más informativas de preferencias. Además, necesitamos un cálculo más preciso sobre los posibles estados del mundo. Para este propósito usaremos *probabilidades*. La teoría de probabilidades ya fue explicada en el capítulo anterior. Aquí, lo único que nos interesa es que un agente tiene una distribución de probabilidades sobre los posibles estados del mundo. Es decir, cada estado E tiene una probabilidad $P(E)$, que es la probabilidad de que

ese estado del mundo ocurra, según la perspectiva del agente. Al mismo tiempo, como se trata de una distribución, la suma de las probabilidades de los estados del mundo (que son exhaustivos y excluyentes) será exactamente 1. De este modo, podemos dar más precisiones sobre la forma en que la teoría de la decisión entiende los actos.

Cada acto A tendrá una utilidad esperada $U(A)$. Esta utilidad es simplemente la ponderación, teniendo en cuenta las probabilidades, de todos los posibles resultados. Es decir, si los estados del mundo son E_1, \dots, E_n , la *utilidad esperada del acto A* se puede calcular así:

$$U(A) = P(E_1) \times U(A/E_1) + \dots + P(E_n) \times U(A/E_n)$$

Por ejemplo, supongamos que la probabilidad de que llueva es 60%, o 0.6, y que nuestra escala de utilidad es la que definimos antes:

	Llueve ($p = 0.6$)	No llueve ($p = 0.4$)
Ir al cine	40	20
Ir a la playa	0	100

Entonces la utilidad de ir al cine será:

$$U(\text{Ir al cine}) = 0.6 \times 40 + 0.4 \times 20 = 24 + 8 = 32$$

Por otro lado, la utilidad de ir a la playa es:

$$U(\text{Ir a la playa}) = 0.6 \times 0 + 0.4 \times 100 = 0 + 40 = 40$$

Es decir, en este escenario, la utilidad de ir a la playa es 40, y la utilidad de ir al cine es 32. La regla de racionalidad más utilizada en la teoría de la decisión nos indica que debemos realizar el acto que nos otorgue la *máxima utilidad esperada*. En este caso, debemos ir a la playa.

Maximización de utilidad esperada: Un agente racional debe realizar el acto que maximiza la utilidad esperada.

La maximización de utilidad esperada es la regla fundamental de la teoría de la decisión, cuando tenemos probabilidades y utilidades. Naturalmente, esta regla implica la regla de dominancia (dejo la prueba a cargo del lector), aunque no necesariamente implica otras reglas que hemos mencionado. Por ejemplo, la regla de maximizar la utilidad esperada no implica *maximin*: si el peor escenario fuera muy improbable, no tenemos por qué focalizarnos en ese escenario al tomar una decisión (en el escenario anterior, *maximin* nos indicaría ir al cine).

Un resultado interesante nos indica que, al tomar una decisión por maximización, la transformación lineal positiva uniforme de las utilidades nos va a arrojar el mismo resultado¹⁷. Es decir, la decisión correcta no depende de los números, sino de las distancias relativas entre los posibles resultados.

Dinero y utilidad

¿Cuál es la relación entre la *utilidad* y los bienes materiales? En teoría de la decisión, generalmente asumimos que el dinero no tiene un valor lineal (y lo mismo aplica, en general, a otros bienes materiales). Es decir, la función de utilidad no es de tal forma que $U(\$x) = x$. Tampoco puede representarse con otra función lineal, como $U(\$x) = 5x + 10$. Decimos que el dinero tiene una *utilidad marginal decreciente*: mientras más dinero tenemos, menos nos significa ganar un poco más.

En otras palabras, si alguien gana 100 dólares mensuales, ganar 120 seguramente le significará una diferencia muy sustantiva. En cambio, si alguien gana 12.000 dólares mensuales, ganar 12.020 no le va a significar una utilidad mucho mayor.

En realidad, la teoría marginal del dinero proviene de una famosa paradoja de la teoría de la decisión, conocida como *Paradoja de*

¹⁷ Esto se debe a que si $x > y$, $ax + c > ay + c$ (asumiendo que a es positivo). Es decir, multiplicar por un número positivo, y sumar/restar lo mismo a ambos lados, no puede cambiar el signo de la inecuación.

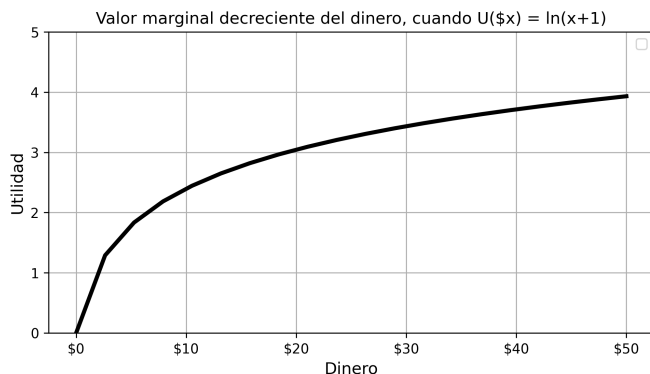
San Petersburgo. La paradoja se basa en un juego: se tira una moneda, y si sale “cara” el agente gana \$2. En cambio, si sale “cruz”, se vuelve a tirar la moneda. Ahora bien, si en el segundo tiro sale “cara”, el agente gana \$4. Si sale “cruz”, se vuelve a tirar. En el tercer tiro, si sale “cara”, el agente gana \$8. Si sale “cruz”, vuelve a tirar. Y así sucesivamente. Es decir, si la moneda sale “cara” en el turno n , el agente obtiene $\$2^n$. Una pregunta que podríamos hacernos ahora es cuánto dinero es razonable pagar por entrar en ese juego. Lo cierto es que el valor monetario esperado de la apuesta es el siguiente:

$$\frac{1}{2} \times \$2 + \frac{1}{4} \times \$4 + \frac{1}{8} \times \$8 + \dots = \$1 + \$1 + \$1 + \dots = \$\infty$$

En otras palabras, el valor monetario de esta apuesta es *infinito*. Guiado por el valor monetario, yo debería pagar todo lo que tengo para jugar al juego de San Petersburgo. Obviamente, esto es absurdo. Por eso la conclusión de esta paradoja es que el valor de una apuesta *no* es su valor monetario. De hecho, el valor de una apuesta suele ser mucho *menor* a su valor monetario. Esto se debe, justamente, a la utilidad marginal decreciente del dinero. Según Daniel Bernoulli (1700-1782)¹⁸, la utilidad del dinero es una función $U(\$x)$ que usualmente se comporta de forma logarítmica. Para simplificar, podríamos decir que $U(\$x) = \ln(x + 1)$. Es decir, la utilidad de \$5 será 1.79, y la de \$10 será 2.39. Pero la utilidad de \$40 será 3.71. Como vemos, multiplicar el dinero no multiplica la utilidad al mismo ritmo. Siempre tener más dinero aumenta la utilidad, pero el aumento no es lineal.¹⁹

¹⁸ La paradoja como tal fue descubierta por Nickolaus Bernoulli (1687-1759), primo de Daniel Bernoulli, y sobrino de Jakob Bernoulli, a quien ya mencionamos como descubridor de la Ley de Grandes Números.

¹⁹ Podríamos ver a la utilidad marginal decreciente del dinero como una ley entre otras que relacionan lo objetivo (cantidad de dinero) con lo subjetivo (cantidad de utilidad) usando funciones logarítmicas. Otra medida de este tipo son los decibeles.



En otros contextos, llamamos *adversas al riesgo* a las personas cuya utilidad es tal que $U(\$x) < x$. Suele asumirse que todos los agentes son, en mayor o menor medida, adversos al riesgo. Pensemos en apuestas como la siguiente:

Acción 1: Me dan mil dólares.

Acción 2: Tiro una moneda. Si sale cara, me quedo sin nada. Si sale cruz, me dan dos mil dólares.

El valor monetario esperado (VM) de la acción 2 es igual al valor monetario esperado de la acción 1, pues $VM(\text{Tirar moneda}) = 0.5 \times \$0 + 0.5 \times \$2000 = \1000 . Sin embargo, en encuestas y contextos experimentales, las personas en su inmensa mayoría prefieren la acción 1 sobre la acción 2. Una forma de explicar este fenómeno es justamente apelar a la noción de aversión al riesgo: la utilidad de la acción 1 es $U(\$1000)$, pero la utilidad de la acción 2 es $0.5 \times U(\$0) + 0.5 \times U(\$2000)$. Y como antes señalamos, dos mil dólares no representan el doble de utilidad que mil dólares.

Vale aclarar, sin embargo, que la utilidad marginal decreciente del dinero no es un principio esencial de la teoría de la decisión. Simplemente es una premisa que suele sostenerse en la mayoría de los contextos. Pero podría haber otros contextos donde las personas, en lugar de ser adversas al riesgo, sean propensas al riesgo.

Ejercicios

1. Debo ir a un congreso en un distrito que queda a 800km, pero anuncian mal clima para esa semana. Me debato sobre qué hacer.

	Día normal ($p = 0.4$)	Lluvias ($p = 0.5$)	Tormenta ($p = 0.1$)
Ir en avión	40	20	-20
Ir en bus	20	40	20
No ir	0	10	150

¿Qué recomendaría la maximización de utilidad esperada?

2. Tengo que elegir entre dos apuestas, ambas tirando un dado.

A1: Si sale 1 o 2, gano \$121.

A2: Si sale 1, 2 o 3, gano \$49.

a. ¿Qué debo elegir si la utilidad del dinero es lineal, es decir, si $U(\$x) = x$?

b. ¿Qué debo elegir si $U(\$x) = \ln(x + 1)$?

c. ¿Qué debo elegir si $U(\$x) = \sqrt{x}$?

3. Llegan los postres y hay muchos chocolates. Me gusta el chocolate, pero tengo problemas estomacales, de modo tal que no debería comer demasiados (más de 8 o 9 ya me traen problemas). El anfitrión me ofrece decidir entre dos acciones: comer una caja grande (9 chocolates) o una pequeña (6 chocolates).

Personalmente, prefiero los chocolates blancos.

Ahora bien: hay 0.6 probabilidades de que me traiga chocolates negros, y 0.4 de que traiga chocolates blancos.

Siendo x la cantidad de chocolates que como, supongo que la utilidad de los blancos es $[10 - (x - 7)^2]$, y la de los negros es $[8 - (x - 7)^2]$. ¿Debo elegir la caja grande o la pequeña?

Parte E: El Teorema von Neumann-Morgenstern

Como señalamos anteriormente, la noción contemporánea de *utilidad* proviene del concepto desarrollado por utilitaristas como Bentham, pero el enfoque es radicalmente distinto. Von Neumann y Morgenstern, en su influyente libro *Theory of Games and Economic Behavior* (1944), mostraron que uno puede *construir* una escala de preferencias cardinal al observar las preferencias de los agentes sobre determinadas “loterías”.

Una lotería, en este sentido técnico, es una situación donde podrían suceder distintos hechos con distintas probabilidades, y todas ellas suman 1. Una lotería podría ser tirar una moneda (50% cara, 50% cruz) o también salir sin paraguas un día nublado (40% llevo mojado al trabajo, 60% llevo seco al trabajo). Mostraremos ahora, de forma intuitiva, cómo construir una escala de preferencias de intervalo a partir de las preferencias sobre loterías.

Supongamos que un agente debe ir al cine a ver una película de Pixar. En principio, sus preferencias son *Toy Story* \succ *Moana* \succ *Cars*. Para comprender la utilidad que le asigna a cada escenario, le ofrecemos las siguientes loterías:

- Ver *Moana* con 100% de probabilidad
- Ver *Toy Story* con 80% de probabilidad, y *Cars* con 20%

Supongamos que ambas loterías le son indiferentes al agente. Con esto, sin suponer ninguna utilidad en particular, podemos armar una suerte de escala:

$$1 \times U(\textit{Moana}) = 0.8 \times U(\textit{Toy Story}) + 0.2 \times U(\textit{Cars})$$

Ahora estipulamos que el mejor resultado vale 100, y el peor vale 0 (esto podemos estipularlo con cualquier número, porque solo nos importan los intervalos). Como señalamos al principio, el mejor resultado para el agente es *Toy Story*, y el peor es *Cars*. Entonces podemos calcular la utilidad de ver *Moana*:

$$U(\textit{Moana}) = 0.8 \times 100 + 0.2 \times 0 = 80$$

La escala entonces queda así:

Toy Story	100
Moana	80
Cars	0

Debería ser claro que no importan los valores que ponga, en tanto y en cuanto use escalas que sean *equivalentes bajo transformación lineal positiva*, es decir, en tanto se mantengan las *distancias* entre los puntos. La escala aquí es (100, 80, 0), pero podría ser (10, 8, 0), o (11, 9, 1), etc.

Veamos ahora cómo podemos armar una escala incluso más informativa. Supongamos que al agente en cuestión le ofrecen ver *Coco*. *Coco* le gusta más que *Moana*, pero menos que *Toy Story*. Pero eso no alcanza para armar una escala de intervalo. Para saber el lugar preciso de *Coco* entre *Moana* y *Toy Story* necesito saber esto: ¿con qué valor de p le daría igual ver *Coco* seguro, o ver *Toy Story* con probabilidad p y *Moana* con probabilidad $(1-p)$? Supongamos que ese valor es 0.7. Es decir:

$$U(Coco) = 0.7 \times U(Toy\ Story) + 0.3 \times U(Moana)$$

Ahora podemos usar la escala de utilidad anterior:

$$U(Coco) = 0.7 \times 100 + 0.3 \times 80 = 70 + 24 = 94$$

Entonces la nueva escala tiene a *Coco* con 94.

Así, puedo establecer la escala de utilidad de un agente a partir de sus preferencias respecto a un conjunto determinado de apuestas. Esta es la esencia conceptual del teorema de Von Neumann y Morgenstern. Ahora veremos el teorema con más detalle.

***Enunciado del Teorema von Neumann-Morgenstern**

En la versión completa del problema, el agente elige entre loterías. Usaremos esta terminología: “ ApB ” es una lotería donde sale A con probabilidad p , y B con probabilidad $(1-p)$.

Hay cuatro axiomas para las preferencias:²⁰

- **Complejitud** (véase sección B de este capítulo)
- **Transitividad** (véase sección B de este capítulo)
- **Independencia**: $A \succ B$ si y sólo si $ApC \succ BpC$.
- **Continuidad**: Si $A \succ B \succ C$, entonces hay probabilidades p y q (mayores a 0 y menores a 1) tales que $ApC \succ B$ y $B \succ AqC$.

El axioma de *Independencia* quiere decir, a grandes rasgos, que agregar opciones irrelevantes no debería cambiar mi preferencia sobre las loterías. Por ejemplo: si prefiero *Cars* a *Pocahontas*, prefiero *Cars* con 80% y *Aladdin* con 20%, a *Pocahontas* con 80% y *Aladdin* con 20%.

El axioma de *Independencia* nos prohíbe elegir sobre loterías en virtud de propiedades *globales*. Por ejemplo, si me fijo en propiedades globales, podría preferir una lotería entre películas del mismo estilo; entonces una lotería entre *Pocahontas* y *Aladdin* (dos películas 2D) sería mejor que una lotería entre *Cars* y *Aladdin* (*Cars* es 3D y *Aladdin* es 2D). Esta consideración de propiedades globales de las loterías violaría *Independencia*.

El axioma de *Continuidad* es más controversial, y una forma de leerlo es que nuestros órdenes de preferencias deben ser “continuos”, sin saltos. Podemos ignorar escenarios negativos si su probabilidad es suficientemente baja, y no podemos ignorarlos si su probabilidad es suficientemente alta.

Un ejemplo del axioma de *Continuidad* es lo siguiente. Si A es obtener \$10.000, B es obtener \$9.000 y C es no obtener nada, obviamente prefiero $A \succ B \succ C$. Ahora veamos cómo opera la continuidad. Supongamos que tengo que elegir entre tirar una moneda y obtener \$10.000 (si sale “cara”) y \$0 (si sale “cruz”), u obtener \$9.000 en mano:

Lotería 1: \$9.000 (B) seguro.

Lotería 2: 50% de chances de \$10.000 (A), 50% de nada (C).

²⁰ También se requiere que las loterías compuestas puedan ser reducidas a loterías simples. Por razones de complejidad, no entramos en detalle sobre este asunto aquí.

La mayor parte de la gente preferirá la lotería 1. Eso podría ilustrar un q donde $B \succ AqC$, siendo $q = 0.5$. Pero ahí la probabilidad de ganar el máximo premio es 0.5. ¿Qué pasaría si la probabilidad de sacar \$10.000 fuera más alta? Lo que nos dice el axioma de Continuidad es que hay una probabilidad $p < 1$ (supongamos, 0.99) según la cual prefiero \$10.000 con probabilidad p antes que \$9.000 seguro. En ese caso, $ApC \succ B$.

Una forma equivalente de formular el axioma de Continuidad es decir que, si $A \succ B \succ C$, existe una probabilidad p tal que $ApC \sim B$. Este principio, como antes vimos, es fundamental para “calibrar” la escala de utilidad.

El principio de Continuidad contradice al “efecto de certeza” del que hablaron psicólogos como Kahneman, quien afirma que una probabilidad de 100% (es decir, algo seguro) tiene un efecto que va más allá del cálculo de las probabilidades. También la Continuidad se opone a la noción general de *maximin*: por más malo que sea el peor resultado, podemos anularlo a fines prácticos cuando la probabilidad es suficientemente baja. Esto no significa que el axioma de Continuidad sea poco realista; por el contrario, nos seguimos yendo de vacaciones a la playa aun cuando sabemos que hay una chance (mínima) de tener un accidente en el camino. Sería irracional dejarnos llevar por el peor resultado posible, que es lo que propone *maximin*.

Ahora podemos mencionar el enunciado del teorema de representación de von Neumann-Morgenstern:

Teorema von Neumann-Morgenstern: Una relación \succ de preferencia sobre loterías satisface los cuatro axiomas (Complejitud, Transitividad, Independencia y Continuidad) si y solo si existe una función u que va de loterías a números reales en $[0,1]$, y que tiene las siguientes propiedades:

1. Una lotería preferible tendrá mayor utilidad:
 $L1 \succ L2$ sii $u(L1) > u(L2)$
2. La utilidad de una lotería se calcula a partir de su utilidad esperada: $u(ApB) = p \times u(A) + (1 - p) \times u(B)$

Asimismo, esa función u será única bajo transformación lineal positiva (i.e., toda función que satisface 1 y 2 es una transformación lineal positiva de u). La prueba no la haremos aquí, porque excede la complejidad de este libro (véase Peterson 2009, Apéndice B). Pero podríamos resumir el resultado de este modo: si un conjunto de preferencias de un agente sobre loterías o actos satisface los axiomas antes mencionados, el agente puede ser representado como maximizador de utilidad.

Ejercicios

1. Supongamos que prefiero helado de chocolate sobre helado de frutilla, y helado de frutilla sobre helado de banana. Pero soy indiferente entre comer helado de frutilla seguro, y una apuesta con 60% de helado de chocolate y 40% helado de banana. ¿Cuál es mi escala de intervalo si la utilidad de comer helado de chocolate es 70 y la de helado de banana es 10?
2. Supongamos que un agente es extremadamente averso al riesgo. El valor de una lotería es el valor del peor resultado posible (sin importar su probabilidad). Probar que viola el axioma de Independencia.
3. Supongamos que un agente tiene una escala *lexicográfica* de preferencias, donde prefiere un celular iPhone a un Android, sin importar el modelo. Asimismo, el agente prefiere un iPhone seguro (por más malo que sea) a cualquier lotería donde sea probable tener un Android. Probar que viola Continuidad.

Parte F: Utilidad conductual y utilidad sustantiva

El teorema Von Neumann-Morgenstern estableció una nueva ortodoxia en la teoría de la decisión, según la cual podemos extraer las utilidades a partir de las preferencias reveladas en la acción. Esto llevó a una aceptación casi universal entre economistas de que la teoría de la decisión tiene un poder esencialmente *explicativo*: observando la conducta ajena, podemos inferir usando la

teoría de la decisión cuáles son las preferencias o creencias de los agentes, y luego predecir su conducta futura. Una versión de este enfoque más común en las ciencias económicas es la Teoría de la Preferencia Revelada, generalmente atribuida a Paul Samuelson (1938). Según Samuelson, la preferencia se revela en la acción, y no hace falta apelar a una teoría psicológica de la utilidad para entender la conducta racional.

Suele llamarse a esta lectura *interpretación conductual* de la teoría de la decisión. La interpretación conductual tiene muchas ventajas: entre otras cosas, para entender las creencias o deseos de los otros, no nos importa tanto lo que las personas *dicen* que quieren o que piensan, sino lo que realmente *hacen*. De ahí podemos inferir sus preferencias. Esto evita tendencias “paternalistas”, donde el teórico decide qué es lo que los agentes prefieren “realmente” sin tomar en cuenta lo que hacen. El enfoque conductual también permite predecir en vista de las elecciones pasadas, en base a la consistencia de las preferencias. Por último, el aparato matemático de la teoría de la decisión me permite inferir patrones de decisión más “finos” que los que pueden expresarse verbalmente (Thoma 2021).

El enfoque conductual recibió numerosas críticas, aunque no todas apuntan a los mismos problemas. Un problema de la teoría conductual (Ostapiuk 2022) es que reduce la irracionalidad a casos de violación de los axiomas (por ejemplo, elegir una manzana entre {manzana, pera} pero una pera entre {manzana, pera, banana}). Sin embargo, “racionaliza” casos cotidianamente considerados irracionales, como una persona que desea dejar de fumar, pero sigue fumando. Dentro del enfoque conductual, si el agente sigue fumando es porque evidentemente prefiere el placer a corto plazo, y actúa racionalmente. Becker y Murphy (1988, p. 675) sostienen que “las adicciones, incluso las más fuertes, son racionales”. Otro problema es que, en ocasiones específicas, las elecciones no parecen revelar una preferencia real. En la paradoja del *asno de Buridan*, un asno debe elegir si tomar una fruta que cuelga a la derecha o la izquierda, pero es indiferente entre ambas. En la paradoja original, el asno (sin poder decidir racionalmente) muere de hambre. Pero supongamos que no es tan

tonto, y toma una de las frutas. Aunque es indiferente entre ambas, cualquiera sea la fruta que tome, su “preferencia revelada” dirá que prefiere una sobre la otra (Sen 1973, p. 248). Por último, la teoría conductual no puede explicar conductas donde violamos los axiomas, como cuando decidimos en base a propiedades “globales” (Dietrich & List 2016).

Una interpretación alternativa de la teoría de la decisión, mucho más común entre filósofos que entre economistas, es la interpretación *sustantiva* (a veces llamada “mentalista” o “realista”). Esta lectura proviene del utilitarismo de Bentham, y del enfoque original de Bernoulli. Para Bentham, como antes mencionamos, la “utilidad” de una acción puede medirse según la cantidad de placer o felicidad que provoca. La versión más usual de la teoría sustantiva sigue siendo subjetivista, y nos dice que la utilidad es un estado mental correlacionado pero no idéntico a las decisiones que tomamos. Por ejemplo, Hausman (2012) propone que la utilidad debería definirse a partir de nuestros juicios subjetivos sobre qué es lo mejor para nosotros. Entonces, maximizar la utilidad consiste en maximizar lo que nos parece mejor (o lo que nos da más placer, felicidad, etc.).

Bermúdez (2009) propone que la teoría sustantiva tiene algunas ventajas explicativas. En primer lugar (p. 49), la teoría nos permite dar una recomendación normativa más clara: debes maximizar tu utilidad. Mientras que la teoría conductual solo prescribe mantener la consistencia con conductas pasadas; pero este criterio no podríamos aplicarlo en situaciones nuevas. En segundo lugar (p. 50), la teoría sustantiva es más compatible con los cambios de preferencias: lo que en un momento nos da placer o felicidad, en otro momento podría dejar de hacerlo.

Por último, el enfoque sustantivo nos permite recoger casos intuitivos de irracionalidad, como el fumador compulsivo que mencioné anteriormente: podríamos decir que el fumador compulsivo es irracional porque no lleva a cabo lo que él mismo considera conveniente (un fenómeno conocido como *akrasia*).

La teoría sustantiva es más común en textos de filosofía, donde se postula a un agente con la capacidad de “entender” sus propias preferencias, y de decidir qué hacer a partir de eso. Pettigrew (2019, p. 15), por ejemplo, adopta la lectura sustantiva porque le

interesa “dar una teoría de la decisión que podríamos usar realmente para deliberar sobre decisiones que tomamos, y que clarifica las razones que motivan y justifican las decisiones que tomamos como resultado de la deliberación”. De hecho, los filósofos suelen adoptar la interpretación sustantiva sin mayores aclaraciones (Okasha 2016).

Algunos autores intentaron reconciliar la preferencia revelada del enfoque conductual con la preferencia “genuina” de las teorías sustantivas. Gauthier (1986, p. 28) sostiene que podemos hacerlo al apelar a la preferencia *considerada*. Una preferencia es “considerada” cuando el agente, además de actuar de acuerdo con ella, también la acepta reflexivamente (esto excluye conductas compulsivas); lo racional, según Gauthier, sería seguir nuestras preferencias consideradas. Otros autores, como Hausman y McPherson (2009), sostuvieron que la conducta no puede identificarse con la utilidad (las personas a veces actúan en contra de su beneficio), pero es la mejor evidencia de lo que cada uno quiere; por lo tanto, el enfoque conductual no es metodológicamente incorrecto para sus aplicaciones usuales.

Parte G: La Paradoja de Allais

En el año 1953, el economista Maurice Allais elaboró una objeción contra la idea de racionalidad como maximización de utilidad esperada. Quizás hasta el día de hoy, la “Paradoja de Allais” sea la mayor paradoja contra la teoría estándar de la decisión.

Lo que hoy conocemos como “Paradoja de Allais” se extrae de un conjunto de charlas y artículos de Allais a principios de los años 50. En uno de esos textos, se incluye una encuesta con muchas preguntas. Allais publicó las preguntas, pero no el resultado de las encuestas en el texto original (Allais 1953*b*). De hecho, lo que se discute actualmente respecto al artículo de Allais es solo un conjunto pequeño de preguntas de ese cuestionario.

Uno debe elegir entre distintas loterías (en el sentido usual de la palabra), cada una de ellas con 100 tickets, que te dan distintos premios según el ticket que salga.

En la primera pregunta, te dan a elegir entre dos loterías A1 y A2. La lotería A1 te da 1 millón en todos los casos, mientras que

la lotería A2 te da 5 millones si salen los tickets del 2 al 11, y 1 millón si salen los tickets del 12 al 100:²¹

	1	2-11	12-100
A1	1 millón	1 millón	1 millón
A2	0	5 millones	1 millón

Mayormente la gente elige A1 sobre A2, porque te da 1 millón de dólares seguro. Esto, en principio, no contradice la maximización de utilidad esperada; por el contrario, solo muestra que la utilidad del dinero es marginalmente decreciente.

En la segunda pregunta, tenemos que elegir entre estas dos loterías:

	1	2-11	12-100
A3	1 millón	1 millón	0
A4	0	5 millones	0

Aquí también, la gente mayormente elige la lotería 4 sobre la 3. Básicamente porque tiene casi las mismas chances de hacerse millonario, pero con el quintuple de dinero. Esta decisión, tomada por sí sola, tampoco viola la teoría estándar.

Sin embargo... *¡tener esas preferencias en conjunto es inconsistente!*

Es decir, no importa cuál sea la utilidad del dinero, uno no puede preferir A1 sobre A2 y A4 sobre A3.

Veamos cómo se calculan las utilidades de A1 y A2.

$$U(A1) = 1/100 \times U(\$1M) + 10/100 \times U(\$1M) + 89/100 \times U(\$1M)$$

²¹ Suponemos que los valores están en dólares. En ciertas monedas un millón podría valer muy poco, y cambiaría el sentido del problema.

$$U(A_2) = 1/100 \times U(\$0) + 10/100 \times U(\$5M) + 89/100 \times U(\$1M)$$

Si $A_1 > A_2$ eso significa que:

$$\begin{aligned} 1/100 \times U(\$1M) + 10/100 \times U(\$1M) + 89/100 \times U(\$1M) > \\ 1/100 \times U(\$0) + 10/100 \times U(\$5M) + 89/100 \times U(\$1M) \end{aligned}$$

Tachando idénticos a ambos lados, obtenemos que:

$$\begin{aligned} (*) \quad 1/100 \times U(\$1M) + 1/10 \times U(\$1M) > \\ 1/100 \times U(\$0) + 1/10 \times U(\$5M) \end{aligned}$$

Pero esto nos impone cierto orden respecto a A_3 y A_4 , pues:

$$\begin{aligned} U(A_3) &= 1/100 \times U(\$1M) + 1/10 \times U(\$1M) + 89/100 \times U(\$0) \\ U(A_4) &= 1/100 \times U(\$0) + 1/10 \times U(\$5M) + 89/100 \times U(\$0) \end{aligned}$$

El tercer término podemos ignorarlo si comparamos A_3 con A_4 . Pero entonces la inequación (*) resultará en que $A_3 > A_4$. Esto contradice lo que elige la mayoría de la gente.

En resumen, lo que establece la paradoja de Allais es que la mayoría de la gente toma decisiones que contradicen el principio de maximización de utilidad. Esto es independiente de la escala de utilidad que tenga el agente.

Una forma de leer la paradoja es como una refutación del principio de Independencia. Si ignoramos los casos de 12-100, ambos escenarios son idénticos (un millón seguro vs. 10/11 de probabilidad de 5 millones). Pero las personas dan respuestas opuestas cuando agregamos un evento adicional (12-100); esto se debe, intuitivamente, a la existencia de un “resultado seguro” en A_1 , que no está presente en A_3 .

Allais pretendía criticar tajantemente la teoría de la decisión racional de Von Neumann, Morgenstern y otros. Para Allais, esta paradoja demuestra que el enfoque “americano” (es decir, el estándar) ignora lo que significa ser un agente *racional*. Entre otras cosas, porque un agente racional o prudente no va a satisfacer axiomas como *Continuidad*: para un agente racional, es mejor algo seguro que la probabilidad de algo mejor, sin importar cuál

sea esa probabilidad. Allais (1953a, p. 504) observa la importancia de la *dispersión* de las posibilidades (por ejemplo, la cantidad de posibilidades que aparecen con un determinado acto), algo que la teoría estándar de la decisión no puede capturar.

Para economistas como Savage, representantes de la escuela americana, la paradoja muestra que el principio de maximización de utilidad no *describe* el comportamiento, sino que establece cómo un individuo racional *debería* comportarse. En ese sentido, que mucha gente haya fallado en el cuestionario muestra lo difícil que es comportarse racionalmente. Esta distinción entre lo normativo y lo descriptivo estará de fondo en la disciplina de la economía del comportamiento, que en aquellos momentos recién estaba naciendo. De hecho, es interesante mencionar que Kahneman y Tversky (1979), los más célebres impulsores de la economía del comportamiento, comprobaron empíricamente las ideas de Allais (para ciertos valores de las apuestas).

En resumen, podríamos decir que Allais *intentó* refutar los principios de la escuela “americana”, pero sus argumentos no fueron convincentes para la mayoría de los economistas. La escuela clásica o americana terminó volviéndose hegemónica, y la paradoja sirvió para distinguir entre una lectura descriptiva (pero difícil de defender) y una lectura normativa de la teoría de la decisión, mucho más aceptada.

Ejercicio

Si la utilidad del dinero fuera $U(x) = \ln(x+1)$, ¿qué elegiría un agente racional en el escenario planteado por Allais?

Parte H: La Paradoja de Ellsberg

Así como la Paradoja de Allais se cuestiona el axioma de Independencia, la Paradoja de Ellsberg (Ellsberg, 1961) nos indica que la teoría de la decisión estándar no es capaz de comprender la incertidumbre sobre probabilidades.

En la paradoja, tengo una urna, donde sé que hay 30 bolas rojas, y otras 60 que algunas son amarillas y otras negras.

$$R = 30$$
$$A \vee N = 60$$

Debo decidir entre dos apuestas (con valores en dólares):

$$A1: \$100 \text{ si sale una bola roja}$$
$$A2: \$100 \text{ si sale una bola amarilla}$$

En general, la gente va a elegir A1, porque seguro la probabilidad es $1/3$ (y no menos).

Después, deberías elegir entre estas dos apuestas:

$$A3: \$100 \text{ si sale una roja o negra}$$
$$A4: \$100 \text{ si sale una amarilla o negra}$$

Aquí, por el mismo criterio, la gente va a elegir A4, porque seguro la probabilidad será $2/3$ (y no menos).

Pero *este conjunto de decisiones es irracional*. Porque, supongamos que a es el número de bolas amarillas en la urna, y que u es la utilidad de \$100. Entonces:

$$U(A1) = 1/3 u$$
$$U(A2) = a/90 u$$

Supongamos que $A1 > A2$.

Entonces $1/3 u > a/90 u$, y esto implica que $a < 30$.

Ahora, sea n la cantidad de bolas negras en la urna.

$$U(A3) = (n/90) u + 1/3 u$$
$$U(A4) = 2/3 u$$

Por el resultado anterior ($a < 30$), n debe ser mayor a 30, y por eso $U(A3)$ será mayor a $2/3 u$. Por lo tanto, debo elegir A3 sobre A4. Lo contrario a lo que hacen las personas.

Según Ellsberg, este caso muestra que la incerteza no siempre puede ser representada como una probabilidad (especialmente en casos de escasez informativa, que él llama “ambigüedad”). La

paradoja fue confirmada por algunos psicólogos experimentales, aunque las decisiones van a depender fuertemente de las probabilidades asignadas (MacCrimmon & Larsson 1978). Luego del desafío planteado por Ellsberg, distintos autores han intentado desarrollar versiones más sofisticadas de la teoría de la decisión, que sean capaces de representar este tipo de “incerteza”.

Teoría del riesgo epistémico

Según la teoría del *riesgo epistémico*, podemos distinguir dos tipos de incertidumbre: cuando conocemos las probabilidades y cuando solo conocemos las *posibilidades*. Por ejemplo, si estamos en un partido de tenis, no es lo mismo *saber* que dos jugadores son igual de buenos, que no tener idea de quién es ninguno de ellos. En ambos casos nos negaríamos a apostar a favor de algún jugador, pero nuestras razones son distintas.

La idea de este enfoque desarrollado en los 70' (Levi 1974), pero muy popular en las últimas décadas, es que podemos representar la meta-incertidumbre usando “rangos” de probabilidades. Por ejemplo, si mi conocimiento sobre los jugadores de tenis es *nulo*, la probabilidad de que gane el Jugador 1 podría ser desde 0 hasta 1. Pero tal vez yo sé que ese jugador es “bueno” y el otro “no es tan bueno”, entonces la probabilidad de que gane el Jugador 1 podría ser un rango entre 0.6 a 1. En cualquier caso, mi estado epistémico queda definido no por una asignación de probabilidades, sino por un *conjunto* de asignaciones de probabilidad.

La pregunta entonces es: ¿cómo debería comportarme en estos casos? La respuesta de Levi es, esencialmente, que tenemos que guiarnos por la aversión al riesgo. Para ejemplificar, podemos distinguir entre tres situaciones epistémicas:

Situación 1: Sé que los dos jugadores son igual de buenos.

Situación 2: Sé que uno es mucho mejor que el otro (0.9 probabilidades de ganar) pero no sé cuál es cuál.

Situación 3: No tengo la más pálida idea de lo que estoy viendo.

La apuesta es la siguiente: Si gana el jugador que elijo, obtengo una utilidad de 2; si pierde el jugador que elijo, pierdo una utilidad de 1. No apostar me da una utilidad de 0.

Entonces, ¿en qué escenarios debería tomar la apuesta?

En la situación 1, la utilidad de la apuesta es $0.5 \times 2 + 0.5 \times (-1) = 1 - 0.5 = 0.5$. Entonces me conviene tomarla. Hasta aquí usamos simplemente la teoría estándar.

En la situación 2, tengo dos asignaciones posibles de probabilidad en mente: que gane el jugador A con .9, o que gane el jugador B con .9. Es decir, podría haber dos casos. Considero solamente el peor:

Caso en que elijo al equivocado: $0.1 \times 2 + 0.9 \times (-1) = .2 - 0.9 = -0.7$

Entonces, *usando el criterio de maximin para utilidades esperadas*, en esta situación no me conviene tomar la apuesta.

En la situación 3, hay infinitos casos, así que simplemente debo usar *maximin* usual. De este modo, me conviene no apostar.

Es decir, la meta-incertidumbre me lleva a actuar racionalmente con una máxima general de aversión al riesgo. Este método puede explicar el comportamiento usual en la paradoja de Ellsberg, donde los agentes eligen la aversión al riesgo en cada caso por separado, incluso cuando es inconsistente con la teoría de la decisión estándar.

Ejercicio

Estoy cuidando dos pacientes gemelos. A veces les doy aspirina para el dolor. Sea cual sea el paciente, la utilidad de calmarle el dolor es 10, la utilidad de que le caiga mal la aspirina es -10, y la utilidad de que siga con dolor (sin tomar la aspirina) es 0. El tema es que a un paciente le cae mal 60% de las veces, y al otro solo 1%. Ahora viene un paciente con dolor y no sé si darle la aspirina o no. Y yo no sé precisamente cuál paciente es el que me está pidiendo la aspirina, dado que son gemelos. ¿Qué debería hacer según la teoría de riesgo epistémico?

Parte I: Psicología de la decisión

Kahneman y Tversky (1979) revolucionaron el mundo de la psicología y la teoría de la decisión, mostrando resultados muy similares a los de Allais. A diferencia de Allais, Kahneman y Tversky tenían formación metodológica como psicólogos, y pudieron demostrar los resultados con mayor rigurosidad. Además, elaboraron algunas hipótesis que explican por qué y cuándo las personas se alejan del principio de maximización de la utilidad.

Efecto de certeza

Las personas deben elegir entre estas dos loterías (en dólares):

\$3.000 seguro	vs	\$4.000 con 80%
\$3.000 con 25%	vs	\$4.000 con 20%

En el experimento de Kahneman y Tversky, la mayoría (80%) prefirieron \$3.000 seguro en la primera, y la mayoría (65%) eligieron \$4.000 con 20% en la segunda ($n = 95$). Esto contradice el principio de maximización de utilidad esperada.

Siguiendo la teoría estándar, los agentes deberían elegir el mismo lado en ambos casos. Esto se debe a que la utilidad de una lotería que me da \$3.000 con 25% de probabilidad es naturalmente $\frac{1}{4}$ de la utilidad de obtener \$3.000 seguro; y la utilidad de obtener \$4.000 con 20% de probabilidad es $\frac{1}{4}$ de la utilidad de \$4.000 con 80%. Es decir, se trata de la misma apuesta, pero las utilidades de la segunda apuesta están multiplicadas por 0.25. Esto no debería cambiar el orden de preferencias.

La idea de Kahneman y Tversky es que aquí opera un *efecto de certeza*. No importa la utilidad del dinero, \$3.000 con certeza es mejor que una *posibilidad* con el mismo valor esperado. Esto no opera en la segunda apuesta porque no hay certeza, es solo una posibilidad contra otra posibilidad.

Probabilidades pequeñas

Otro fenómeno que detectaron Kahneman y Tversky es este:

\$6.000 con 45%	vs	\$3.000 con 90%
-----------------	----	-----------------

\$6.000 con 0.01% vs \$3.000 con 0.02%

En el experimento que realizaron los autores, la mayoría (86%) elige \$3.000 con 90% antes que \$6.000 con 45%. Mientras que la mayoría (73%) prefiere \$6.000 con 0.01% antes que \$3.000 con 0.02% ($n = 66$). Esto viola el principio de maximización de utilidad, por las mismas razones que en la apuesta anterior: las probabilidades de arriba y las de abajo tienen la misma proporción. Cualquiera sea la utilidad de \$6.000 o de \$3.000, la preferencia en estas apuestas debe ser igual.

Según Kahneman y Tversky, lo que está operando aquí es la incapacidad para operar con probabilidades pequeñas. En particular, aquí los agentes perciben que 90% es el doble que 45% pero no que 0.02% es el doble que 0.01%, dado que ambos números son pequeños. Esto también había sido notado por Allais (1953b, p. 54), que planteó que las probabilidades objetivas se “distorcionan” subjetivamente.

Efecto reflejo

La última encuesta de Kahneman y Tversky que tomaremos es la siguiente. En realidad, es un par de decisiones que deben tomar los agentes:

1. Te dieron \$1000. Ahora debes elegir entre \$500 extra seguros, o tirar una moneda entre \$1000 o nada.
2. Te dieron \$2000. Ahora debes elegir entre que te saquen \$500 seguro, o tirar una moneda entre que te saquen \$1000 o nada.

Si leemos con atención, veremos que ambas apuestas son *exactamente iguales* (elegimos entre \$1500 seguros, y tirar la moneda entre \$1000 y \$2000); solamente cambia la forma de describirlas. Aquí ni siquiera opera algún tipo de transformación de utilidades o probabilidades. La gente (84%), sin embargo, tiende a elegir los \$500 extra seguros en el primer caso ($n = 70$), pero también eligen (69%) tirar la moneda en el segundo caso ($n = 68$).

Lo que indica esto para Kahneman & Tversky es una asimetría entre las conductas respecto a pérdidas y ganancias. En particular, la gente prefiere irse con \$1500 en el primer caso, que es lo seguro, y además “gana” algo respecto al *statu quo* de \$1000. Mientras que, en el segundo caso, la gente piensa que lo mejor sería arriesgarse a quedarse con los \$2000 (es decir, mantener el *statu quo*), porque perder \$500 es visto como una *pérdida*.

En términos un poco más técnicos, a partir de cierto punto percibimos una diferencia entre “pérdida” y “ganancia”. Respecto a las ganancias, nos comportamos de forma más aversa al riesgo; pero respecto a las pérdidas, preferimos tomar riesgos. Por eso este efecto se llama *efecto reflejo* (*reflection effect*). Este efecto es una versión más general de un fenómeno conocido como “efecto del marco” (*framing effect*), donde las decisiones se toman según la forma en que está presentado el problema.

Teoría de prospectos

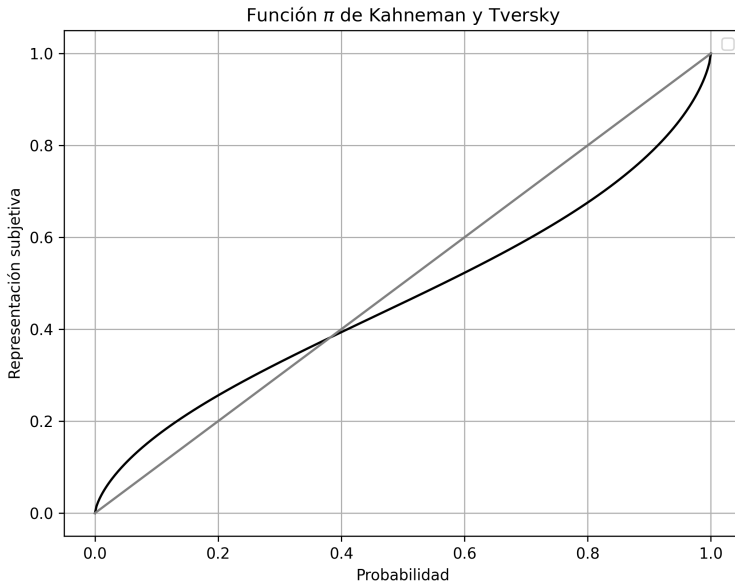
Usando algunas de estas consideraciones, Kahneman & Tversky plantearon cierta modificación de la teoría de la decisión, donde la utilidad de un acto se define así:

$$U(A) = \pi(p_1) \times v(u_1) + \dots + \pi(p_n) \times v(u_n)$$

Ahora veremos precisamente qué significan esa v y esa π .

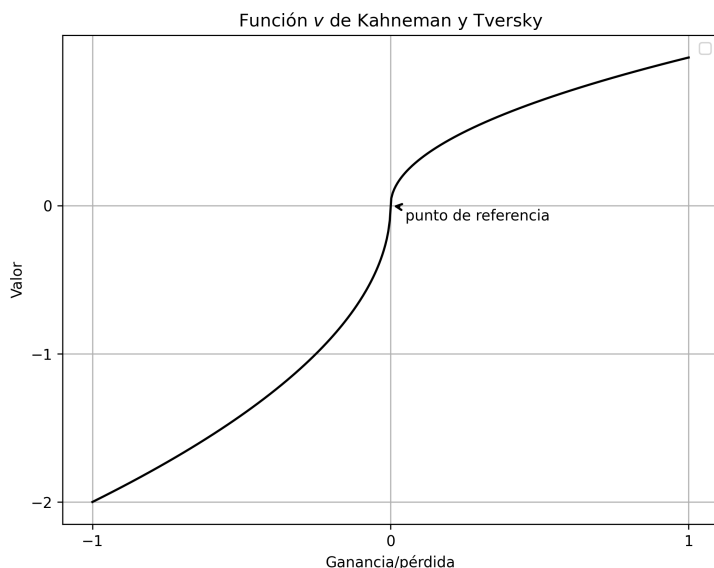
Función π sobre probabilidades

La función π sobre probabilidades no es exactamente lineal, porque *sobreestima* las probabilidades pequeñas y *subestima* las grandes. La idea es que, si escuchamos que algo es 99% probable, exageramos la importancia de ese 1% faltante, aunque sea realmente bajo. Subjetivamente, un 1% significa más que solo un 1% (por estas razones, una hamburguesa siempre cuesta \$5.99, y no \$6). La curva en este gráfico representa nuestra representación subjetiva de las probabilidades:



Función v sobre utilidades

Respecto a la función v , lo que hace es mostrar el valor subjetivo de determinadas ganancias o pérdidas (relativas a un punto de referencia). La idea, brevemente, es que en un marco de decisión el agente encuentra un punto de referencia a partir del cual ganar cuenta como “ganancia”, y perder cuenta como “pérdida”. Solemos pensar que ese punto es \$0, pero esto depende del escenario descrito. Según estos autores, las pérdidas importan más (proporcionalmente) que similares ganancias. Y por otro lado, los agentes tienden a adoptar conductas más riesgosas (*risk-seeking*) cuando tienen la pérdida asegurada; mientras que la conducta de aversión al riesgo es más común cuando hay alguna ganancia segura. Esto nos arroja una particular función en forma de S:



La función v no solo explica el efecto reflejo, sino también el efecto de certeza. Porque la certeza nos da un punto de referencia, a partir del cual las pérdidas son vistas como mucho peores (proporcionalmente) a las ganancias. En cambio, cuando el problema es planteado en términos de loterías con distintas probabilidades, las personas no tienen ese punto de referencia.

En resumen, para Kahneman y Tversky los agentes violan sistemáticamente los axiomas de la teoría de la decisión estándar, pero podemos buscar una modificación no demasiado sustantiva de la teoría original que describe mejor el comportamiento de los agentes.

Ejercicio

En la película *Don't Look Up* (2021), los astrofísicos encarnados por Leonardo di Caprio y Jennifer Lawrence le plantean a la presidenta (Meryl Streep) que el mundo se va a acabar por la caída de un meteorito. La presidenta les critica que en ciencia nada es

exacto, y ellos admiten que la probabilidad es 99.78%. El asistente de la presidenta se alegra de que “no es 100%”, y la presidenta dice “Ok, entonces digamos que es 70% y sigamos”. ¿Cómo explicamos esto usando la teoría de prospectos?

Parte J: Buchak y el riesgo como factor

En su libro *Risk and Rationality* (2013), Lara Buchak propone una teoría alternativa de la decisión que puede explicar mejor algunas situaciones específicas, como las presentadas por Allais. Buchak quiere explicar un tipo de aversión al riesgo distinta a la planteada por la teoría estándar, y más independiente de la utilidad marginal decreciente del dinero.

En la teoría de Buchak, un agente tiene un índice r , que indica su aversión al riesgo. El índice r va de probabilidades a números entre 0 y 1 (podríamos decir que la función r “distorsiona” las probabilidades).

Una función de riesgo $r: [0,1] \rightarrow [0,1]$, tiene las siguientes propiedades:

- $r(0) = 0$
- $r(1) = 1$
- $0 \leq r(x) \leq 1$
- r no es decreciente (si $x \leq y$, entonces $r(x) \leq r(y)$).

Por ejemplo, puedo representar la máxima aversión al riesgo o *maximin* si asumo que $r(1) = 1$, y para todos los otros x , $r(x) = 0$. Otra asignación posible es $r(x) = x^2$.

Ahora supongamos que voy al cine a ver una película de un director que últimamente ha hecho películas relativamente buenas, otras excelentes, y también algunas malas:

	Buena ($p = 0.5$)	Excelente ($p = 0.3$)	Mala ($p = 0.2$)
Cine	11	15	5

En la teoría estándar, calculamos la utilidad del acto de este modo:

$$U(\text{Cine}) = 0.5 \times 11 + 0.3 \times 15 + 0.2 \times 5 = 5.5 + 4.5 + 1 = 11$$

La versión de Buchak funciona distinto. Primero, ordeno los resultados de peor a mejor: $5 < 11 < 15$. Luego calculo las utilidades de forma “acumulativa”.

La idea de un cálculo acumulativo es que lo peor es “seguro”. Tengo cierta probabilidad de obtener lo segundo peor; y cierta probabilidad de obtener lo tercero peor (que si hay tres opciones, es lo mejor); etc. Entonces, si hay tres opciones, calculamos la utilidad de este modo:

$$U(A) = U(\text{peor}) + P(2^\circ \text{ peor o mejores}) \times [U(2^\circ \text{ peor}) - U(\text{peor})] + P(\text{mejor}) \times [U(\text{mejor}) - U(2^\circ \text{ peor})]$$

Más generalmente, cuando hay n opciones:

$$U(A) = U(\text{peor}) + P(2^\circ \text{ peor o mejores}) \times [U(2^\circ \text{ peor}) - U(\text{peor})] + \dots + P(\text{mejor}) \times [U(\text{mejor}) - U((n-1)^\circ \text{ peor})]$$

Para el ejemplo en cuestión, podemos calcularlo así:

$$\begin{aligned} U(\text{Cine}) &= 5 + 0.8 \times (11 - 5) + 0.3 \times (15 - 11) \\ &= 5 + 0.8 \times 6 + 0.3 \times 4 \\ &= 5 + 4.8 + 1.2 \\ &= 11 \end{aligned}$$

Como podemos ver, esta lectura “acumulativa” nos da el mismo resultado que la teoría estándar; solo es otra forma de plantear el mismo concepto.

Pero en la teoría de Buchak, podemos aplicarle el índice r a las probabilidades, y así recogemos la idea de aversión al riesgo como elemento autónomo. La fórmula resultante es esta, que es igual a la anterior, pero aplicamos la r a cada aparición de las probabilidades (Buchak 2013, p. 53):

$$U(A) = U(\text{peor}) + r(P(2^\circ \text{ peor o mejores})) \times [U(2^\circ \text{ peor}) - U(\text{peor})] + \dots + r(P(\text{mejor})) \times [U(\text{mejor}) - U((n-1)^\circ \text{ peor})]$$

Por ejemplo, supongamos que $r(x) = x^2$. En este caso, nuestra utilidad para el ejemplo anterior será:

$$\begin{aligned} U(\text{Cine}) &= 5 + r(0.8) \times (11 - 5) + r(0.3) \times (15 - 11) \\ &= 5 + r(0.8) \times 6 + r(0.3) \times 4 \\ &= 5 + 0.64 \times 6 + 0.09 \times 4 \\ &= 5 + 3.84 + 0.36 \\ &= 9.2 \end{aligned}$$

Es decir, la utilidad del acto es menor, porque cuenta el riesgo. Una ventaja fundamental de la teoría de Buchak es que nos permite explicar la paradoja de Allais sin recurrir a la irracionalidad. Por ejemplo, recordemos esta apuesta:

	1	2-11	12-100
A1	1 millón	1 millón	1 millón
A2	0	5 millones	1 millón

Aquí, la mayoría de la gente prefiere A1 sobre A2, porque A2 es *seguro*. Pero las personas prefieren A4 sobre A3 aquí:

	1	2-11	12-100
A3	1 millón	1 millón	0
A4	0	5 millones	0

La teoría estándar de la decisión es incapaz de explicar esta situación. Pero la teoría de Buchak lo explica fácilmente. Supongamos que uno valora al dinero de forma lineal (es decir, $U(\$x)$

= x), pero tiene cierta aversión al riesgo, representable con $r(x) = x^3$. Entonces:

$$U(A1) = 1M.$$

$$U(A2) = 0 + (0.99)^3 \times 1M + (0.1)^3 \times 4M = 970.000 + 4.000 = 974.000$$

$$U(A3) = 0 + (0.11)^3 \times 1M = 1331$$

$$U(A4) = 0 + (0.1)^3 \times 5M = 5000$$

Aquí podemos ver claramente que $A1 > A2$ y $A4 > A3$, solo haciendo uso de la función de riesgo r .

Ejercicio

Debes elegir cuánto tiempo ir de vacaciones a la playa. Tienes dos opciones, ir una semana o ir dos semanas.

	Lindo Clima ($p = 0.5$)	Mal clima ($p = 0.4$)	Clima normal ($p = 0.1$)
Una semana	7	4	6
Dos semanas	10	1	9

1. ¿Qué deberías hacer según la teoría de la decisión tradicional? Probar que da igual calcular las utilidades de forma usual o de forma “acumulativa”.
2. ¿Qué deberías hacer según la teoría de la decisión de Buchak, suponiendo que $r(x) = x^2$?

Parte K: Experiencia transformadora

En su libro *Transformative Experience*, la filósofa L.A. Paul (2014) propone un desafío para la teoría de la decisión. Según la autora, algunas situaciones no permiten la aplicación de la teoría de la decisión estándar. Esto sucede especialmente en situaciones de *experiencia transformadora*.

Llamamos “experiencia transformadora” a una experiencia que cambiará de forma drástica nuestra escala de utilidades. Por

ejemplo, supongamos que no quiero tener hijos, pero pienso que, si eventualmente me convirtiera en padre, realmente valoraría la paternidad más que cualquier otra cosa. Un caso más polémico es meterse en una secta religiosa. Yo en principio soy ateo, y lleno de conflictos; pero sé que, si me meto en una secta, lo que más va a importarme es el amor a Dios, y voy a estar contento con eso. ¿Debería meterme en la secta?

Según Paul, hay dos problemas para este tipo de decisiones. El primer problema, que podríamos llamarlo “epistémico”, es que *no sé cuáles son las utilidades* de la experiencia antes de realizar la experiencia. Esto aplica a casos de experiencia transformadora más clásicos: por ejemplo, no sé cuán bella es la paternidad antes de ser padre. Y también aplica a otros casos más sencillos: por ejemplo, no sé cuánto me gustará el mole mexicano antes de probarlo, etc. Según Paul, la teoría de la decisión presupone el *conocimiento* de las utilidades futuras, pero esto es inaccesible en estos casos.

Una respuesta obvia es que puedo preguntarles a mis amigos o conocidos cómo se siente, y hacerme una idea. Pero según Paul, tomar una decisión en base a las utilidades de mis amigos le quita *autenticidad* a la decisión.

El segundo problema, que podríamos llamar “existencial”, es que *la experiencia cambia mi punto de vista*. Por ejemplo, sé que el mismo hecho de convertirme en padre cambiará mi punto de vista. Posiblemente, al convertirme en padre, la paternidad me va a interesar mucho más de lo que me interesa ahora. Pero si ahora la paternidad no me interesa, ¿por qué debo privilegiar a ese “yo futuro” respecto a mi yo actual? La teoría de la decisión presupone que el sujeto y sus utilidades no cambian con el tiempo, pero el sujeto y sus utilidades de hecho cambian y mucho.

Por esta razón, Paul dice que *no hay decisiones racionales o irracionales respecto a experiencias transformadoras*. En estos casos, la decisión debo plantearla en términos de si quiero tomar ese riesgo o no.

El problema de este tipo de decisiones había sido planteado anteriormente. Un caso conocido es el existencialismo francés: en su célebre libro *El existencialismo es un humanismo* (1946), Sartre nos plantea la decisión de un joven que no sabe si ir a la guerra

o quedarse cuidando a la madre. Aún dentro de la escuela analítica, Ullman y Morgenbesser (1977) sostienen que en estos casos de decisiones existenciales uno no puede realmente *elegir*, sino simplemente *escoger* (“pick”), sabiendo que no hay razones decisivas a favor de ninguna opción. No podemos optimizar la acción porque no hay un acto que maximice utilidad desde nuestra perspectiva actual, pero podemos tomar decisiones razonables, al escoger alguna opción que no esté claramente dominada por otras.

En la filosofía reciente, distintos autores han intentado explicar las decisiones transformadoras usando variaciones del esquema clásico de la teoría de la decisión. Entre ellos podríamos mencionar a Pettigrew (2019), quien propone que uno debería decidir a partir de la consideración y la ponderación de distintos “yoes” (pasados, presentes y futuros). Callard (2018) también propone que las decisiones transformadoras pueden ser racionales, porque las personas podemos aspirar a ser distintos; y muchas veces, ese “nuevo yo” está presente (aunque sea en forma de ideal) en nosotros mismos.

Soluciones para el capítulo 2

PARTE A

El problema es que enviar más soldados a la guerra no es independiente de ganarla, sino que lo hace más probable. Una forma de reformular la matriz de decisión es que las acciones propias sean la cantidad de soldados que uno envía y los “estados del mundo” sean los soldados que envía el ejército enemigo.

PARTE B

- a. No hay acciones estrictamente dominadas. Ir a Recoleta está *débilmente* dominada por Netflix en casa, porque $c \sim e$ y $d > f$.
- b. Maximax nos indicaría ir al Parque Rivadavia: puede ocurrir el evento preferido a .
- c. Maximin nos indicaría ver Netflix en casa: lo peor que puede suceder es c , que es mejor que f y b .

PARTE C

- a. Agentes 3 y 4.
- b. Agentes 1 y 2. La función es $y = 5x$.
- c. Agentes 1 y 3. La función es $y = 2x + 100$.

PARTE D

1.

$$U(\text{Avión}) = 0.4 \times 40 + 0.5 \times 20 + 0.1 \times -20 = 16 + 10 - 2 = 24$$

$$U(\text{Bus}) = 0.4 \times 20 + 0.5 \times 40 + 0.1 \times 20 = 8 + 20 + 2 = 30$$

$$U(\text{No ir}) = 0 + 0.5 \times 10 + 0.1 \times 150 = 5 + 15 = 20$$

Debo ir en bus.

2.

a. A1.

b. ¿Qué debo elegir si $U(\$x) = \ln(x + 1)$?

$$U(A1) = 1/3 (4.8) = 1.6$$

$$U(A2) = 1/2 (3.91) = 1.95$$

Debo elegir A2.

c. ¿Qué debo elegir si $U(\$x) = \sqrt{x}$?

$$U(A1) = 1/3 (11) = 3.66$$

$$U(A2) = 1/2 (7) = 3.5$$

Debo elegir A1.

3.

¿Debo elegir la caja grande o la pequeña?

$$\begin{aligned} U(\text{Grande}) &= 0.6 \times [8 - (9 - 7)^2] + 0.4 \times [10 - (9 - 7)^2] \\ &= 0.6 \times 4 + 0.4 \times 6 = 2.4 + 2.4 = 4.8 \end{aligned}$$

$$\begin{aligned} U(\text{Pequeña}) &= 0.6 \times [8 - (6 - 7)^2] + 0.4 \times [10 - (6 - 7)^2] \\ &= 0.6 \times 7 + 0.4 \times 9 = 4.2 + 3.6 = 7.8 \end{aligned}$$

Me conviene la caja pequeña.

PARTE E

1.

$$U(\text{frutilla}) = 0.6 \times 70 + 0.4 \times 10 = 42 + 4 = 46$$

2.

Si $A \succ B \succ C$, entonces se da que $A \succ B$, pero $A \# C \sim B \# C$. Esto viola Independencia.

3.

Supongamos que $[\text{iPhone nuevo} > \text{iPhone viejo}] > [\text{Android viejo}]$. El axioma de Continuidad nos dice que existe alguna p tal que $[\text{iPhone nuevo}]p[\text{Android viejo}] > [\text{iPhone viejo}]$. Y esto es lo que rechaza un orden lexicográfico.

PARTE G

Alcanza con analizar los escenarios 1-11 (el escenario 12-100 es idéntico para ambas acciones). Aquí, la utilidad de A1 es $0.11 \times \ln(1.000.001) = 1.51$. Mientras que la utilidad de A2 es $0.1 \times \ln(5.000.001) = 1.54$. Entonces conviene A2, y también A4, aunque la diferencia es mínima.

PARTE H

En el peor escenario, la utilidad sería $0.6 \times -10 + 0.4 \times 10 = -6 + 4 = -2$. Entonces mejor no le doy la aspirina (utilidad 0).

PARTE I

Sesgo de las probabilidades: las personas interpretan 0.99 como algo mucho menor a 0.99.

PARTE J

a. ¿Qué deberías hacer según la teoría de la decisión tradicional? Probar que da igual calcular las utilidades de forma usual o de forma “acumulativa”.

Según la tradicional:

$$U(\text{Una}) = 0.5 \times 7 + 0.4 \times 4 + 0.1 \times 6 = 3.5 + 1.6 + 0.6 = 5.7$$

$$U(\text{Dos}) = 0.5 \times 10 + 0.4 \times 1 + 0.1 \times 9 = 5 + 0.4 + 0.9 = 6.4$$

Debo ir dos semanas.

Usando la versión acumulativa da el mismo resultado:

$$U(\text{Una}) = 4 + 0.6 \times 2 + 0.5 \times 1 = 4 + 1.2 + 0.5 = 5.7$$

$$U(\text{Dos}) = 1 + 0.6 \times 8 + 0.5 \times 1 = 1 + 4.8 + 0.5 = 6.4$$

b. ¿Qué deberías hacer según la teoría de Buchak, suponiendo que $r(x) = x^2$?

$$U(\text{Una}) = 4 + 0.36 \times 2 + 0.25 = 4 + 0.72 + 0.25 = 4.97$$

$$U(\text{Dos}) = 1 + 0.36 \times 8 + 0.25 = 1 + 2.88 + 0.25 = 4.13$$

Ahora me conviene ir solo una semana.

CAPÍTULO 3: TEORÍA DE JUEGOS

El propósito de este capítulo es introducir algunos conceptos fundamentales de la *teoría de juegos*. Solemos pensar a los juegos como situaciones lúdicas y simplemente recreativas. El ejemplo más obvio de un juego es un partido de ajedrez o de fútbol. En la teoría de juegos, sin embargo, nos interesan muchas situaciones, y los “juegos” entendidos a la forma usual son solamente un subconjunto de esas situaciones.

La teoría de juegos es principalmente una teoría sobre las interacciones racionales entre agentes: aquí, un juego es un escenario donde las consecuencias de las acciones de un agente dependen de lo que hagan otros agentes. De hecho, casi todas las interacciones humanas (desde invitar a persona a salir, hasta manejar en una autopista) se podrían modelar como un juego.

Parte A: Juegos estratégicos

De modo similar a lo que sucedía con la teoría de la decisión, empezaremos trabajando con escenarios donde las probabilidades no están presentes. En los juegos más sencillos, los agentes tienen distintas *acciones* a disposición. Por ejemplo, supongamos que los agentes juegan “piedra, papel o tijera”. Asumimos que el lector conoce este juego. En un juego, el resultado de cada acción ya no va a depender del estado del mundo, sino de lo que hagan otros agentes. Podemos partir por un escenario de solamente dos agentes. Supongamos que estamos en un partido de piedra, papel o tijera, donde ambos apuestan \$1, y el que gana se lleva el “pozo” (es decir, se lleva una ganancia neta de \$1). Si ambos juegan lo mismo (por ejemplo, ambos juegan Tijera), se devuelve el dinero a ambos. Podemos representar el escenario de este modo:

	Piedra	Papel	Tijera
Piedra	\$0, \$0	-\$1, \$1	\$1, -\$1
Papel	\$1, -\$1	\$0, \$0	-\$1, \$1
Tijera	-\$1, \$1	\$1, -\$1	\$0, \$0

En la tabla, las filas representan las acciones posibles del Agente 1. Mientras que las columnas representan las acciones posibles del Agente 2. En general, todos los textos de teoría de juegos usan esta misma convención: Fila es Agente 1 y Columna es Agente 2. Si un juego involucra más agentes, ya no puede representarse con una tabla (se suelen usar varias tablas); pero en este capítulo nos restringimos a juegos de dos agentes.

El resultado del juego será ahora un *par* de resultados, es decir, un resultado para cada jugador. Lo escribimos de este modo: “(Resultado para el Agente 1, Resultado para el Agente 2)”. Por ejemplo, (\$1, -\$1) significa que el Agente 1 recibe \$1 y el Agente 2 pierde \$1. Los paréntesis podrían obviarse para facilitar la lectura, como en la tabla anterior. Este juego pertenece a un tipo específico: es un juego de *suma-cero*, donde lo que un jugador pierde, lo gana otro jugador. Pero no todos los juegos son así.

Como señalamos anteriormente, la teoría de juegos no solamente representa *juegos* literalmente, sino también cualquier situación de interacción. Por ejemplo, esta tabla nos permite representar a dos autos en una autopista:

	Carril izquierdo	Carril derecho
Carril izquierdo	0, 0	1, 1
Carril derecho	2, 2	0, 0

La idea de esta tabla es la siguiente. El Auto 1 y el Auto 2 quieren ir rápido, y definitivamente no quieren ir por el mismo carril. Pero el Auto 1 no quiere ir *tan* rápido, y prefiere ir por el carril

derecho (el carril “lento”). Mientras que el Auto 2 prefiere la velocidad, y le gustaría más ir por el lado izquierdo (el carril “rápido”). ¿Qué terminarán haciendo estos autos? Uno querría pensar que elegirán la mejor opción *para ambos*: que el Auto 1 vaya por el carril derecho, y el Auto 2 vaya por el carril izquierdo. Este es un *juego cooperativo*, porque los agentes podrían salir beneficiados de llegar a un acuerdo, y no tienen razones para romper ese acuerdo (a diferencia de otros juegos que veremos luego). Para entender estos conceptos, más adelante usaremos la noción de *equilibrio*.

Filosóficamente, un elemento importante de la teoría de juegos es que los agentes *conocen las utilidades de los otros*. Esto permite que los agentes puedan actuar de forma estratégica, considerando lo que harán los demás. El conocimiento de las utilidades de los otros puede ser un elemento realista (por ejemplo, en una apuesta yo sé cuánto dinero obtendrían los demás), o simplemente una idealización.

La matriz que recién dibujamos es un *esquema de juego*. Un esquema de juego dibujado en *forma estratégica* se compone de lo siguiente:

- Un conjunto de n jugadores.
- Un conjunto S_i de estrategias para cada jugador i (por ahora, una estrategia es una posible acción).
- Un conjunto S de perfiles de estrategia, que describen acciones de cada agente (por ejemplo: <Auto 1 va por carril izquierdo, Auto 2 va por carril izquierdo>).
- Una función que, a cada perfil de estrategia, asigna una n -tupla de utilidades (p. ej., la utilidad de ir ambos por la izquierda es $(0,0)$).
- Un ranking \geq_i de preferencia entre resultados para cada jugador i . Suponemos que este ranking es completo y transitivo. En nuestros ejemplos usamos números para expresar la utilidad, por lo cual será innecesario aclarar el orden (un número más alto representa algo preferible).

Para los juegos que veremos por ahora (y en casi todo este capítulo), las escalas serán ordinales. Podemos poner las utilidades que queramos en tanto preserven el orden.

Utilidad, dinero y egoísmo

Antes de seguir, también hace falta aclarar que en un contexto colectivo la utilidad no necesariamente sigue al dinero. Es decir, un escenario donde yo obtengo \$1000 no necesariamente es mejor para mí que un escenario donde yo obtengo \$900. De hecho, dado que los resultados involucran a *varios agentes*, la utilidad de los resultados dependerá de la disposición de cada uno: hay agentes más altruistas, otros más egoístas, otros igualitarios, etc. La idea de que el dinero tiene una utilidad marginal decreciente es correcta en la escala individual, pero a escala colectiva podría haber otros factores en juego.

Por ejemplo, supongamos que tenemos esta situación. Dos amigos pueden aplicar a un pequeño trabajo de \$20, y lo harán sin saber qué hace el otro. Si ambos aplican, se reparten el trabajo y ganan \$10 cada uno. Si uno aplica y el otro no, el que aplicó se quedará con \$15, y le dará \$5 al otro. Y si ninguno de los dos aplica, no ganarán nada.

	Aplicar	No aplicar
Aplicar	\$10, \$10	\$15, \$5
No aplicar	\$5, \$15	\$0, \$0

Si ambos sujetos son egoístas y ambiciosos (es decir, solo se preocupan por el dinero que ellos mismos obtienen), la utilidad podría verse así:

	Aplicar	No aplicar
Aplicar	2, 2	3, 1
No aplicar	1, 3	0, 0

Sin embargo, podría ser que el Agente 1 fuera muy igualitarista, y lo que más le interesa es la igualdad (sólo se preocupa por la

cantidad si ambos reciben lo mismo, pero rechaza cualquier escenario no igualitario). Entonces la matriz podría verse así:

	Aplicar	No aplicar
Aplicar	2, 2	0, 1
No aplicar	0, 3	1, 0

De todas formas, en los ejemplos siguientes vamos a trabajar directamente con utilidades. Y a falta de aclaración, vamos a asumir que los agentes son de hecho egoístas y ambiciosos; usualmente se llama *homo-economicus* a este modelo de ser humano. Pero el egoísmo respecto al dinero o los bienes, valga remarcar, *no* es un elemento esencial de la teoría de juegos.²²

Dominancia débil y estricta

Lo primero que podríamos asumir en un juego es que los agentes racionales no tomarán decisiones dominadas; esta idea proviene de la teoría de la decisión.

Los conceptos de dominancia se definen de forma muy similar a como lo hacemos en teoría de la decisión:

(Dominancia estricta) Una estrategia *A* *domina estrictamente* otra estrategia *B* para el jugador *n* si y sólo si para toda posible estrategia de los otros jugadores, la estrategia *A* del jugador *n* en conjunción con las estrategias de los demás da un *mejor* resultado para *n* que el que daría la estrategia *B*.

Decimos también que una estrategia es *estrictamente dominante* si domina estrictamente a todas las otras.

Para ver si una estrategia es dominante no hace falta ver el juego entero, sino los resultados para el jugador en cuestión. Por ejemplo:

²² Una crítica clásica al *homo economicus* se encuentra en Sen (1977).

	E	F	G
A	4, ...	3, ...	2, ...
B	2, ...	2, ...	1, ...
C	5, ...	0, ...	2, ...
D	5, ...	1, ...	2, ...

Vista esa matriz, podemos preguntarnos si hay alguna estrategia estrictamente dominada por otra. Y la respuesta es que sí: para el jugador 1, la estrategia B está estrictamente dominada por la estrategia A. Cualquier cosa que haga el jugador 2, al jugador 1 le conviene hacer A en vez de B.

Por otro lado, existe una noción más débil de dominancia:

(Dominancia débil) Una estrategia A *domina débilmente* otra estrategia B para el jugador n si y sólo si:

- En todos los casos, la estrategia A del jugador n da *igual o mejor* resultado para n que lo que daría la estrategia B;
- En algunos casos, la estrategia A del jugador n da *mejor* resultado para n que lo que daría la estrategia B.

Por ejemplo, en la matriz anterior, la estrategia D domina débilmente a la C para el jugador 1. Observemos que, si una estrategia domina estrictamente a otra, también la domina débilmente.

Ahora podemos introducir el concepto de *equilibrio*. Intuitivamente, un “equilibrio” es una situación donde cada jugador hizo lo mejor que pudo. Decimos que un perfil de estrategias (o sea, una celda) es un *equilibrio de estrategias estrictamente dominantes* cuando para cada jugador i , s_i es una estrategia estrictamente dominante. (Análogamente podemos hablar de equilibrios de estrategias débilmente dominantes).

En el juego de Aplicar al Trabajo, cuando ambos agentes son egoístas y ambiciosos, el equilibrio de estrategias estrictamente

dominantes es que ambos apliquen. Esto se debe a que Aplicar es estrictamente dominante para ambos agentes individualmente (marco las estrategias dominantes con negrita):

	Aplicar	No aplicar
Aplicar	2, 2	3, 1
No aplicar	1, 3	0, 0

Es natural pensar que, si existe un equilibrio de estrategias estrictamente dominantes, los agentes racionales tomarán esa opción. Por ejemplo, podríamos predecir que todos los agentes racionales jugarán <Aplicar, Aplicar>. Aunque este supuesto suele mantenerse en muchos casos, también tuvo cuestionamientos, como veremos luego.

Procedimientos de borrado iterado

Como vimos en teoría de la decisión, el principio fundamental de la racionalidad es la *dominancia*: para que un agente sea considerado racional, es necesario que no tome decisiones dominadas por otras. De ahí viene la importancia del equilibrio de estrategias estrictamente dominantes.

El *procedimiento de borrado iterado (de estrategias estrictamente dominadas)* podría verse como una extensión del concepto de equilibrio de estrategias estrictamente dominantes. La idea aquí es que, en un juego, no solo podemos asumir que los agentes racionales evitan acciones dominadas: también podemos asumir que los otros jugadores saben que existe esta presunción de racionalidad.

Este proceso de borrado iterado nos dice que, si encontramos una estrategia estrictamente dominada, podemos borrarla (porque todos los agentes saben que el agente no tomará esa estrategia). Ahora obtenemos un nuevo juego. Si encontramos una nueva estrategia estrictamente dominada, podemos borrarla. Así sucesivamente, hasta que no haya estrategias estrictamente dominadas.

Si llegamos a una celda única, diremos que encontramos una *solución* al juego.

Para el ejemplo de Aplicar al Trabajo, donde ambos agentes son egoístas y ambiciosos, partimos del juego entero:

	Aplicar	No aplicar
Aplicar	2, 2	3, 1
No aplicar	1, 3	0, 0

Aplicar al trabajo es dominante para el Agente 1, entonces borramos la otra estrategia.

	Aplicar	No aplicar
Aplicar	2, 2	3, 1

Ahora podemos ver que No Aplicar está dominado para el Agente 2, entonces borramos su estrategia de No Aplicar.

	Aplicar
Aplicar	2, 2

Y bien, ahora tenemos una solución al juego.

Esto no solo sucede con juegos simétricos. Supongamos que estamos en la versión del juego donde el Agente 1 es estrictamente igualitarista, y las utilidades se ven así:

	Aplicar	No aplicar
Aplicar	2, 2	0, 1
No aplicar	0, 3	1, 0

Aquí, el Agente 1 no tiene una estrategia dominante. Pero el Agente 2 sí tiene: puede Aplicar. Entonces podemos eliminar No Aplicar para el Agente 2:

	Aplicar
Aplicar	2, 2
No aplicar	0, 3

Pero el Agente 1 “sabe” que el Agente 2 hará esto. Por ende, también decide aplicar:

	Aplicar
Aplicar	2, 2

Esta es la solución al juego. Conceptualmente, la idea sería la siguiente: en el juego Aplicar al Trabajo, el Agente 2 puede pensar “Haga lo que haga el Agente 1, me conviene Aplicar”. Entonces va a aplicar. Mientras que el Agente 1 podría pensar “el Agente 2 va a aplicar, porque le conviene en ambos casos; entonces me conviene también aplicar”. Es decir, cada agente actúa suponiendo qué es lo que racionalmente hará el otro.

Entonces, el comportamiento racional en este juego será igual para el egoísta o el igualitarista, si sabe que el otro es egoísta. Podemos practicar el método con una matriz más compleja:

	D	E	F
A	7, 3	5, 4	2, 1
B	1, 2	6, 3	3, 4
C	3, 0	3, 6	6, 7

El Agente 1 no tiene estrategias dominadas. Pero sí el Agente 2: en el primer paso, borramos la D, que está dominada por la E.

	E	F
A	5, 4	2, 1
B	6, 3	3, 4
C	3, 6	6, 7

Ahora podemos eliminar la A, que está dominada por la B:

	E	F
B	6, 3	3, 4
C	3, 6	6, 7

Ahora podemos borrar la E, que está dominada por la F:

	F
B	3, 4
C	6, 7

Por último, eliminamos la B, que está dominada por la C. Así llegamos a una solución: CF, con utilidad (6,7).

Podríamos interpretar la situación del siguiente modo. En el primer paso, la D está dominada por la E. Es decir, E es mejor que D en cualquier situación, para el Agente 2. Entonces racionalmente el Agente 2 no va a elegir esa opción. Pero todos conocen el juego. Y esto es sabido por el Agente 1. Entonces el Agente 1 va a restringir su atención al juego sin D. Así restringido, el Agente 1 no va a hacer A, porque está dominada por B. Pero esto lo sabe el Agente 2, que va a restringir su atención al juego sin

A. Entonces, el Agente 2 hará F. Conociendo este razonamiento, el Agente 1 hará C.

El procedimiento de borrado iterado de estrategias estrictamente dominadas presupone el concepto de *conocimiento común de racionalidad*. La idea es que todos creen que todos son racionales, pero a la vez que todos creen que todos creen que todos creen que todos son racionales, y así *ad infinitum*. Esta es una idealización clásica en la teoría de juegos.

El procedimiento de borrado iterado no siempre arroja una solución (i.e., una celda única). Muchas veces, “sobreviven” muchas celdas. Por ejemplo, si quisiéramos aplicar este método al juego de los carriles, no podríamos eliminar ninguna fila o columna.

Equilibrio de Nash

Como vimos, no todos los juegos tienen un equilibrio de estrategias estrictamente dominantes, o una solución de borrado iterado. Un ejemplo es el juego de los carriles mencionado antes. Para esos casos, noción de *Equilibrio de Nash* es especialmente útil. El concepto fue creado por el matemático John Nash (1950), y le valió el Premio Nobel de Economía en 1994.

Intuitivamente, un estado es un equilibrio de Nash cuando cada jugador hace lo mejor posible, dejando fijo lo que hicieron los otros jugadores.

Más formalmente, cuando hay dos agentes i y j , un estado (S_i, S_j) es un *Equilibrio de Nash* si y sólo si:

- a) La utilidad para i de (S_i, S_j) es mayor o igual a la utilidad para i de (S_i^*, S_j) para cualquier otra estrategia S_i^* de i .
- b) La utilidad para j de (S_i, S_j) es mayor o igual a la utilidad para j de (S_i, S_j^*) para cualquier otra estrategia S_j^* de j .

Naturalmente, un juego puede tener varios equilibrios de Nash. Por ejemplo, imaginemos el siguiente juego:

	D	E	F
A	3, 6	0, 0	1, 1
B	1, 0	2, 5	4, 4
C	1, 2	3, 7	3, 5

Una forma de encontrar equilibrios de Nash es la siguiente. Primero, busco los ganadores de cada columna para el Agente 1 (es decir, la mejor acción del Agente 1, para cada acción del Agente 2), y elimino a los perdedores (los marco con una “x”).

Por ejemplo, en la columna D gana la acción A, que le da 3 puntos al Agente 1. Las otras no pueden ser equilibrio de Nash.

	D	E	F
A	3, 6	0, 0 (x)	1, 1 (x)
B	1, 0 (x)	2, 5 (x)	4, 4
C	1, 2 (x)	3, 7	3, 5 (x)

Luego de marcar los ganadores para cada columna, buscamos los ganadores de cada fila para el Agente 2. Por ejemplo, si el Agente 1 hace A, al Agente 2 le conviene hacer D. Marcamos con una “x” las demás. De este modo obtenemos los dos equilibrios de Nash, que son AD y CE:

	D	E	F
A	3, 6	0, 0 (x)	1, 1 (x)
B	1, 0 (x)	2, 5 (x)	4, 4 (x)
C	1, 2 (x)	3, 7	3, 5 (x)

Para probarlo, miremos AD. Dado que el Agente 1 eligió A, si el Agente 2 cambiaba de columna obtenía menos (0 o 1). Y dado que el Agente 2 eligió D, si el Agente 1 cambiaba de fila obtenía menos (1). Por eso AD es un equilibrio de Nash. El mismo razonamiento podríamos aplicarlo a CE.

Hay distintas formas de interpretar el concepto de “equilibrio de Nash”. Por ejemplo, podemos apelar al *arrepentimiento*: un estado es un equilibrio de Nash cuando, al ver lo que hicieron los otros, ningún individuo se arrepiente de lo que hizo. También podemos apelar al concepto de *acuerdo autoimpuesto*: un estado X es un equilibrio de Nash cuando, si todos los jugadores arreglaran antes en hacer X, ninguno tendría incentivos en desviarse de lo arreglado.

En algunas aplicaciones, el equilibrio de Nash suele tener también un valor *predictivo*. Por ejemplo, en el juego de los carriles, podemos asumir que, con suficiente tiempo, los agentes (si son racionales) van a terminar tomando carriles distintos. La idea es que los agentes racionales terminan ajustando sus acciones hacia alguna estrategia colectivamente estable. Más adelante volveremos a este tema (especialmente en la parte D).

Podemos enfocarnos ahora en la relación entre los equilibrios de Nash y los otros tipos de equilibrio. Naturalmente un equilibrio de estrategias estrictamente dominantes es también un equilibrio de Nash. Pero la relación no se da a la inversa: en el juego de los carriles, tanto <Izquierdo, Derecho> como <Derecho, Izquierdo> son equilibrios de Nash. Sin embargo, ninguno de los dos es un equilibrio de estrategias dominantes.

El procedimiento de borrado iterado de estrategias estrictamente dominadas, cuando arroja una solución (única), también nos da un equilibrio de Nash. Y cuando el procedimiento no arroja un resultado único, los equilibrios de Nash estarán entre las celdas que “sobrevivieron”.²³

Del mismo modo, podríamos decir provisoriamente que *no todo juego tiene un equilibrio de Nash*, al menos si entendemos el

²³ Puede encontrarse una prueba sencilla en Bonanno (2015), p. 52.

equilibrio de Nash de la forma sencilla que describimos anteriormente. Un ejemplo es “Piedra, papel o tijera”. Aquí, cualquiera sea el resultado, alguno de los dos jugadores va a desear haber jugado distinto.

Nash estableció que en realidad sí existe un equilibrio para todo juego, pero usando “estrategias mixtas”. Este es un concepto bastante más complejo, que veremos más adelante.

Ejercicios

1. Realice el procedimiento de borrado iterado de estrategias estrictamente dominadas en esta matriz de decisión:

	E	F	G	H
A	3, 2	5, 4	1, 1	3, 9
B	5, 1	6, 3	3, 2	5, 1
C	5, 0	3, 2	6, 1	4, 0
D	2, 9	2, 8	2, 3	6, 1

2. Encuentre los equilibrios de Nash en este juego:

	E	F	G	H
A	3, 10	5, 4	1, 1	2, 9
B	3, 7	2, 9	2, 2	5, 3
C	2, 0	3, 7	1, 1	4, 5
D	1, 9	3, 8	2, 3	6, 10

3. Antes señalamos que todo equilibrio de Nash “sobrevive” al borrado de estrategias estrictamente dominadas. ¿Se mantiene el

resultado si borramos estrategias *débilmente* dominadas? Piense la respuesta a partir de este ejemplo:

	C	D
A	1, 1	0, 0
B	0, 0	0, 0

Parte B: Juegos dinámicos

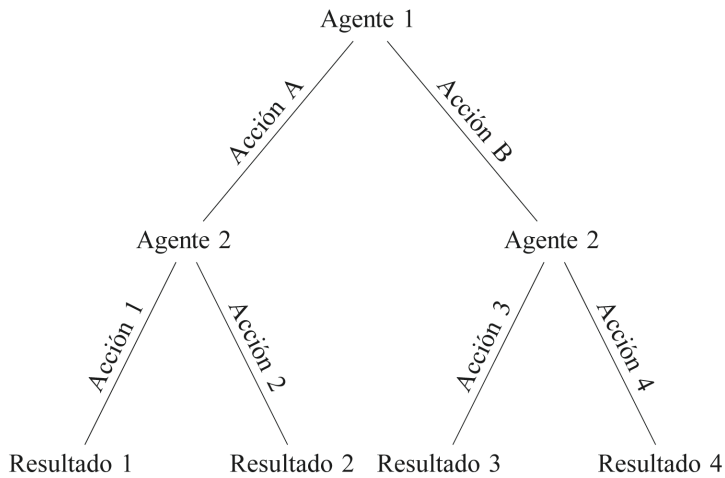
Algunos juegos no son simultáneos (como el “piedra, papel o tijera”) sino secuenciales. A estos juegos los podemos llamar *dinámicos* o *extendidos*. Por ahora nos concentraremos en juegos de información perfecta, donde cada movimiento de cada jugador es visible para los demás jugadores. Dos casos de juegos dinámicos con información perfecta son el ajedrez y las damas. Podemos representar los juegos dinámicos usando árboles. Introducimos ahora el concepto de *árbol* para luego mostrar ejemplos. Un *árbol con raíz dirigido* tiene esta forma:

- La raíz del árbol (el nodo que está arriba de todo) no tiene ninguna flecha que la señale, mientras que todos los otros nodos del árbol son señalados por una flecha.
- A partir de cada nodo, hay *un solo camino* que lleva a la raíz.
- Los nodos que no flechan a otros se llaman *terminales*, mientras que los demás se llaman *nodos de decisión*.

Un juego puede verse como un árbol finito con raíz dirigido, donde:

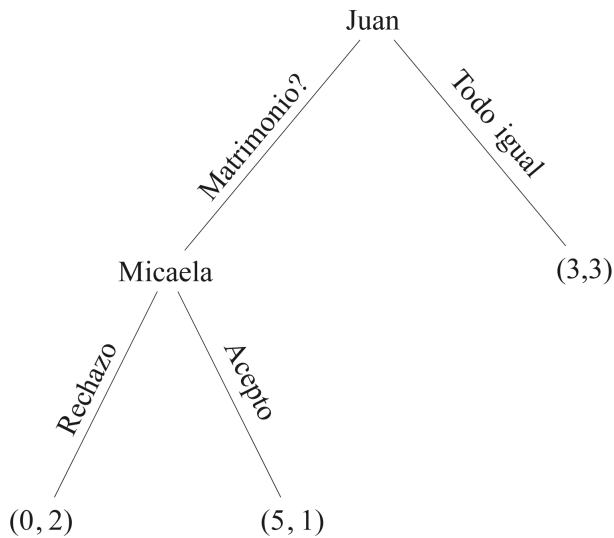
- Cada nodo está asociado a un agente (un agente puede tener asociados varios nodos, uno por cada decisión que toma).
- Cada flecha está asignada a una acción.
- Cada nodo terminal está asociado con un resultado.

Un árbol de decisión podría verse así (obviamos la “punta” de la flecha, porque siempre apuntan hacia abajo):



Ejemplo: el enamorado realista

Juan y Micaela están saliendo hace un par de meses. Juan está locamente enamorado, y no puede imaginar algo más hermoso que casarse con Micaela. Micaela, por su parte, quiere ir de a poco. Esto ya está charlado entre ellos.



Aquí, cada nodo terminal representa el resultado para ambos, donde (x, y) representa el resultado para Juan (x) y el resultado para Micaela (y). ¿Qué debería hacer Juan? Obviamente prefiere casarse. Pero sabe que, si le propone matrimonio a su novia, ella lo va a rechazar. Entonces, tomando eso en cuenta, lo mejor es que siga todo igual.

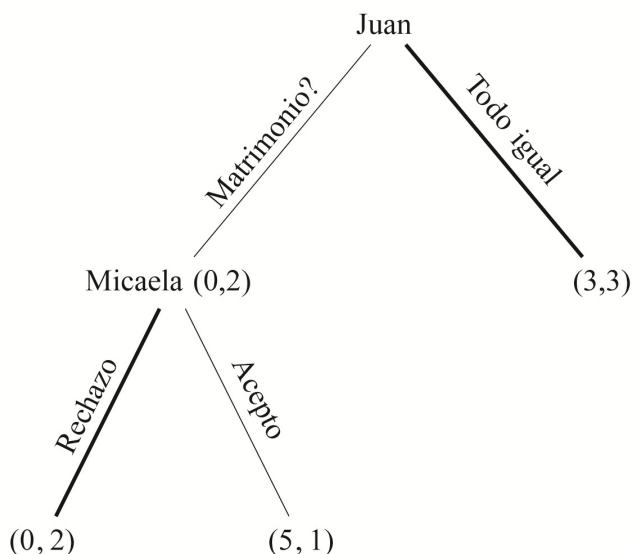
Este tipo de razonamiento práctico, donde los agentes actúan en vistas de lo que piensan que harán los otros agentes, se llama *inducción hacia atrás*, y lo explicaremos en más detalle a continuación.

Inducción hacia atrás

El proceso de *inducción hacia atrás* (*backward induction*) nos permite determinar una decisión racional en un escenario de juegos dinámicos (esto puede servir para tomar decisiones, o predecir decisiones de otros). Una de las primeras versiones de este método (para juegos de suma-cero) apareció en el citado libro de Von Neumann y Morgenstern (1944), que suele considerarse como el punto de partida de la teoría de juegos moderna.

El proceso funciona en etapas. Para cada nodo de decisión, marcamos la opción que le conviene al agente que toma la decisión, y anotamos el resultado de la opción en el nodo decisorio. Empezamos desde los nodos terminales y vamos subiendo hasta la raíz, es decir, hasta la primera decisión.

La idea filosóficamente es que cada uno va a decidir qué hacer suponiendo que los otros (y uno mismo) van a decidir lo que sea mejor para ellos (es decir, un supuesto básico de racionalidad), y que los otros van a asumir también nuestra racionalidad y la de todos los demás (y saben que nosotros asumimos su racionalidad, etc.). Por ejemplo, en el caso del “enamorado realista” hacemos lo siguiente:



En la primera etapa, comienzo marcando el “Rechazo” porque es lo mejor para Micaela: le da 2 utilidades en vez de 1. Entonces Juan va a decidir entre (0, 2) y (3, 3), y obviamente va a decidir que todo siga igual, que le da 3 utilidades en vez de 0. De forma coloquial: Juan no va a proponerle matrimonio a Micaela, porque sabe que, si lo hiciera, ella lo rechazaría.

Vale mencionar que, aunque no lo mostramos aquí, un mismo juego podría tener distintas soluciones de inducción hacia atrás, en caso de que haya un empate en algún nodo. Si algún jugador es indiferente entre una opción u otra, los otros jugadores deben tener en cuenta todas esas opciones.

Relación entre juegos dinámicos y estratégicos

En juegos dinámicos también podemos hablar de *estrategias*. Una estrategia es intuitivamente una planificación como esta: “Si el otro hace A, yo haré C, y si el otro hace B, yo haré D”. Es decir, las estrategias son posibles respuestas a todas las movidas de los demás jugadores.

En el árbol, una estrategia de un jugador es simplemente *una lista de decisiones*, una para cada nodo de decisión de ese jugador.

Por ejemplo, en el esquema anterior, Micaela podría tener la estrategia (Si me propone, acepto), mientras que Juan podría tener la estrategia (Sigue todo igual).

Podemos escribir ese juego dinámico como un juego estratégico:

	Acepto	Rechazo
¿Matrimonio?	5, 1	0, 2
Todo igual	3, 3	3, 3

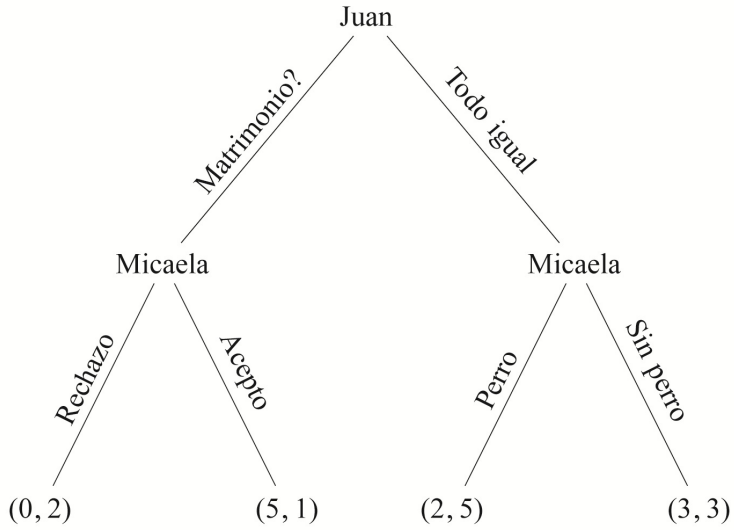
Al pasar de un árbol a una tabla, las estrategias podrían ser un poco redundantes, por razones meramente técnicas. Esto se debe a que, si Juan no le ofrece matrimonio, no tiene mucha importancia si Micaela acepta o no.

Podríamos analizar ahora los distintos equilibrios. En principio, no hay equilibrios de estrategias estrictamente dominantes. Pero sí hay un equilibrio de Nash: (Todo igual, Rechazo). Este es el mismo equilibrio encontrado usando inducción hacia atrás.

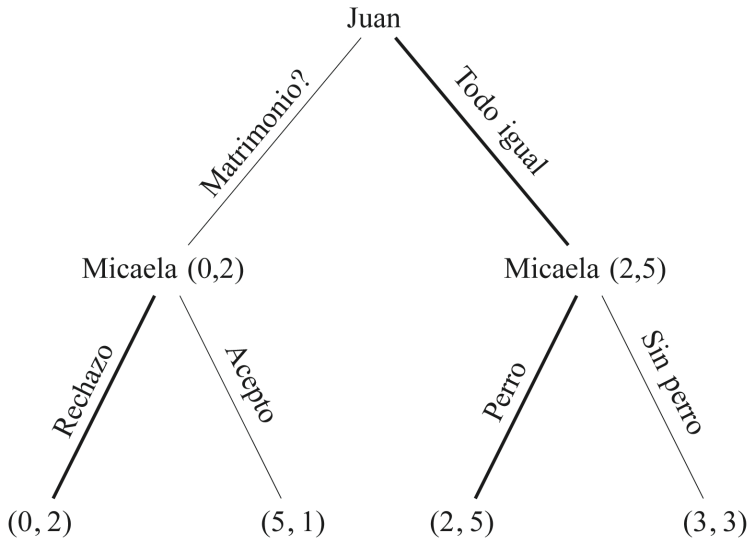
Este ejemplo ilustra un resultado general: toda solución de inducción hacia atrás en forma dinámica será equilibrio de Nash en forma estratégica. Esto no se da a la inversa (véase Ejercicios).

El paso de juegos dinámicos a la forma estratégica se vuelve más complejo cuando tenemos más opciones de decisión. Por ejemplo, supongamos que en caso de que Juan no le proponga matrimonio, Micaela no tendrá que ahorrar para la boda, y podrá volver a pensar su plan original de adoptar un perro. Adoptar un perro le resulta muy atractivo a ella, pero no tanto a Juan, que sentirá que el perro se llevará más atención que él.

El árbol decisorio será ligeramente distinto al original.



La solución por inducción hacia atrás quedaría así:



Es decir, incluso sabiendo que al no proponer matrimonio terminará conviviendo con el perro de su novia, Juan decide actuar así, porque proponer matrimonio y ser rechazado es aún peor.

En términos de estrategias (es decir, utilizando matrices), necesitaremos una representación más compleja. Una estrategia es una forma de decidir por anticipado para cada posible situación de decisión, porque no podemos anticipar lo que harán los demás; es decir, una estrategia decide en cada uno de los nodos decisivos del agente en cuestión.

Por eso, ahora Micaela tendrá cuatro estrategias posibles:

- Si Juan se propone, acepto; si no se propone, adopto un perro.
- Si Juan se propone, rechazo; si no se propone, adopto un perro.
- Si Juan se propone, acepto; si no se propone, no adopto un perro.
- Si Juan se propone, rechazo; si no se propone, no adopto un perro.

Esto se traduce en una matriz de decisión más compleja:

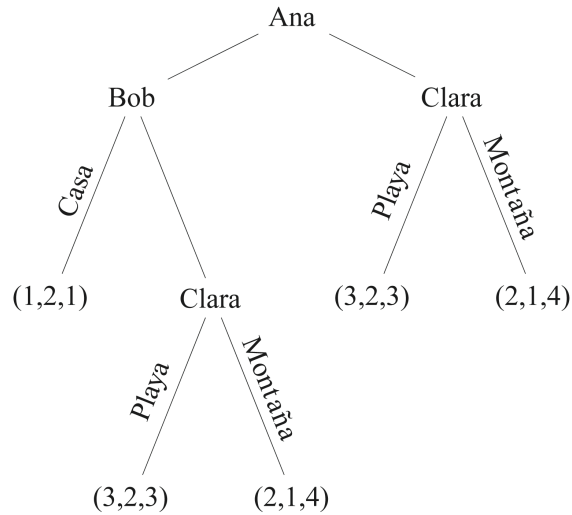
	Acepto/ Perro	Rechazo/ Perro	Acepto/ Sin perro	Rechazo/ Sin perro
Propongo Matrimonio	5, 1	0, 2	5, 1	0, 2
Todo igual	2, 5	2, 5	3, 3	3, 3

Sin embargo, como podemos ver, el equilibrio de inducción hacia atrás (Todo sigue igual, Rechazo/Perro) sigue siendo un Equilibrio de Nash. De hecho, es el único equilibrio de Nash.

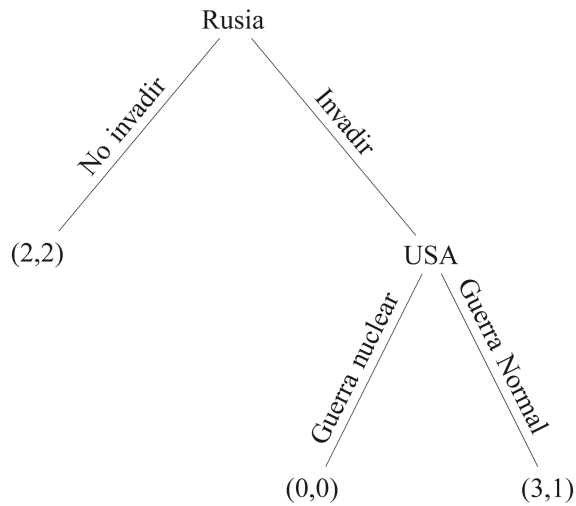
Ejercicios

1. Luego de años de peleas, los padres de Ana, Bob y Clara (Agentes 1, 2 y 3) deciden dejar que ellos elijan a donde ir de vacaciones. Clara, la hermana menor, nunca toma decisiones, así que los padres prefieren que sea ella quien decida esta vez. Ana puede hacer que decida Clara o darle un poder de veto a Bob, que

puede hacer que todos se queden en casa. Resuelva el juego por inducción hacia atrás:



2. Rusia (Agente 1) está considerando la invasión de un país de la OTAN, protegido por Estados Unidos (Agente 2).



- a. Resolverlo por inducción hacia atrás.
- b. Escribir su forma estratégica correspondiente y hallar sus equilibrios de Nash.
- c. ¿Qué podemos concluir sobre la relación entre equilibrios de inducción hacia atrás y equilibrios de Nash?

Parte C: Cooperación

El Dilema del Prisionero

El *Dilema del Prisionero* es un ejemplo problemático para los conceptos de equilibrio vistos anteriormente, porque se da un equilibrio de estrategias estrictamente dominantes (y por ende, un equilibrio de Nash), pero la solución es poco satisfactoria.

El escenario planteado es de dos ladrones que son atrapados por la policía, y pueden confesar (traicionando al otro) o cooperar entre ellos, manteniendo el silencio. Si cooperan ambos, salen libres. Si ambos confiesan (es decir, se traicionan mutuamente), tienen una pena razonable. Y si uno confiesa y el otro no, el que no confiesa tiene una pena muy grave, y el “traidor” sale libre con un considerable premio económico.

Así se ve el juego, con las utilidades correspondientes:

	Cooperar	Traicionar
Cooperar	2, 2	0, 3
Traicionar	3, 0	1, 1

Aquí, para el Agente 1 traicionar es estrictamente dominante. Lo mismo sucede para el Agente 2. Entonces (Traicionar, Traicionar) es el único equilibrio de estrategias estrictamente dominantes (y también, el único equilibrio de Nash).

Pero en el Dilema del Prisionero ocurre algo particular: intuitivamente, (Traicionar, Traicionar) es peor que (Cooperar, Cooperar). Para entenderlo, usaremos algunos términos nuevos:

(Pareto Superior) Decimos que un resultado o_1 es *estrictamente Pareto superior* a un resultado o_2 si y sólo si o_1 es mejor para todos los jugadores que o_2 .

Lo que sucede en el Dilema del Prisionero es que la conjunción de dos estrategias dominadas (la cooperación mutua) es estrictamente Pareto superior al equilibrio de estrategias estrictamente dominantes (que ambos traicionen). Autores como Peterson (2009, p. 215), lo interpretan de este modo: “lo que es óptimo para el grupo no es óptimo para el individuo”.

Este problema ha generado mucha discusión entre filósofos, economistas y politólogos. De hecho, el Dilema del Prisionero es el problema filosófico más estudiado en torno a la teoría de juegos, y existen miles de artículos académicos sobre el tema. Esto se debe a que involucra problemas de cooperación, leyes, castigo, normas sociales y convenciones, que son importantes para diferentes disciplinas.

Interpretaciones del Dilema del Prisionero

Hay muchas interpretaciones filosóficas del Dilema del Prisionero. Mencionaré dos lecturas clásicas y contrapuestas del dilema: la de Ken Binmore y la de David Gauthier.

Binmore (2015) interpreta el dilema del modo más literal posible: el Dilema del Prisionero muestra que, en ciertos escenarios, cooperar es irracional. El Dilema del Prisionero se compara con otros dilemas sociales. En muchos casos, una persona puramente egoísta obtendría beneficios en ser la única que no coopera, mientras los demás sí cooperan: por ejemplo, si todos se vacunan de sarampión, pero yo no me vacuno (porque prefiero evitar un dolor en el brazo), quedará protegido por la inmunidad de rebaño. En general se llaman a estos sujetos “aprovechadores” (*free riders*). Claro que, si todos se aprovecharan, terminaríamos en resultados indeseables: esto se conoce como la “Tragedia de los comunes”. Ahora bien, el Dilema del Prisionero plantea el problema de si, individualmente, estas acciones son irracionales.

Según Binmore, es imposible justificar la cooperación en estos casos: “Es verdad que sería malo si todos se comportaran antisocialmente, pero yo no soy todos; yo soy yo” (2015, p. 16). Según este autor, los agentes racionales deben *traicionar* en el Dilema del Prisionero, como bien indica la teoría de juegos estándar.

Para Binmore, los defensores de la idea que cooperar es racional cometen varias falacias. Entre ellas, la *Falacia de los Gemelos*. Esto consiste en asumir que, como el juego es simétrico, el otro va a hacer lo mismo que yo. Bajo ese supuesto, voy a cooperar. Pero este supuesto no está justificado: lo correcto es asumir que el otro va a traicionarme (aunque yo coopere). De forma similar, el *Imperativo Categórico* (leído coloquialmente) diría que lo racional es hacer lo que es mejor que *todos* hagan al mismo tiempo; si así fuera, veré solo las “diagonales” del juego, y cooperaré. Pero el imperativo no se justifica por sí solo: bajo la regla de dominancia, lo racional es traicionar (aunque el otro coopere). Binmore propone comparar el Dilema del Prisionero con el juego de la Caza del Ciervo:

	Cazar ciervo	Cazar conejitos
Cazar ciervo	15, 15	0, 5
Cazar conejitos	5, 0	5, 5

La idea de este juego (inspirado en un fragmento de Rousseau) es que dos agentes se debaten sobre qué cazar. Pueden cooperar y cazar un ciervo. O pueden ir individualmente a cazar conejos. El peor escenario es intentar cazar el ciervo individualmente; por hipótesis, decimos que eso es imposible, porque se requiere de dos personas para cazarlo. Aquí, hay dos equilibrios de Nash: que ambos vayan a cazar conejos individualmente, o que cooperen. A diferencia del caso del Prisionero, aquí la cooperación (ir a cazar el ciervo entre ambos) sí es un equilibrio de Nash; por eso, se trata de un *juego cooperativo*. Los agentes terminarán cooperando si creen que los demás van a cooperar también. Esto

es obviamente distinto al caso del Dilema del Prisionero, donde si asumimos que los demás cooperan, nos conviene traicionar.

En resumen, para Binmore no hay nada especial en el Dilema del Prisionero: es un juego en donde cualquier agente racional debe traicionar. Si queremos entender la estructura de la cooperación, debemos estudiar otros juegos, como la Caza del Ciervo (o el Dilema del Prisionero iterado, que veremos luego).

A diferencia de Binmore, Gauthier (2015) propone que lo racional en el Dilema del Prisionero es cooperar. Porque no deberíamos identificar la racionalidad con la maximización de utilidad esperada individual. La situación de Cooperar-Cooperar es *Pareto superior* a la de Traicionar-Traicionar, porque da más utilidad a ambos jugadores. Asimismo, la situación de Cooperar-Cooperar es *Pareto óptima* (o *Pareto eficiente*), porque ninguna otra situación es mejor para ambos.

En la interacción social, según Gauthier, el objetivo es sacar suficiente provecho de una situación para *todos*. Por eso, en vez de buscar equilibrios de Nash, tenemos que buscar situaciones Pareto óptimas. Si esta situación Pareto óptima además es Pareto superior a los equilibrios de Nash (como sucede en el Dilema del Prisionero), entonces definitivamente es la elección más racional. Según Gauthier: “en las interacciones, la irracionalidad no consiste en la falla de los individuos para obtener lo mejor en su situación particular, sino en la falla de todos los partícipes de la interacción en obtener lo mejor de su situación conjunta” (p. 39). Gauthier contrapone su visión cooperativa, donde los agentes “interactúan” entre sí, con la visión estándar, donde los agentes solamente “responden” uno al otro, y se ven entre sí como obstáculos para la búsqueda de bienestar individual. El enfoque cooperativo refleja la naturaleza social del ser humano, capaz de cumplir compromisos y actuar en beneficio a otros; el enfoque estándar está dogmáticamente convencido de que todas nuestras decisiones buscan principalmente el bienestar individual.

En conclusión, para Gauthier, el Dilema del Prisionero muestra un choque entre dos conceptos de racionalidad: equilibrio de Nash y Pareto-optimalidad. En ocasiones ambos criterios coinciden, pero otras veces no. Además del choque entre los dos conceptos de racionalidad, el Dilema del Prisionero muestra que el

concepto de Pareto-optimalidad es mejor que el de Equilibrio de Nash para situaciones de interacción, porque trae mejores resultados.

Convenciones y normas

Algunos filósofos utilizaron los juegos cooperativos y el Dilema del Prisionero para ilustrar el surgimiento de convenciones y normas sociales. En su conocido libro *Convention* (1969), David Lewis propone que las convenciones son un tipo de *juego de coordinación*:

	Ir al parque	Ir al cine
Ir al parque	1, 1	0, 0
Ir al cine	0, 0	1, 1

Un juego de coordinación es un tipo de juego cooperativo con distintos equilibrios de Nash, donde (en la forma más usual) los agentes salen beneficiados de hacer lo mismo. En este ejemplo, los dos agentes son viejos amigos, y solo se preocupan por estar juntos: no tienen preferencias estrictas por ir al parque o al cine. Otras convenciones podrían favorecer a algún agente:

	Hablar español	Hablar inglés
Hablar español	2, 1	0, 0
Hablar inglés	0, 0	1, 2

Esta tabla (conocida como “Guerra de los Sexos”) podría ilustrar una situación donde el Agente 1 es hispanohablante nativo, y el Agente 2 es angloparlante nativo. Los dos pueden hablar ambos idiomas decentemente, y deben definir en qué idioma conversar: cada uno prefiere hablar en su idioma nativo.

Este tipo de convenciones iluminan un problema clásico de la teoría de juegos: la necesidad de *elegir* entre equilibrios. La multiplicidad de equilibrios muestra que las convenciones son, en cierta medida, arbitrarias. La elección de equilibrios va a depender de factores como la costumbre o las relaciones de poder. Binmore (1998) y más recientemente Vanderschraaf (2019) exploraron la idea de justicia como la búsqueda de equilibrios justos. Según estos autores, centrar la filosofía política en la búsqueda de equilibrios (en vez de situaciones ideales) permite llegar a escenarios más sostenibles en el tiempo.

Lo importante de las convenciones, entonces, es que las personas prefieren cumplirlas, en tanto suponen que los otros la cumplen. Esto distingue a las convenciones de las normas sociales, que traen cierto “costo” para los agentes. Según Bicchieri (2005), una norma social surge de un juego de “motivación mixta”, como el Dilema del Prisionero, donde hay motivación para cooperar y también para traicionar.²⁴ Las normas sociales incluyen actos como tirar la basura en el tacho (y no en el suelo), pagar propina, o cumplir promesas. Aquí, podríamos vernos tentados a incumplir la norma. Para Bicchieri, la razón por la que cumplimos con las normas sociales es que sabemos que la mayor parte de la gente las cumple (“expectativa empírica”), y sabemos que se espera que nosotros también cumplamos la norma (“expectativa normativa”), a veces bajo pena de castigo. En las convenciones, en cambio, basta que exista una “expectativa empírica” para seguirlas.

Dilema del Prisionero iterado

Muchos filósofos sostienen que las personas cooperan porque piensan a futuro: si hoy ayudo, mañana me ayudarán a mí. El mismo Hume (1739, p. 698) describió esta tendencia: “Aprender a prestar servicios a otra persona sin sentir por ella ningún afecto real, porque preveo que ésta me devolverá el favor esperando que

²⁴ La diferencia entre juegos de coordinación y juegos de motivación mixta fue desarrollada por Schelling (1960). Ese libro fue pionero en aplicar la teoría de juegos a contextos políticos y militares, y le valió el premio Nobel en 2005.

yo realice otro de la misma clase”. Un siglo antes, Hobbes (1651) propuso el famoso “argumento del necio”, donde sostiene que incumplir un pacto es auto-destructivo, porque excluye a ese “necio” de otros futuros pactos sociales. En torno a esa idea, distintos autores exploraron qué sucede si el Dilema del Prisionero se juega repetidamente: ¿sigue siendo conveniente traicionar?

En términos analíticos, puede probarse que, si la repetición es finita, y los agentes saben cuántas rondas tiene el juego, el equilibrio será traicionar siempre (Luce & Raiffa 1957, p. 99).²⁵

Axelrod (1984) tuvo la interesante idea de investigar el problema usando simulaciones (un método bastante innovador para su época). Su proyecto es investigar qué pasa si ponemos a distintos jugadores a jugar repetidamente el Dilema del Prisionero, de modo que no solo pueden tener una estrategia como Cooperar o Traicionar, sino estrategias más complejas como “Cooperar 2 veces, Traicionar 2 veces, y así sucesivamente”, o incluso estrategias *relacionales*, que dependen de lo que el otro haga. Se puede organizar un “torneo” entre distintas estrategias, y analizar a cuál le va mejor. Los partidos obviamente no son infinitos, pero los agentes no saben cuántas rondas tendrá cada partido (esto evita estrategias del estilo “Cooperar hasta la última ronda”).

Algunas estrategias utilizadas fueron las siguientes:

- Vengativo: Cooperar, mientras el otro no me traicione. Si el otro traiciona, después traiciono siempre.
- “Ojo por ojo” o Reciprocador: Primero cooperar. Después, imitar lo que hizo el otro en la última jugada.
- Traidor: Traicionar siempre.

Naturalmente, el Traidor no puede *perder* un partido (como mucho, puede empatar). Sin embargo, en un “torneo” uno puede evaluar cuántos puntos hace cada uno, y analizar en forma más global cuáles estrategias terminan obteniendo más puntos. Así como en el fútbol, no solo importan las victorias, sino también los goles que cada uno hace.

²⁵ La prueba es algo compleja, y presupone el concepto de equilibrio perfecto en subjugos, que veremos más adelante.

No es difícil elaborar un “torneo” entre distintos jugadores con distintas estrategias. Armar un “partido” de dos es sencillo: en cada ronda, un jugador puede jugar (Coopera o Traiciona) teniendo en cuenta lo que hizo su oponente en rondas anteriores (excepto en la primera ronda, donde no tiene información previa). De acuerdo con lo que hacen los jugadores en esa ronda, se computan los puntajes (siguiendo la tabla de puntajes definida anteriormente). Al final del partido sumamos los puntos de cada ronda y obtenemos el puntaje final de cada jugador. El jugador con más puntos gana ese partido.

Para recrear las ideas de Axelrod, podemos pensar un torneo de los siguientes jugadores:

- Cooperador: siempre coopera.
- Reciprocador: empieza cooperando, pero en las rondas siguientes hace lo que hizo el oponente.
- Traidor: siempre traiciona.

Ahora veremos cómo funcionó el torneo, jugando rondas de 80 partidos uno contra uno:

	Cooperador	Reciprocador	Traidor
Cooperador	160-160		
Reciprocador	160-160	160-160	
Traidor	240-0	82-79	80-80

Vemos que, si bien el traidor gana siempre (aunque en dos casos, con bajo puntaje), y el cooperador puede ser destruido por el traidor, el reciprocador tiene puntajes altos en dos de los tres casos, y un puntaje razonable en el otro caso.

Axelrod considera que el reciprocador es “el ganador” en el juego iterado del prisionero; de hecho, esta estrategia tuvo un resultado notable en los “torneos” organizados por el autor.

El resultado fue usado por Axelrod como una posible explicación de por qué existe la cooperación humana. Si la población fuera

reciprocadora, va a sacar ventajas de la cooperación, y va a poder anular a eventuales traidores. En cambio, si fuera traidora, nunca obtendrá los beneficios de la cooperación. Y si fuera simplemente cooperadora, no podrá neutralizar a eventuales traidores, que finalmente sacarán toda la ganancia. Por eso, desde un punto de vista “evolucionista”, la estrategia reciprocadora va a terminar imponiéndose a nivel poblacional.

Esta idea de Axelrod también ha sido utilizada para argumentar que la cooperación no requiere de sentimientos altruistas, porque puede justificarse de forma egoísta (como sugiere la cita de Hume al principio de esta sección). Otros autores sostienen que el Dilema del Prisionero Iterado muestra, contra Hobbes, que el estado de naturaleza no es siempre un estado de caos, y que es posible llegar a equilibrios razonables sin un “Leviatán”.

De hecho, a diferencia de lo que sucede en el Dilema del Prisionero de una ronda (o de una cantidad finita de rondas), en la versión iterada indefinidamente dos agentes cooperativos que aplican “Ojo por Ojo” sí están en equilibrio. Podríamos pensarlo de este modo: si soy Reciprocador y el otro agente es Reciprocador, obtengo 2 en cada ronda. No saco ventaja cooperando más, porque de hecho coopero en todos los pasos, y tampoco saco ventaja traicionando ocasionalmente, porque cada traición será castigada en el paso siguiente.²⁶

Las ideas de Axelrod fueron cuestionadas por autores posteriores. Rapoport *et al.* (2015) indican que el reciprocador no es el “ganador” natural de un torneo: quién será el “ganador” depende de quién juega el torneo. Por otro lado, Binmore (1998, p. 186) observa que una comunidad de reciprocadores no es el *único* equilibrio posible para el Dilema del Prisionero iterado (veremos otros posibles equilibrios en los ejercicios); por eso no podemos evitar el problema de la elección entre equilibrios, mencionado anteriormente.

²⁶ Para una prueba completa véase Binmore (2007), cap. 11. Técnica-mente también existe una diferencia entre repetición infinita y repetición indefinida; en esta instancia, los tratamos como equivalentes.

Ejercicios

1. Hacer un torneo del Dilema del Prisionero Iterado con estas características, haciendo que cada jugador juegue un partido contra cada uno de los otros jugadores (pero no contra sí mismo):

- Cada partido tiene 4 jugadas (iteraciones).
- Un jugador es Alternador: empieza cooperando, después traiciona, después coopera, etc.
- Un jugador es Vengativo: coopera, pero si lo traicionan, de ahí en adelante siempre traiciona.
- Otro jugador es Reciprocador: empieza cooperando, después hace lo mismo que hizo el otro jugador en la última jugada.

2. En el Dilema del Prisionero iterado indefinidamente:

- a. ¿(Traicionar Siempre, Traicionar Siempre) es un equilibrio?
- b. ¿(Cooperar Siempre, Cooperar Siempre) es un equilibrio?
- c. ¿(Vengativo, Vengativo) es un equilibrio?

Parte D: Paradojas y experimentos

Del mismo modo que sucede con la teoría de la decisión estándar, la teoría de juegos enfrentó cuestionamientos a partir de ciertas paradojas y de resultados experimentales.

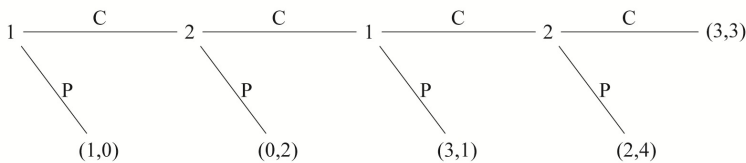
Podríamos ver al Dilema del Prisionero como una paradoja respecto al equilibrio de estrategias dominantes (el concepto más sólido de la teoría de juegos). Sin embargo, como vimos, autores como Binmore sostienen que no es una verdadera paradoja. Incluso los resultados empíricos muestran que, en un juego de una sola ronda, las personas suelen traicionar, tal como indica la teoría de juegos (Andreoni & Miller 1993).

El *Dilema del Viajero* (Basu 1994) es una interesante paradoja para los equilibrios de Nash. En este escenario, dos personas deben decidir cuánto dinero van a obtener por un bien. Cada uno debe decir un número entre \$2 y \$100. Si dicen ambos el mismo número, obtienen ese dinero. Si uno dice menos que el otro, el que dijo menos obtiene ese dinero, con un premio de \$2 (por honestidad). El que dijo el número más alto obtiene el número más bajo, menos un castigo de \$2. Por ejemplo, si a dice \$40 y b dice \$90, a recibe \$42 y b recibe \$38.

Lo curioso del escenario es que, si bien no hay estrategias dominantes, el único equilibrio de Nash es que ambos digan \$2. Porque (a) si ambos dicen el mismo número n , ambos querrán haber dicho un número menos, y obtener $n+1$; y (b) si dicen números distintos, el que dijo el número mayor preferiría haber dicho el número menor n , y obtener n en vez de $n-2$. Si dicen \$2 no pueden arrepentirse, porque no podían decir un número menor (y si decían un número mayor, obtenían \$0).

El equilibrio (2, 2) es muy poco intuitivo: si ambos decían \$100, recibirían \$100. De hecho, en escenarios experimentales, los agentes dicen números altos (Basu *et al.* 2011). Este juego puede usarse como crítica al equilibrio de Nash, en su rol descriptivo (porque los agentes no eligen así), y también normativo (porque no parece una forma razonable de elegir). Un modo de responder a esta paradoja es decir que los agentes piensan el juego de forma más “indefinida”: las opciones son “decir un número alto” y “decir un número bajo”, y con esa matriz, que ambos digan números altos es un equilibrio de Nash (Basu 1994).

Así como el equilibrio de Nash debe enfrentar la paradoja del viajero, el método de inducción hacia atrás se enfrenta a otras paradojas, como el *juego del ciempiés* (Rosenthal 1981):



En este juego, los agentes tienen que decidir si “continuar” o “parar”; si ambos continúan hasta el final, ganan (3, 3). Aquí, el único equilibrio por inducción hacia atrás es que el jugador 1 se retire en la primera ronda. Sin embargo, eso lo deja con una utilidad de 1. Un resultado muy poco satisfactorio.

Binmore (1987) sostiene, a partir de esta paradoja, que el método de inducción hacia atrás no debe identificarse *siempre* con la elección racional. En versiones “largas” del juego del ciempiés (por ejemplo, con 100 pasos), la mínima probabilidad de que el

otro jugador no se comporte “racionalmente”, me llevará a continuar, al menos en las primeras rondas.

Los resultados experimentales arrojan más luz sobre esta idea. Las personas, en general, deciden “continuar” en las primeras rondas. Pero como muestran Rapoport *et al* (2003), esto depende de los premios en juego. Si los premios monetarios son altos, las personas van a estar más cerca del equilibrio, es decir, de parar en la primera ronda. Esto rara vez puede evaluarse en un laboratorio, donde no hay dinero para dar premios altos. Asimismo, si el juego es corto, las personas adultas suelen razonar de acuerdo con la inducción hacia atrás. Esto no necesariamente se correlaciona con el éxito: en el juego del ciempiés corto, los niños obtienen más ganancia que los adultos (Brocas & Carrillo 2025). Estos resultados parecen confirmar la idea de que, si el juego es corto y los estímulos monetarios son altos, la inducción hacia atrás es un buen modelo del comportamiento racional.

Además de la paradoja del viajero y el juego del ciempiés, otros argumentos se han presentado para evaluar los métodos y principios de la teoría de juegos. Actualmente existe la Teoría de Juegos Conductual, que estudia la (conflictiva) relación entre la teoría de juegos y el comportamiento humano. Como pudimos ver, los resultados paradójicos son utilizados por algunos para rechazar los métodos principales de la teoría de juegos (o al menos, rechazar su aplicabilidad general), mientras que otros intentan hacer encajar la teoría con los resultados, a partir de teorías refinadas, o de lecturas alternativas de los escenarios.

Ejercicio

En el *juego del ultimátum*, el Agente 1 debe decidir cómo reparte 100 dólares (en billetes de \$1) entre él y el Agente 2. El Agente 2 tiene el poder de rechazar la oferta, y dejar a ambos sin nada.

1. ¿Cómo va a repartir el dinero el Agente 1 según la teoría de juegos (más precisamente, según el razonamiento por inducción hacia atrás)?

2. ¿Cómo crees que se comportaría la gente en casos experimentales?

Parte E: Juegos de información incompleta

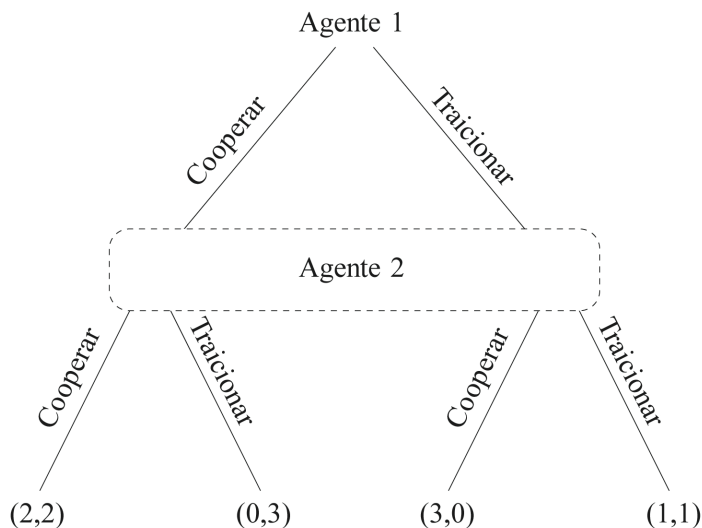
Hasta aquí, discutimos juegos *estratégicos* y juegos *dinámicos*. Los juegos dinámicos que vimos, sin embargo, tienen una característica particular: se trata de juegos de información *perfecta*. En esta clase de juegos, los agentes siempre saben las jugadas que hace el otro agente. Un tipo de juego de información perfecta en el mundo real es el Tres en Raya (*tic-tac-toe*)²⁷: cada jugada es visible, porque los jugadores deben escribir “X” o una “O” en la grilla de 3 por 3. Otros juegos podríamos entenderlos como de información *incompleta*: un ejemplo muy sencillo es el “piedra, papel o tijera”, porque tenemos que jugar sin saber qué estrategia eligió el otro jugador.

En cierto sentido, podemos entender los juegos “estratégicos” vistos antes en este capítulo (el Dilema del Prisionero, “Piedra, papel o tijera”, etc.) como juegos de información incompleta, donde cada jugador debe elegir qué hacer sin saber qué hizo el otro jugador.

En teoría de juegos, la información incompleta se suele escribir usando “conjuntos de información” (*information sets*). Un conjunto de información es un conjunto de nodos, tales que el agente que debe tomar la decisión no puede saber si está en un nodo o en el otro. Todos los juegos de información incompleta tienen, al menos en cierto lugar, un conjunto de información.

Por ejemplo, podríamos escribir el dilema del Prisionero como un juego dinámico con información incompleta:

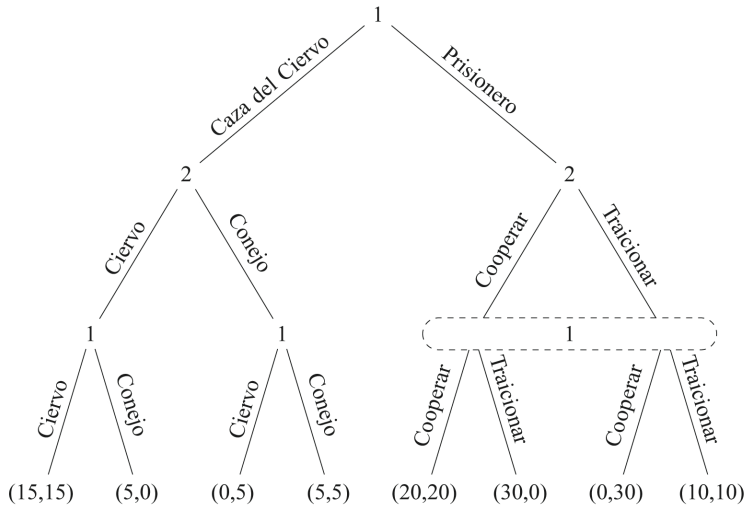
²⁷ Este juego tiene distintos nombres, como “gato” (Chile), “tatetí” (Argentina), “michi” (Perú) y “tiqui” (Colombia).



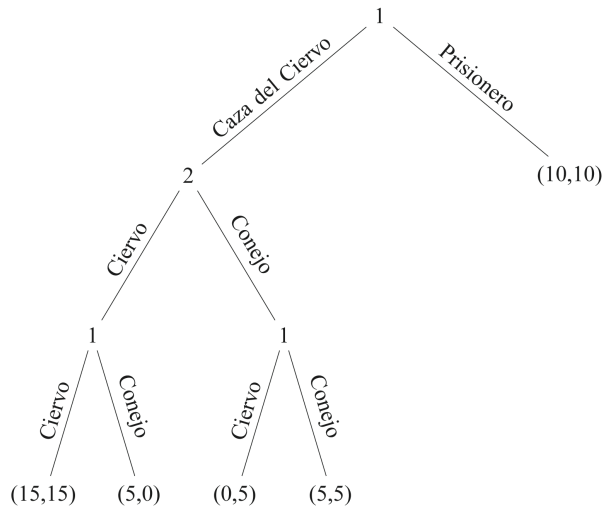
El conjunto de información es señalado con la línea de puntos. Esto indica que el Agente 2 no puede saber en cuál de esos nodos está realmente.

Los juegos de información incompleta también tienen un tipo de “solución”, que combina la inducción hacia atrás y los equilibrios de Nash. Para eso es necesario introducir el concepto de “subjuego”. Conceptualmente, un *subjuego* es un “juego dentro de un juego”. Es decir, es un subconjunto del árbol original que tiene (i) un nodo central donde empieza el subjuego, y (ii) nodos terminales que indican los resultados del juego. Naturalmente, un subjuego no puede “cortar” elementos dentro de un conjunto de información; por lo tanto, un subjuego siempre va a incluir, no cortar, los conjuntos de información.

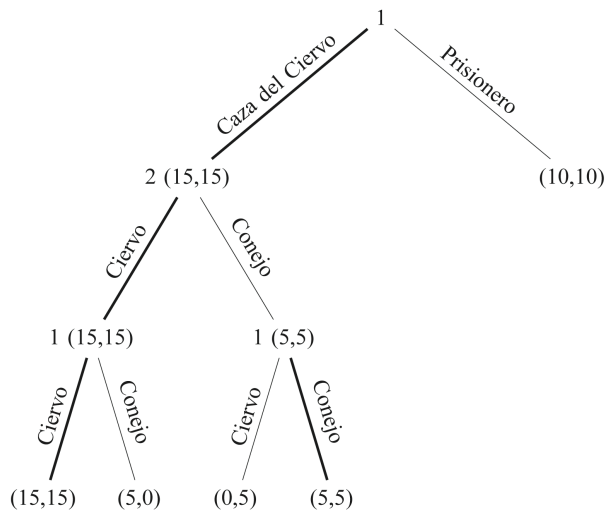
Supongamos, por ejemplo, que tenemos un juego de información imperfecta, donde un agente debe decidir si jugar el Dilema del Prisionero (con información incompleta) o la Caza del Ciervo de forma secuencial (es decir, con información perfecta). Usaremos una escala de utilidad más tentadora para el Dilema del Prisionero. El árbol del juego se vería así:



¿Cómo podríamos resolver este juego de información incompleta? Del lado izquierdo, parece fácil: usaremos inducción hacia atrás. Pero del lado derecho no funciona nuestra receta original, porque el Agente 1 no sabe en cuál nodo está al tomar la última decisión. Incluso si pensamos que su última decisión será Traicionar, no sabemos si su última decisión resultará en $(30, 0)$ o $(10, 10)$. Por eso, para estos casos apelamos al *Equilibrio de Nash*. Es decir, cuando un subjuego dentro de un juego se puede resolver mediante un Equilibrio de Nash, podemos suponer que los agentes serán racionales y llegarán a este equilibrio. Como hemos mostrado repetidas veces, en el Dilema del Prisionero, el único equilibrio de Nash es la traición de ambos. Entonces podemos pensar el juego de este modo:



Ahora sí, podemos hacer la inducción hacia atrás como corresponde. Como veremos, el Agente 1 va a preferir jugar a la Caza del Ciervo de forma secuencial.

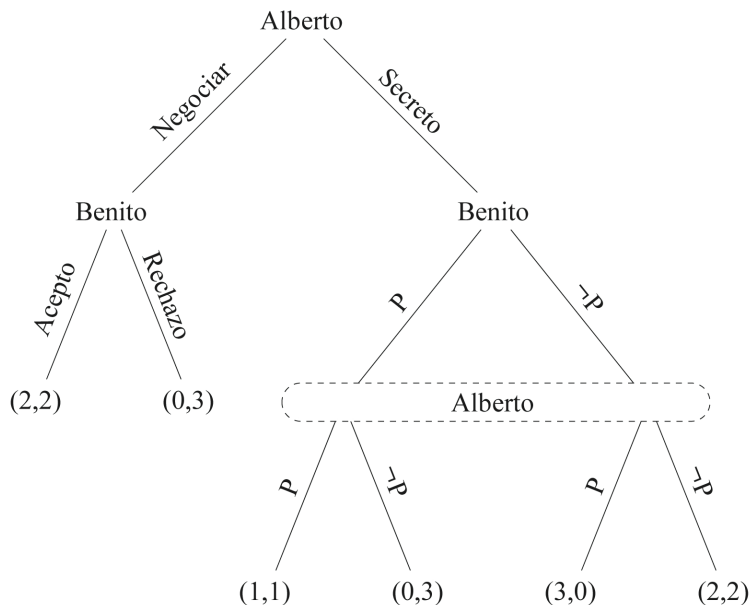


La palabra técnica para esta solución es “equilibrio perfecto en subjugos” (Selten 1975), porque consiste en resolver el juego a partir de la solución de todos sus subjugos. Es una extensión del equilibrio por inducción hacia atrás: el equilibrio perfecto en subjugos, si la información es completa, equivale a la inducción hacia atrás (cada paso de la inducción hacia atrás resuelve un subjuego). Aquí también sucede que todo equilibrio perfecto en subjugos será un equilibrio de Nash (aunque no al revés).

Por último, puede notarse algo interesante respecto a la comparación entre la Caza del Ciervo y el Dilema del Prisionero. La Caza del Ciervo tiene *dos* equilibrios de Nash: que ambos vayan a cazar el ciervo (15,15) y que ambos vayan a cazar conejos (5,5). De modo tal que, sin información contextual, no podemos saber si los agentes racionales harán una cosa o la otra. El Dilema del Prisionero, en cambio, tiene solo un equilibrio (la traición mutua). En el Dilema del Prisionero, jugar de forma secuencial no altera la solución correcta: de modo secuencial o simultáneo, ambos agentes van a suponer que los otros agentes van a traicionar, y en consecuencia van a traicionar. En la Caza del Ciervo, que es un juego cooperativo, jugar de forma secuencial permite que el segundo agente sepa qué hizo el primero, y motiva una acción cooperativa. Autores como Skyrms (2002) exploraron la importancia de la *información* en estos juegos cooperativos.

Ejercicio

Alberto y Benito son dos hermanos políticos que se alternan el gobierno de la misma provincia hace décadas. En la nueva elección provincial, parte del pueblo les pide que se retiren. Si ninguno de los dos presenta candidaturas, esto será visto como un acto honrado. Sin embargo, si uno de los dos presenta candidaturas (y el otro no), seguramente va a ganar, y destruirá al otro. Y si ambos se presentan por separado, van a perder. Alberto se plantea qué hacer: ¿juega secretamente o negocia con el hermano una lista conjunta? Resuelva esta situación a partir de este juego de información incompleta (donde Alberto es el Agente 1 y Benito el Agente 2):



*Parte F: Estrategias mixtas y probabilidades

La última parte de este capítulo sobre teoría de juegos va a extender un poco el enfoque respecto a lo visto anteriormente. En las secciones anteriores señalamos que ciertos juegos no tienen un equilibrio de Nash entendido de la forma usual. Por ejemplo, en el “Piedra, papel o tijera”, no hay ninguna situación en que *ambos* jugadores estén conformes con lo jugado: si uno ganó, el que perdió desearía haber jugado otra cosa, y si empataron, ambos desearían haber jugado otra cosa.

Sin embargo, existe un posible equilibrio de Nash para estos juegos. Para entenderlo hace falta combinar la teoría de juegos con la idea de maximización de utilidad, proveniente de la teoría de la decisión. Una forma de jugar la “mejor” estrategia posible, cuando no existe tal cosa literalmente, es usar estrategias *mixtas*. Por ejemplo, supongamos que yo sé que el Agente 2 va a jugar de este modo: primero tira un dado, y luego juega Tijera si sale 1 o 2, Piedra si sale 3 o 4, y Papel si sale 5 o 6. Si yo (Agente 1)

supiera que el otro agente va a jugar esto, ¿qué estrategia me conviene adoptar? Podría jugar Tijera y ganaría con un tercio de probabilidades (cuando el otro juega Papel). Lo mismo sucedería si juego Papel, o también si juego Piedra. En resumen, no va a importar mucho qué estrategia pura decida usar finalmente.

Ahora supongamos que yo (Agente 1) también tiro un dado (no el mismo dado que tira el otro jugador), y luego juego Tijera si sale 1 o 2, Piedra si sale 3 o 4, y Papel si sale 5 o 6. ¿Me conviene usar esta estrategia “mixta”? Veamos. La utilidad esperada de esta acción (llamémosla “A”) podríamos calcularla así:

$$U_1(A) = 1/3 \times U_1(\text{Piedra}) + 1/3 \times U_1(\text{Tijera}) + 1/3 \times U_1(\text{Papel})$$

Dado que, como dijimos anteriormente, $U_1(\text{Piedra}) = U_1(\text{Tijera}) = U_1(\text{Papel})$, entonces podemos deducir que $U_1(A) = U_1(\text{Piedra}) = U_1(\text{Tijera}) = U_1(\text{Papel})$. De modo tal que esta estrategia mixta no es mejor, pero tampoco es peor, que las estrategias puras.

Ahora podemos interpretar este fenómeno de forma más global. Supongamos que el otro agente adopta una estrategia randomizadora, como tirar el dado, y yo también. En ese caso, *ambos estaríamos haciendo lo mejor posible, suponiendo lo que hace el otro*. Por lo tanto, estamos en un equilibrio de Nash. De este modo, incluso los juegos que no tienen equilibrios en estrategias puras pueden tener equilibrios en estrategias mixtas.

Para encontrar un equilibrio de estrategias mixtas a partir de una matriz (es decir, a partir de un juego en forma estratégica), deberíamos primero descartar las estrategias dominadas. Una vez que hago eso, puedo determinar estrategias mixtas para los jugadores, que resulten en un equilibrio de Nash. Una forma sencilla de encontrar equilibrios de estrategias mixtas es hallar escenarios donde, como jugador, nos resulte indiferente hacer cualquier opción (como vimos en el caso del “Piedra, papel o tijera”).

Por ejemplo, supongamos que el juego estratégico tiene esta forma:

	C	D
A	4, 2	5, 5
B	3, 4	7, 3

Este juego no tiene estrategias dominadas, ni tampoco tiene equilibrios de Nash en estrategias puras. Entonces podemos obtener un equilibrio en estrategias mixtas.

Primero debemos determinar en qué casos al Agente 1 le daría lo mismo hacer A que hacer B. Sería un caso donde $U_1(A) = U_1(B)$. Supongamos que p fuera la probabilidad de que el Agente 2 haga C, y $(1 - p)$ sea la probabilidad de que el Agente 2 haga D. Entonces:

$$U_1(A) = p \times 4 + (1 - p) \times 5$$

$$U_1(B) = p \times 3 + (1 - p) \times 7$$

Para que ambas opciones le sean indiferentes al Agente 1, podemos calcular p mediante una ecuación, asumiendo que $U_1(A) = U_1(B)$:

$$p \times 4 + (1 - p) \times 5 = p \times 3 + (1 - p) \times 7$$

$$4p + 5(1 - p) = 3p + 7(1 - p)$$

$$p = 2(1 - p) \Leftrightarrow p = 2 - 2p \Leftrightarrow 3p = 2 \Leftrightarrow p = 2/3$$

Ya tenemos entonces un lado del equilibrio. El Agente 2 hace C con probabilidad $2/3$ y D con probabilidad $1/3$.

Ahora hace falta obtener el otro lado, aunque es simétrico. Sea q la probabilidad de que el Agente 1 haga A, y $(1 - q)$ la probabilidad de que haga B:

$$U_2(C) = 2q + 4(1 - q)$$

$$U_2(D) = 5q + 3(1 - q)$$

Para que al Agente 2 le sea indiferente hacer C o D, la probabilidad q podemos calcularla así, asumiendo que $U_2(C) = U_2(D)$:

$$2q + 4(1 - q) = 5q + 3(1 - q)$$

$$1 - q = 3q \Rightarrow 1 = 4q \Rightarrow q = 1/4$$

Es decir, ahora ya encontramos nuestro equilibrio de estrategias mixtas: la estrategia del Agente 1 es hacer A con $1/4$ y B con $3/4$, mientras que la estrategia del Agente 2 es hacer C con $2/3$, y hacer D con $1/3$.

Un resultado importante en la teoría de juegos es que cualquier juego en forma estratégica tiene al menos un equilibrio de Nash, ya sea en forma pura o en forma mixta. Esa es, de hecho, la principal contribución de la obra de Nash.

Antes de terminar esta sección, hago dos observaciones. En primer lugar, el hecho de que un juego tenga equilibrios “puros” es compatible con que también tenga equilibrios “mixtos” (veremos un caso en los ejercicios). Por otro lado, un juego podría tener múltiples equilibrios mixtos; esto será más común en juegos más complejos, con más jugadores o más estrategias. De hecho, encontrar los equilibrios de un juego puede ser muy difícil. Aquí solamente desarrollamos una aproximación a este concepto para juegos sencillos de dos jugadores.

Ejercicio

1. Encuentre un equilibrio de estrategias mixtas para este juego

	C	D
A	5, 4	3, 5
B	4, 6	4, 3

2. Encuentre un equilibrio de estrategias mixtas para el juego de Hablar inglés o español (conocido como “Guerra de los Sexos”):

	Hablar español	Hablar inglés
Hablar español	2, 1	0, 0
Hablar inglés	0, 0	1, 2

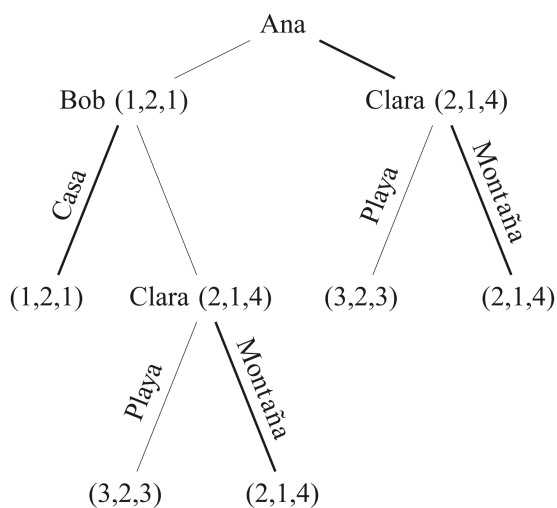
Soluciones para el capítulo 3

PARTE A

1. Queda BF.
2. AE y DH.
3. Este juego tiene dos equilibrios de Nash: AC y BD. Sin embargo, si borramos las estrategias débilmente dominadas, debemos borrar B y también D. Nos queda solo AC. Esto muestra que el borrado de estrategias débilmente dominadas puede borrar algunos equilibrios de Nash.

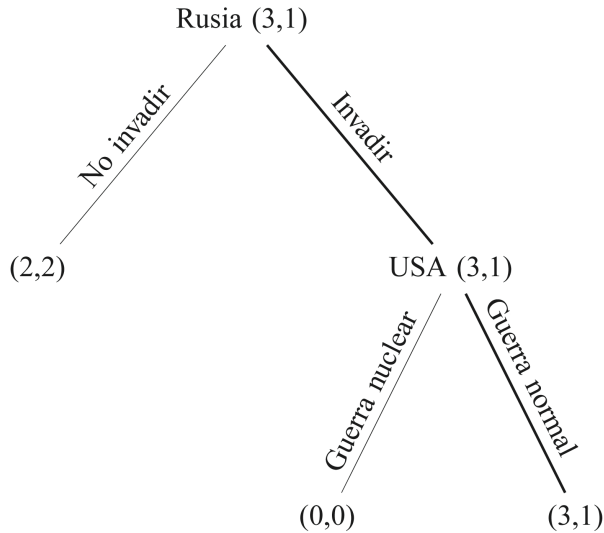
PARTE B

1. La solución es dejar que decida Clara:



- 2.

a. La solución es (Invadir, Guerra normal):



b. La forma estratégica asociada es:

	Guerra nuclear	Guerra normal
No invadir	2, 2	2, 2
Invadir	0, 0	3, 1

Hay dos equilibrios de Nash: (No invadir, Guerra nuclear) y (Invadir, guerra normal).

c. Si bien los equilibrios por inducción hacia atrás son siempre equilibrios de Nash, algunos equilibrios de Nash (como <No invadir, Guerra nuclear>) no son equilibrios de inducción hacia atrás. Aquí, el equilibrio incluye una guerra nuclear autodestructiva, que es muy poco razonable para Estados Unidos. Estos casos se conocen como “amenazas no creíbles”, un concepto popularizado por Schelling (1960). Schelling insistió en que la racionalidad requiere tomar en cuenta las amenazas no-creíbles, y esto

sirvió para fundamentar la estrategia de la “disuasión nuclear”, predominante en la Guerra Fría.

PARTE C

1.

Alternador vs. Vengativo:

2, 2 | 3, 0 | 0, 3 | 1, 1

Resultado: Empate 6-6

Alternador vs. Reciprocador:

2, 2 | 3, 0 | 0, 3 | 3, 0

Resultado: gana Alternador 8-5 (si había otra ronda, empatan)

Reciprocador vs. Vengativo:

2, 2 | 2, 2 | 2, 2 | 2, 2

Resultado: empate 8-8

2.

a. (Traicionar Siempre, Traicionar Siempre) *sí* es un equilibrio porque si el otro traiciona siempre, no saco ningún provecho en cooperar a veces; me conviene traicionar siempre.

b. (Cooperar Siempre, Cooperar Siempre) *no* es un equilibrio, porque si el otro coopera siempre, yo sacaré ventaja traicionando siempre (o incluso, traicionando a veces).

c. (Vengativo, Vengativo) *sí* es un equilibrio. Los agentes van a cooperar siempre. Y uno no sacará ventajas en traicionar en una ronda: solo generaría la traición eterna del otro jugador.

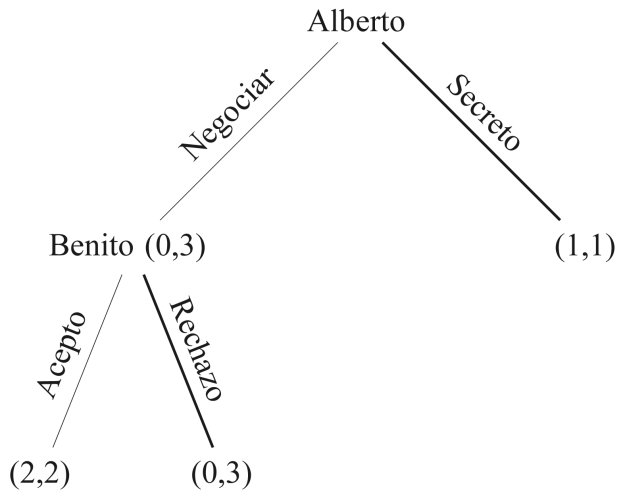
PARTE D

1. Si los agentes son egoístas y ambiciosos, el único equilibrio (por inducción hacia atrás) es que el Agente 1 le ofrezca \$1 al Agente 2 (es decir, 1% del total), y que el Agente 2 lo acepte. Después de todo, el Agente 2 debe elegir entre \$1 o nada.

2. Sin embargo, en experimentos, los agentes suelen ofrecer alrededor del 40% del pozo. Y las ofertas muy bajas son rechazadas en la mitad de los casos (Bicchieri 2005, p. 104). Esto suele explicarse bajo la idea de que las personas, además del apetito por dinero, tienen sentimientos de altruismo o justicia, y también sentimientos de venganza hacia los agentes injustos.

PARTE E

La solución es jugar en secreto, asumiendo que Benito rechazaría la negociación:



PARTE F

1. El juego está representado por esta matriz:

	C	D
A	5, 4	3, 5
B	4, 6	4, 3

Supongamos que la probabilidad de que el agente 2 haga C es p , mientras que la probabilidad de que haga D es $(1 - p)$. Entonces:

$$\begin{aligned}U_1(A) &= 5p + 3(1 - p) \\U_1(B) &= 4p + 4(1 - p) = 4\end{aligned}$$

Entonces $U_1(A) = U_1(B)$ significa que:

$$5p + 3(1 - p) = 4 \Leftrightarrow 5p + 3 - 3p = 4 \Leftrightarrow 2p = 1 \Leftrightarrow p = \frac{1}{2}$$

Por otro lado, siendo q la probabilidad de que el agente 1 haga A, y $1-q$ la probabilidad de que el agente 1 haga B, podemos definir:

$$\begin{aligned}U_2(C) &= 4q + 6(1 - q) \\U_2(D) &= 5q + 3(1 - q)\end{aligned}$$

Entonces si $U_2(C) = U_2(D)$ podemos saber que:

$$\begin{aligned}4q + 6(1 - q) &= 5q + 3(1 - q) \\3(1 - q) &= q \Leftrightarrow 3 - 3q = q \Leftrightarrow 3 = 4q \Leftrightarrow q = \frac{3}{4}\end{aligned}$$

Obtenemos un equilibrio donde 1 hace A con $\frac{3}{4}$ y B con $\frac{1}{4}$, mientras que 2 hace C con $\frac{1}{2}$ y D con $\frac{1}{2}$.

2.

Con un razonamiento similar al del ejercicio anterior, obtenemos que el Agente 1 hace A con probabilidad $\frac{1}{3}$ y B con probabilidad $\frac{2}{3}$; y el Agente 2 hace C con probabilidad $\frac{2}{3}$ y D con probabilidad $\frac{1}{3}$. Lo curioso en este caso es que en cualquier acto (por ejemplo, A), los agentes obtendrán una utilidad esperada de $\frac{2}{3}$; les convendría coordinar en algún equilibrio puro antes que usar esta estrategia mixta, porque se garantizan obtener al menos una utilidad de 1.

CAPÍTULO 4: ELECCIÓN SOCIAL

Parte A: Antecedentes y teoría del voto

Desde el origen de la civilización, los grupos han tenido que tomar decisiones colectivas, y para eso diseñaron distintos mecanismos. La *teoría de la elección social* estudia estos mecanismos y sus propiedades. Un método muy popular para decidir grupalmente si hacer determinada acción es el siguiente:

- 1. Los que están a favor levantan la mano.
- 2. Contamos la cantidad.
- 3. Si la cantidad es mayor a la mitad de las personas presentes, realizamos la acción.

Este método, conocido como *voto por mayoría*, es útil para decidir entre dos opciones. Pero los mecanismos de voto a veces requieren decisiones más complejas, como ordenar las preferencias entre opciones. Y ahí aparecen muchos problemas. Un método apto para elecciones entre distintas opciones fue propuesto por Condorcet (1785). Supongamos que una comisión de tres personas (las llamamos sin mucho esfuerzo “1”, “2” y “3”) debe elegir entre cuatro candidatos para un puesto: Juan, Alicia, Pablo y Matilda. Luego de revisar el currículum y realizar entrevistas, los miembros del comité tienen las siguientes preferencias (las mejores opciones son las que están más arriba):

1	2	3
Juan	Pablo	Matilda
Alicia	Alicia	Pablo
Pablo	Juan	Alicia
Matilda	Matilda	Juan

Supongamos que el objetivo de la comisión no es solo elegir *un* ganador, sino también determinar un segundo puesto (en caso de que el ganador no pueda o no quiera aceptar). Es decir, la comisión no debe elegir un candidato, sino generar un ranking nuevo.

¿Cómo debería proceder? El método de Condorcet nos sirve para resolver este tipo de casos. Es una versión del voto por mayoría que puede aplicarse cuando hay varias opciones.

(Criterio de Condorcet)²⁸

En el orden grupal, $x \geq y$ sii $x \geq_i y$ para la mitad o más de los agentes i .

El criterio de Condorcet nos arrojará el siguiente resultado:

Alicia > Juan, por el voto de 2 y 3
Pablo > Alicia, por el voto de 2 y 3
Pablo > Juan, por el voto de 2 y 3
Todos > Matilda, por el voto de 1 y 2

Entonces el ranking colectivo queda de este modo:

Pablo > Alicia > Juan > Matilda

Este método para generar rankings colectivos es muy útil, y aun actualmente suele usarse en distintos comités.

El problema es que, como descubrió Condorcet, este método intuitivo se ve afectado por paradojas. La paradoja más importante fue conocida como “Paradoja del Voto” o “Paradoja de Condorcet”. Supongamos que hay que elegir entre tres eventos (por ejemplo, tres gustos de helado). Y hay tres agentes votando, cuyas preferencias son las siguientes (llamaremos “perfil” a un conjunto de preferencias individuales de distintos agentes):

1	2	3
a	b	c
b	c	a
c	a	b

¿Qué ranking tendrá el grupo, a partir de este perfil?

²⁸ Terminológicamente, usaremos subíndices del tipo \geq_i para hablar de preferencias de un individuo i , y al no usar subíndices nos referimos a las preferencias grupales.

Como veremos, el método de Condorcet resulta problemático en este caso particular, porque nos da el siguiente resultado:

$a > b$ por el voto de agentes 1 y 3
 $b > c$ por el voto de agentes 1 y 2
 $c > a$ por el voto de agentes 2 y 3

Se genera entonces un ciclo de preferencias colectivas:

(Ciclo) $a > b, b > c, c > a$

Es otras palabras, las preferencias “grupales” formadas mediante el método de Condorcet generan un ciclo, algo que generalmente prohibimos para cualquier escala de preferencia. Individuos consistentes generan un grupo inconsistente.

¿Cuánto revela la Paradoja de Condorcet sobre la naturaleza del voto? Esperamos que el lector pueda tener un juicio propio hacia el final de este capítulo. Como resulta obvio, el método de Condorcet no es el único método de decisión colectiva. Más adelante, veremos otros métodos de voto que no generan ciclos, como el de Borda. Sin embargo, estos métodos tienen otros problemas.

Ejercicios

1. Encuentre el ranking de Condorcet en este perfil de votos:

1	2	3
a	b	b
b	c	a
c	a	c

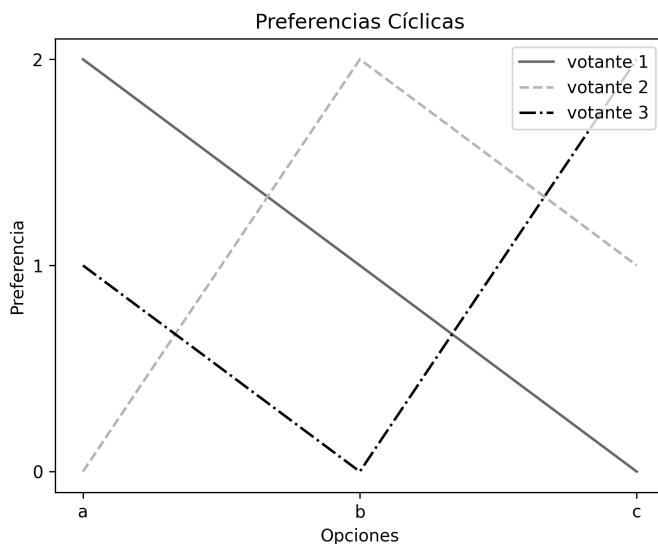
2. Siglos antes de Condorcet, el filósofo español Ramón Lull (1232-1316) ideó un método de elección “de a pares”, publicado en su texto *De Arte Electionis*.²⁹ El método de Lull es el siguiente: primero se vota la opción (x, y) por mayoría; luego,

²⁹ Puede encontrarse una traducción crítica del texto de Lull en Barenstein (2013).

quien gane esa votación se enfrenta a z ; luego, quien gane esa votación se enfrenta a w ; y así sucesivamente, hasta encontrar el último ganador. El método de Lull no produce ciclos, pero tiene otro problema, ¿cuál es?

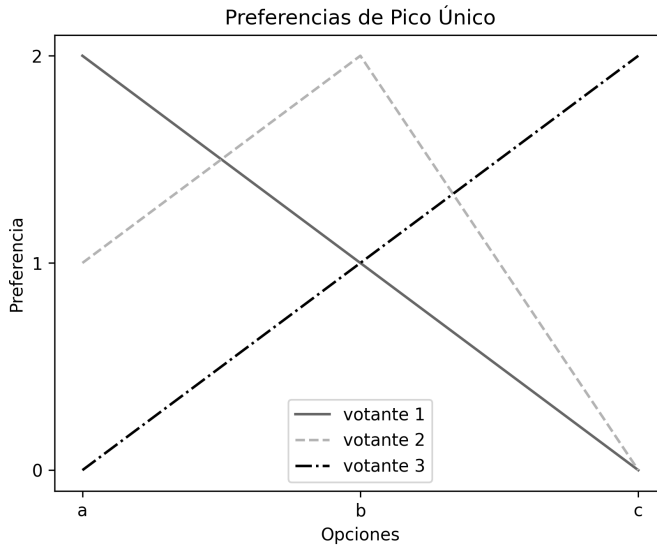
***Parte B: Blindar el voto por mayoría**

Algunos autores notaron que lo que falla en la paradoja de Condorcet es que no hay acuerdos básicos: las personas desacuerdan sobre qué es lo mejor y qué es lo peor.



Duncan Black (1948) formalizó esta idea en términos gráficos. Según este autor, el problema con los perfiles como el de la Paradoja de Condorcet es que no hay un “pico único”. Es decir, no hay forma de dibujar los votos de los agentes $\{1,2,3\}$ de modo tal que cada agente tenga un único pico. Siempre algún agente va a quedar en “V”. Por ejemplo, en el gráfico anterior (que representa los votos en la Paradoja de Condorcet), el votante 3 tiene dos “picos”. El lector puede comprobar que, sin importar cómo se dibuje el gráfico, algún agente tendrá dos picos.

Estos resultados llevaron a algunos autores a proponer una “restricción de dominio”, es decir, condicionar de antemano los perfiles individuales o grupales que se admiten. Por ejemplo, podríamos aceptar solamente perfiles donde todos los rankings tengan un “pico único” (*single-peaked preferences*). Los perfiles con pico único se ven así:



Podríamos escribir este gráfico como un perfil de votos:

- 1: $a > b > c$
- 2: $b > a > c$
- 3: $c > b > a$

Ahora, usando el método de Condorcet, el perfil resultante será $b > a > c$.

Estos perfiles representan la idea de que hay cierto acuerdo básico entre los agentes sobre qué es (o no es) lo mejor (o lo peor). Por ejemplo, este perfil representa el acuerdo en que b no es la peor opción. Si nos restringimos de antemano a este tipo de perfiles, los ciclos paradójicos se evitan. Este tipo de soluciones, sin

embargo, suelen ser criticadas por dos razones. En primer lugar, su complejidad. Detectar que un perfil no puede representarse con “picos únicos” no es fácil, y se vuelve mucho más complejo cuando aumentan los agentes o las opciones. Una segunda crítica a esta estrategia es su falta de realismo: ¿Qué clase de método de voto pondría condiciones sobre los perfiles admisibles?

Sin embargo, se ha explorado una forma realista de obtener perfiles como estos. Típicamente, una forma de obtener rankings de “pico único” es pedir a los agentes que elijan “aspectos”. Podemos entenderlo con el siguiente ejemplo: la gente seguramente no puede ponerse de acuerdo sobre qué línea del subte (metro) de Buenos Aires es mejor que otra. Por ejemplo, algunos prefieren la línea A (moderna, útil), otros la línea C (conecta dos estaciones centrales), otros la línea D (bella estéticamente), etc. Pero si votan sobre aspectos específicos, como por ejemplo cuál línea de subte es más limpia, los desacuerdos van a minimizarse. Seguramente los agentes tienen ciertos acuerdos sobre qué línea de subte es más limpia, o cuál tiene los mejores vagones, etc. Si se vota sobre *aspectos*, es muy probable que los perfiles tengan un “pico único”. Ahora el problema será cómo integrar estos distintos aspectos en un voto único.

Parte C: Arrow y funciones de bienestar social

Si Condorcet fue el impulsor de la teoría clásica del voto en el siglo XVIII, la moderna Teoría de la Elección Social (*Social Choice Theory*) nació con la obra de Kenneth Arrow (1951). Este autor formalizó las características centrales de los métodos de voto, probó algunos teoremas fundamentales y dio impulso a un nuevo desarrollo científico sobre estos temas.

En el esquema de Arrow, cada agente aporta un ranking de preferencia individual (un elemento ya conocido en teoría de la decisión). Un ranking de preferencia individual R debe ser *completo* (para todo a y b , aRb o bRa), *transitivo* (si aRb y bRc , entonces aRc) y *reflexivo* (para todo a , aRa). Intuitivamente, aRb o $a \geq b$ se lee como “prefiero a sobre b , o me dan igual”. Definimos $a > b$ o aPb (“prefiero a sobre b ”) como $(aRb \ \& \ \neg bRa)$, y definimos $a \sim b$ o aIb (“soy indiferente entre a y b ”) como $(aRb \ \& \ bRa)$.

Un ranking de preferencias podría verse así:

$$\begin{array}{c} c \\ b, d \\ e \\ a, f \end{array}$$

donde $x > y$ se representa como que x está “arriba” de y , mientras que $x \sim y$ se representa como que x está al mismo nivel que y . Hasta aquí no hay nada nuevo respecto a las ideas de preferencia en la teoría de la decisión.

Pero ahora entramos en los aspectos colectivos de la preferencia. Un *perfil* es una n -tupla de rankings de preferencias, $\langle R_1, \dots, R_n \rangle$, que representa las preferencias individuales de un *conjunto* de personas $\{1, \dots, n\}$.

Por ejemplo, un perfil para los agentes 1-5 podría verse así:

1	2	3	4	5
c	c	a	f	a
b, d	a	b	a	c
e	b	c	b	d, e, f
a, f	d	e	c, d	b
	e, f	d	e	
	f			

A partir de un perfil, podemos “votar”. Una *función de bienestar social* es una función que a cada *perfil* dentro de un *dominio* le asigna un ranking de preferencias colectivo (que vendría a representar el resultado del voto). Por ejemplo, una función *dictatorial* asigna a cada perfil $\langle R_1, \dots, R_n \rangle$ la función R_i para cierto agente i (podría ser el agente 3).

Teorema de Arrow

¿Qué propiedades debería tener una función de bienestar social para que la consideremos “razonable”? Podríamos pensar en propiedades básicas, como la siguiente: si *todos* individualmente prefieren a sobre b , entonces el grupo prefiere a sobre b (esta condición será conocida como Pareto Débil).

Naturalmente, hay distintos criterios. Para Arrow, una función de bienestar social f es satisfactoria cuando (además de arrojar un ranking completo y transitivo) cumple estas propiedades:

(Pareto Débil) Si todos los agentes prefieren a sobre b , entonces la preferencia colectiva pone a sobre b .

(Dominio Irrestricto) No hay ninguna restricción sobre los perfiles admitidos en el dominio de la función.

(Independencia de Alternativas Irrelevantes) Para determinar si colectivamente $a > b$, $b > a$ o $a \sim b$, lo único relevante es cómo votan los agentes sobre a y b .

(No-Dictadura) No hay ningún dictador; i.e., no hay ningún agente i tal que para todo perfil P , $f(P) = R_i$.

La condición de Pareto Débil es fácil de entender, porque se trata de una regla de unanimidad: si todos los individuos del grupo prefieren cerveza antes que vino, el grupo prefiere cerveza antes que vino. El Dominio Irrestricto nos dice que el sistema de voto no depende de las preferencias de los agentes: en principio, cualquier preferencia debería ser admitida, en tanto cumpla con las propiedades fundamentales de transitividad y completitud. La Independencia de Alternativas Irrelevantes propone que el cálculo del voto sea “enfocado”: si queremos establecer la preferencia social entre a y b , debería bastar con fijarnos en las preferencias individuales sobre estas dos opciones, y no sobre otras. La condición de No-Dictadura también es auto-explicativa: un buen sistema de elección no debería permitir que la elección social se base en hacer lo que un agente específico quiera hacer.

Y aquí es cuando Arrow prueba el teorema principal de la teoría de la elección social, luego conocido como el “Teorema de Imposibilidad de Arrow”:

(Teorema de Imposibilidad de Arrow)

Si hay tres o más opciones, ninguna función de bienestar social puede satisfacer Pareto Débil, Dominio Irrestricto,

Independencia de Alternativas Irrelevantes y No-Dictadura.

El teorema muestra que ningún método de voto satisface las cuatro propiedades deseables al mismo tiempo. Una prueba sencilla del teorema aparece en la próxima sección.

Se considera que el Teorema de Arrow es el principal resultado en la teoría de la elección social, y muchos (como veremos más adelante) lo han leído como un golpe fatal contra la idea misma de *democracia*.

Ejercicio

Muestre que una dictadura, es decir, un método donde el grupo decide el ranking del jugador i (el dictador), satisface Pareto Débil e Independencia de Alternativas Irrelevantes, además de arrojar un ranking completo y transitivo.

*Parte D: Prueba del Teorema de Arrow

Hay muchas pruebas del Teorema de Arrow. No nos interesa necesariamente la formulación original. Distintas pruebas fueron desarrolladas, y quizás las más conocidas son la original de Arrow (1951) y la versión de Sen (1970*b*).

La prueba que haremos es más sencilla, y se basa en el texto “A straightforward proof of Arrow’s theorem”, de Mark Fey (2014). La estructura de la prueba, tal como la presentamos aquí, es la siguiente:

- Primero, determinamos que respecto a dos opciones a/b , hay un agente que “voltea” las preferencias sobre $a > b$. (Paso 1)
- Luego, vamos mostrando que ese agente es *decisivo* sobre otras opciones. Un agente es *decisivo* sobre $x > y$ cuando, si el agente prefiere $x > y$, entonces el grupo entero prefiere $x > y$. (Pasos 2-5)
- Luego de varios pasos similares, llegamos a la conclusión de que ese agente es decisivo sobre el par (a, b) . (Paso 6)
- Finalmente, apelamos a los lemas previos para probar que, si el agente es decisivo sobre un par, es decisivo sobre todos. Por lo tanto, es un dictador. (Paso final)

Enunciado del Teorema:

Partimos de un método de elección f que satisface:

- Pareto Débil: si todos los agentes j prefieren $x \succ_j y$, entonces el grupo prefiere $x \succ y$.
- Dominio Irrestricto: El dominio de la función f es irrestricto, es decir, cualquier perfil es admisible.
- Independencia de Alternativas Irrelevantes (IAI): si tenemos un perfil P donde $f(P)$ pone a x sobre y , y tenemos otro perfil P' donde los agentes tienen las mismas preferencias respecto al par (x, y) , pero distintas preferencias respecto a otros eventos, entonces $f(P')$ también pone a x sobre y .

Con esas premisas, establecemos que hay un agente i que es un *dictador*: para todas las opciones x e y , si $x \succ_i y$, entonces $x \succ y$ (es decir, si i prefiere x sobre y , la sociedad prefiere x sobre y).

Prueba:

Paso 1: *hay un agente i que voltea las preferencias sobre a/b .*

Supongamos el siguiente perfil:

R_1	...	R_n
a		a
b		b
...		...

Aquí, por Pareto Débil, el grupo decide $a \succ b$.

Ahora supongamos otro perfil:

R_1	...	R_n
b		b
a		a
...		...

Aquí, por Pareto Débil, el grupo decide $b \succ a$.

Ahora vamos generando perfiles “intermedios”. A partir del primer perfil, donde todos votan $a \succ b$, damos vuelta a/b en cada

agente, empezando por el Agente 1. En algún momento el grupo pasará a votar $b > a$ (esto podría pasar al final, el tema es que en algún momento va a pasar).

Detectamos el agente que “da vuelta” a la preferencia. Es decir, el agente i donde:

Perfil P1

R_{i-}	R_i	R_{i+}
b	a	a
a	b	b
...

nos da que $a > b$, pero...

Perfil P2

R_{i-}	R_i	R_{i+}
b	b	a
a	a	b
...

nos da que $b > a$.

(“ R_{i-} ” son todos los agentes hasta i , y “ R_{i+} ” son todos los agentes después de i).

Paso 2: Para todo c , el individuo i es decisivo sobre $b > c$.

Usamos el perfil P1, pero ubicando una c :

R_{i-}	R_i	R_{i+}
b	a	a
c	b	b
a	c	c
...

Por Independencia de Alternativas Irrelevantes (de ahora en más, “IAI”) con el perfil P1, colectivamente sucede que $a > b$.

También tenemos $b > c$ por Pareto Débil.

Por Transitividad obtenemos $a > c$.

Ahora bien, ¿por qué el agente i decide sobre $b > c$?

Supongamos que ahora tenemos este perfil:

R_{i-}	R_i	R_{i+}
b/c	b	a
a	a	b/c
...	c	...
...

La terminología “ b/c ” significa que los otros agentes pueden preferir $b > c$, $c > b$ o $b \sim c$.

Aquí, el grupo va a seguir prefiriendo $b > c$. Esto se debe a que, por IAI respecto al perfil P2, el grupo establece $b > a$, y por IAI respecto al perfil anterior, también $a > c$ (dado que la ubicación relativa de a y c no ha cambiado). Entonces por Transitividad, obtenemos $b > c$.

Esto es independiente de lo que prefieran otros agentes sobre (b , c), por eso el agente i es decisivo sobre $b > c$.

El mismo estilo de prueba lo repetiremos para otros casos.

Paso 3. *Para todo c , el individuo i es decisivo sobre $a > c$.*

Supongamos que tenemos este perfil:

R_{i-}	R_i	R_{i+}
a/c	a	a/c
b	b	b
...	c	...
...

Por Pareto Débil, sabemos que $a > b$.

Dado que i es decisivo sobre $b > c$ (Paso 2), sabemos que $b > c$.

Entonces por Transitividad $a > c$.

Esto es independiente de lo que prefieran otros agentes sobre (a , c), por eso el agente i es decisivo sobre $a > c$.

Paso 4: *Para todo c , el individuo i es decisivo sobre $c > a$.*

Supongamos que tenemos este perfil:

R_{i-}	R_i	R_{i+}
b	c	c
c	a	a
a	b	b
...

Por IAI respecto al perfil P1, sucede que $a > b$.

Por Pareto Débil tenemos $c > a$.

Por Transitividad tenemos $c > b$.

Ahora consideremos perfiles de esta forma:

R_{i-}	R_i	R_{i+}
b	c	a/c
a/c	b	b
...	a	...
...

Por IAI respecto al perfil P2, tenemos $b > a$.

Por IAI respecto al perfil anterior, tenemos $c > b$.

Por Transitividad tenemos $c > a$.

Esto es independiente de lo que prefieran otros agentes sobre (c, a) , por eso el agente i es decisivo sobre $c > a$.

Paso 5: Para todo c , el individuo i es decisivo sobre $c > b$.

Ahora consideremos perfiles de esta forma:

R_{i-}	R_i	R_{i+}
a	c	a
b/c	a	b/c
...	b	...
...

Por Pareto Débil, sabemos que $a > b$.

Por Paso 4 (el agente i es decisivo sobre $c > a$) sabemos que $c > a$.

Por Transitividad tenemos que $c > b$.

Esto es independiente de lo que prefieran otros agentes sobre (c, b) , por eso el agente i es decisivo sobre $c > b$.

Paso 6: *El individuo i es decisivo sobre $b > a$ y sobre $a > b$.*
Supongamos este perfil:

R_{i-}	R_i	R_{i+}
a/b	a	a/b
...	c	...
...	b	...
...

Por Pareto Débil, sabemos que $a > c$.

Por el Paso 5 (i es decisivo sobre $c > b$), sabemos que $c > b$.

Por Transitividad obtenemos $a > b$.

Lo mismo ocurre para $b > a$, usando este perfil:

R_{i-}	R_i	R_{i+}
a/b	b	a/b
...	c	...
...	a	...
...

Aquí, por Pareto Débil sabemos que $b > c$.

Por Paso 4 (i es decisivo sobre $c > a$) sabemos que $c > a$.

Entonces por Transitividad obtenemos $b > a$.

Paso final: *El agente i es un dictador.*

Por el Paso 6, sabemos que i es dictador sobre (a, b) , i.e. es decisivo sobre $a > b$ y sobre $b > a$.

Probaremos que es dictador sobre cualquier otra opción (x, y) .

Consideremos este perfil:

R_{i-}	R_i	R_{i+}
x/y	x	x/y
a/b	a	a/b
...	b	...
...	y	...
...

Por Paso 4 (i es decisivo sobre $c > a$), tenemos $x > a$.

Por Paso 6 (i es decisivo sobre a/b), tenemos $a > b$.

Por Paso 2 (i es decisivo sobre $b > c$), tenemos $b > y$.

Por Transitividad, tenemos $x > y$.

(Para probar que es decisivo sobre $y > x$, solo cambiamos de lugar x con y).

QED.

Ejercicios

1. Asumiendo Dominio Irrestricto, muestre que si hay tres opciones (a, b, c), no puede suceder que un agente i sea decisivo sobre $a > b$, otro agente j sea decisivo sobre $b > c$, y otro agente k sea decisivo sobre $c > a$.

2. Asumiendo Dominio Irrestricto y Pareto Débil, muestre que si hay tres opciones (a, b, c), no puede suceder que un agente i sea decisivo sobre $a > b$ y otro agente j sea decisivo sobre $b > c$.

Parte E: Lecturas del Teorema de Arrow

¿Cómo podríamos interpretar el Teorema de Arrow? Como suele suceder en otros debates de Filosofía Política, la lectura que hagamos va a depender de nuestra inclinación ideológica.

En principio, el teorema establece la incompatibilidad entre cuatro propiedades: No-Dictadura, Independencia de Alternativas Irrelevantes, Dominio Irrestricto y Pareto Débil. Muchas lecturas se concentran en la incompatibilidad entre No-Dictadura e Independencia de Alternativas Irrelevantes. Como veremos adelante con más detalle, la Independencia de Alternativas Irrelevantes garantiza la no-manipulabilidad del voto. Entonces, una forma relativamente neutral de leer el teorema es decir que, si queremos un sistema democrático, tenemos que aceptar cierto nivel de manipulabilidad.

Uno de los libros más importantes en la tradición de la Teoría de la Elección Social es *Liberalism against populism* (1987) de William Riker. Según Riker, el Teorema de Arrow muestra los límites del “populismo”, entendido como la tradición política que se nutre de Rousseau, donde el gobierno surge de una voluntad

popular expresada mediante el voto. En el enfoque de Riker, el teorema de Arrow muestra que este tipo de gobierno es imposible, porque no hay una “voluntad popular” (si la hubiera, sería inconsistente o absurda). Y el sistema de voto popular solo favorece los mecanismos de manipulación.

En cambio, según Riker, el teorema funciona como justificación de una versión muy mínima del liberalismo republicano. Es decir, de la corriente que él atribuye a Madison, según la cual la función del voto es muy limitada, y se reduce a *controlar* y *vetar* lo que hacen los gobernantes: “Todo lo que las elecciones hacen o deben hacer es permitir que la gente se deshaga de sus gobernantes” (Riker 1987, p. 244).

Podríamos ver la teoría de Riker como inspirada parcialmente en Joseph Schumpeter, un economista austríaco que previamente a Arrow defendió un tipo de democracia muy limitada. Para Schumpeter, en su clásico libro *Capitalismo, Socialismo y Democracia* (1949), la democracia “populista” era peligrosa, y lo mejor que podemos esperar de la democracia es un voto para elegir *élites*. Es decir, la democracia es una forma de elegir entre proyectos políticos sobre los que, más adelante, no tendremos demasiada influencia. Si se quiere, podríamos leer el modelo de Schumpeter como una visión aristocrática (o pesimista) de la democracia. Un politólogo más actual que defiende las ideas de Schumpeter es Przeworski (2010). Según este autor, no podemos esperar *mucho* de la democracia, más allá de un sistema recurrente de votación.

Esta corriente “pesimista” se contrapone al enfoque de *democracia deliberativa*, que postula un sistema democrático donde la participación ciudadana va más allá de la representación indirecta en el congreso o el gobierno. Más adelante veremos cómo otro resultado de Condorcet se suele usar a favor de este tipo de sistema político.

Recién entrado el siglo XXI, muchos autores empezaron a cuestionar las lecturas “antidemocráticas” del Teorema de Arrow, especialmente la lectura dominante de Riker. Un autor que se tomó el trabajo de responder punto por punto a Riker fue Mackie (2003). Según Mackie, la visión democrática no defiende que el voto expresa la voluntad popular directamente, sino que es una

forma *aproximada* de expresarla. Rousseau (1762, p. 60) ya señalaba que la voluntad general no puede identificarse con la “voluntad de todos”, entendida como la suma de voluntades individuales; e incluso Arrow (1951) parece sugerir algo similar en el último capítulo de su libro (aunque esto abre el problema de cómo hallar la voluntad general, si no es por medio del voto). Mackie señala también que autores como Riker se enfocan demasiado en el efecto de la manipulación y los posibles votos incoherentes, pero argumenta (a partir del análisis de casos históricos) que tales fenómenos no se dan tan frecuentemente. Otros autores como Christian List y John Dryzek (2003) proponen una “reconciliación” entre la Teoría de la Elección Social y la democracia deliberativa, las dos corrientes que consideran “dominantes” (aunque usualmente contrapuestas) en la teoría política contemporánea. La idea general de List y Dryzek es que el aspecto deliberativo nos permite tener preferencias mejor “ordenadas”, y de ese modo evitar ciclos. Por ejemplo, una buena deliberación podría llevarnos a generar perfiles de “pico único”.

Parte F: Soluciones al Teorema de Arrow

Como podríamos esperar, hay distintas estrategias para bloquear al Teorema de Arrow. Las estrategias suelen atacar conceptualmente alguna de las premisas del teorema.

Voto sin Independencia: el método de Borda

La respuesta más común al Teorema de Arrow es el rechazo de la Independencia de Alternativas Irrelevantes. Esto puede hacerse buscando un método que, a diferencia del de Condorcet, no satisfaga esta propiedad.

Hay muchos métodos alternativos al de Condorcet. Por ejemplo, uno podría aplicar el *método de Borda*, propuesto por el matemático Jean-Charles de Borda en 1770:

(Método de Borda para órdenes estrictos)

Cada persona tiene un ranking de n objetos. Los primeros puestos valen n . los segundos puestos valen $n-1 \dots$ así

sucesivamente, hasta que los últimos puestos valen 1. El ranking final se determina sumando los puntos de cada opción.

Por ejemplo, supongamos que el perfil es P1:

1	2	3
<i>a</i>	<i>a</i>	<i>c</i>
<i>b</i>	<i>b</i>	<i>d</i>
<i>c</i>	<i>d</i>	<i>b</i>
<i>d</i>	<i>c</i>	<i>a</i>

El puntaje de *a* será $4 + 4 + 1 = 9$. El puntaje de *b* será $3 + 3 + 2 = 8$. El puntaje de *c* será $2 + 1 + 4 = 7$. Y el puntaje de *d* será $1 + 2 + 3 = 6$. Entonces el ranking final será: $a > b > c > d$.

El método de Borda no puede generar ciclos por razones puramente matemáticas (lo único que hace este método es darle un puntaje a cada opción). Cuando el método de Condorcet genera ciclos, el método de Borda genera empates.

Una característica importante del método de Borda (y otros similares) es que, aunque sí satisface Pareto Débil (la prueba queda al lector), no satisface Independencia de Alternativas Irrelevantes. Por ejemplo, imaginemos este perfil P2:

1	2	3
<i>a</i>	<i>a</i>	<i>b</i>
<i>b</i>	<i>b</i>	<i>d</i>
<i>c</i>	<i>d</i>	<i>c</i>
<i>d</i>	<i>c</i>	<i>a</i>

Este perfil es casi idéntico al anterior, pero el agente 3 intercambió *b* por *c*. Ahora el ranking final cambia: El puntaje de *a* es $4 + 4 + 1 = 9$; el puntaje de *b* es $3 + 3 + 4 = 10$, el de *c* es $2 + 1 + 1 = 5$, y el de *d* es $1 + 2 + 3 = 6$. Es decir, ahora queda $b > a > d > c$. Aquí vemos que no se cumple la Independencia de Alternativas Irrelevantes, porque los votos sobre $\{a, b\}$ son idénticos a los

del perfil P1 (1 prefiere *a* sobre *b*, 2 prefiere *a* sobre *b*, y 3 prefiere *b* sobre *a*), pero el resultado final respecto a estas dos opciones es distinto.

La violación de Independencia de Alternativas Irrelevantes permite un fenómeno llamado *manipulabilidad*. La forma más sencilla de entender la manipulabilidad es a partir de una *regla de elección social*. Una regla de elección social es una función que, a partir de un perfil, arroja una opción ganadora (o un conjunto de opciones ganadoras). Por ejemplo, la regla de Borda puede verse como una función de elección social donde gana la opción que obtuvo más puntos: en el perfil P1 gana la opción *a* y en P2 gana la opción *b*. Decimos que un método de elección social es *manipulable* cuando un agente podría mentir sobre sus preferencias reales, y obtener un ganador más preferible que el que hubiese obtenido votando honestamente. Para ilustrarlo, podríamos imaginar que las preferencias reales son las de P1, pero en una votación el agente 3 mintió sobre sus preferencias (y presentó su orden de P2) para lograr un resultado más preferible para él (que gane *b* en vez de *a*).

La manipulabilidad está en todos nuestros sistemas de voto: por ejemplo, cuando votamos un candidato que no es nuestro preferido (pero que tiene chances de ganar), para que no gane otro candidato que nos gusta aún menos.

Problemas similares se aplican a nociones *cardinales* de decisión colectiva. Aunque las escalas cardinales suelen aportar información suficiente para generar un ranking colectivo sin ciclos, se suele defender que estas escalas dan *demasiada* información. Así, abren la puerta para estrategias variadas de manipulación.

Imaginemos unos padres que deben darle de comer a dos hijos. Uno de los hijos quiere comer pastas, y el otro quiere comer carne. Si solo cuentan los rankings, es imposible tomar una decisión racionalmente. Sin embargo, uno de los hermanos se pone a llorar y gritar, diciendo “¡odio la carne, amo las pastas, por favor hagan carne o voy a sufrir!”. Los padres, al darse cuenta lo importante que es esto para el niño, deciden hacer carne. Luego de comer, el niño revela que sus berrinches fueron puro teatro. Así los padres aprenden el significado de la manipulabilidad. La

próxima vez, les dirán a sus hijos: “digan qué prefieren, pero por favor, sin llorar”.

Un problema relacionado es la dificultad para *comparar* utilidades de distintas personas. Esto afecta principalmente a los puntos de vista más utilitaristas. Si el objetivo de una decisión colectiva es maximizar el bienestar del grupo, un joven quisquilloso que *sufre muchísimo* por cualquier decisión que no satisfaga sus deseos individualistas tendría prioridad sobre otras personas que, con preferencias igualmente respetables, sufren menos. Este problema porque aparece incluso si asumimos la sinceridad de los individuos. Esta es otra razón por la que en la teoría de la elección social se suele trabajar con escalas ordinales.

Restricciones de dominio

Podemos encontrar soluciones muy variadas al Teorema de Arrow. Por ejemplo, hay contextos donde parece más justificada una restricción de dominio. En la parte B mencionamos un tipo de restricción de dominio: la restricción a perfiles de pico único. Aquí nos centraremos en otra opción.

En un artículo ya clásico, Okasha (2011) sostuvo que el Teorema de Arrow también aplica al amalgamiento de evidencia científica: no hay modo racional de decidir entre hipótesis cuando distintos experimentos resultan en preferencias distintas entre hipótesis, porque deberían cumplir los criterios de Arrow y entonces ninguna opción será satisfactoria. En Cresto & Tajer (2020) sostenemos que, para el caso de la ciencia empírica, muchas veces podemos apelar a restricciones de dominio. Para eso usamos la tesis de Duhem-Quine: según esta conocida tesis sobre el funcionamiento de la ciencia, las hipótesis centrales son evaluadas en conjunto con hipótesis auxiliares. Como consecuencia, no cualquier orden es individualmente racional. Con esas restricciones, probamos que hay funciones que satisfacen Pareto Débil e Independencia de Alternativas Irrelevantes.

Para dar una idea de cómo el rechazo de Dominio Irrestricto puede salvarnos del problema, pensemos en el siguiente escenario. Hay que elegir qué comer, entre cuatro opciones:

- Nuggets en McDonald's

- Hamburguesa en McDonald's
- Ensalada César en Green Salad
- Sandwich de atún en Green Salad

Ahora supongamos que no todos los perfiles son admisibles. Pensamos, por ejemplo, que si te gusta la comida basura vas a preferir McDonald's, y si te gusta la comida sana vas a preferir Green Salad. Entonces solo admitimos perfiles donde esas preferencias vienen en *bloques* ordenados estrictamente (sin empaques entre bloques), como estos:

1	2	3
Nuggets, hamburguesa	ensalada	sandwich de atún
Sándwich de atún	sandwich de atún	ensalada
Ensalada	hamburguesa	nuggets
	nuggets	hamburguesa

El Agente 1 prefiere la comida chatarra a la comida sana, y los Agentes 2 y 3 prefieren lo contrario. Dentro de cada bloque, uno puede preferir lo que quiera.

Ahora pensemos en la siguiente regla de elección:

(Semi-dictadura)

En un grupo de $n > 2$ personas, decimos que el Agente 1 decide sobre los bloques, y el Agente 2 decide sobre los órdenes internos.

Si fueran los agentes del grupo anterior, el resultado sería:

hamburguesa > nuggets > ensalada > sandwich de atún

Podemos ver que obtenemos un orden coherente que es distinto al de cualquiera de los tres agentes. Es decir, no hay Dictadura. También hay Pareto Débil e Independencia de Alternativas Irrelevantes (la prueba queda al lector, pero es fácil).

Estas reglas *no* podrían funcionar en dominios irrestrictos, porque generarían inconsistencias; es necesario que los perfiles es-

tén ordenados “en bloques”. Nótese que el resultado no es terriblemente prometedor, porque si bien no tenemos dictadores, ahora tenemos semi-dictadores.

Ejercicios

1. Determine el ranking de Borda a partir de este perfil.

1	2	3
<i>a</i>	<i>b</i>	<i>d</i>
<i>b</i>	<i>c</i>	<i>c</i>
<i>c</i>	<i>a</i>	<i>b</i>
<i>d</i>	<i>d</i>	<i>a</i>

2.* Diseñe alguna estrategia de manipulación para el jugador 3.

2. Muestre que el método de semi-dictadura en dominios restringidos “por bloque” satisface Pareto Débil y es completo.

3.* Muestre que el método de semi-dictadura en dominios restringidos “por bloque” satisface Transitividad.

Parte G: Paradoja Del Liberal Paretiano

Otra paradoja muy conocida en Teoría de la Elección Social, mucho más sencilla que la de Arrow, fue mostrada por Amartya Sen (1970*a*). Más adelante se llamó a esta paradoja “Paradoja del Liberal Paretiano”.

La idea de Sen es, a grandes rasgos, mostrar cierta tensión entre derechos individuales y decisiones colectivas. En el ideario liberal (sea de derecha o de izquierda), hay un conjunto de decisiones que pertenecen a la “esfera privada” y sobre eso no pueden decidir los otros (por ejemplo, el color de mi ropa, o el próximo libro que voy a leer); y hay otro conjunto de decisiones que pertenecen a la “esfera pública”, y sobre eso podemos decidir colectivamente (por ejemplo, cuánto voy a pagar de impuestos). La paradoja cuestiona esa distinción.

Los derechos individuales son entendidos de este modo:

(Liberalismo) Para cada agente i , hay un par de alternativas a y b sobre las cuales el agente i es dictador.

La idea es que un agente es dictador sobre algunas cosas dentro de su esfera privada, como su color de pelo, o la película que va a ver esta noche.

La paradoja del liberal paretiano nos dice lo siguiente:

(Paradoja del liberal paretiano)

Ninguna función de bienestar social puede satisfacer Liberalismo, Dominio Irrestricto y Pareto Débil.

Prueba: Supongamos que el Agente 1 es dictador sobre (a, b) y el Agente 2 es dictador sobre (c, d) . Obviamente (a, b) y (c, d) no pueden ser el mismo par, porque si lo fueran, habría contradicción en caso de que $a >_1 b$ y $b >_2 a$. Ahora bien, (a, b) y (c, d) no pueden tener un elemento en común. Porque supongamos que lo tuvieran, por ejemplo, $a = d$. Entonces los pares son (a, b) y (c, a) . Ahora supongamos que $a >_1 b$ y $c >_2 a$. Y que para toda la sociedad $b > c$. Entonces (por Pareto Débil) se genera un ciclo y no puede haber función de bienestar social.

Ahora, supongamos que no tienen ningún elemento en común. Y ahora decimos que $a >_1 b$ y $c >_2 d$. Pero todos los agentes pueden votar unánimemente $b > c$ y $d > a$, generando otro ciclo (por Pareto Débil). QED.

El teorema tuvo distintas lecturas y generó una extensa discusión. Tal como sucedió en el caso de Arrow, la interpretación del mismo Sen fue algo ambigua: su intención era mostrar una incompatibilidad entre el paretianismo y los derechos individuales, sin proponer ninguna salida en especial. En textos siguientes (Sen 1975), parece haber defendido un rechazo o restricción del principio de Pareto Débil.

Quizás la respuesta más influyente al desafío fue la de Nozick (1974, p. 166). Este autor, representante del pensamiento liber-

tario, propone que el voto viene *después* de los derechos individuales. Es decir, aquellos eventos (a, b) sobre los cuales decide el Agente 1 no pueden formar parte de un conjunto de eventos sobre los que decide el resto de la sociedad. Entonces si $a \succ_1 b$ y $c \succ_2 d$, la sociedad no puede decidir libremente sobre eventos como b y c . Aunque Nozick no presentó los detalles formales de su idea, una forma natural de entenderla es como un rechazo (o restricción) de Pareto Débil (Sen 1996, p. 156). Una propuesta similar fue desarrollada por Saari (1997, p. 92): su nuevo principio, “Pareto Relajado”, propone que Pareto Débil se aplique entre opciones (a, b) que no pertenecen a ninguna esfera de derechos individual.

Otros autores proponen restringir el principio aquí llamado “Liberalismo”. Por ejemplo, Gibbard (1974) intenta formalizar la idea intuitiva de que las personas pueden decidir sobre su esfera privada cuando no entra en conflicto con la esfera privada de los demás. Esto requiere especificar en qué casos un derecho es “abandonado” por conflictuar con los derechos de otros (los detalles formales exceden la complejidad de este libro). Muchos autores han seguido esta tendencia, al intentan formalizar el concepto de “derecho” para el marco de la elección social.

Ejercicios

1.

Gibbard (1974) presentó una paradoja similar a la de Sen, pero sin usar Pareto Débil.

Alberto y Carolina (Agente 1 y 2) tienen una cita, y cada uno pide una copa, de vino tinto (T) o blanco (B). El resultado será un par; por ejemplo, (T, B) significa que Alberto toma vino tinto y Carolina toma vino blanco. Cada uno es decisivo sobre su esfera de influencia: Alberto (Agente 1) puede decidir (por ejemplo) entre (T, B) y (B, B) . El problema es que Alberto quiere tomar dos vinos distintos, y Beatriz quiere que ambos tomen el mismo. Supongamos que estas fueran sus preferencias:

Alberto: $(B, T) \succ (T, B) \succ (T, T) \succ (B, B)$

Carolina: $(T, T) \succ (B, B) \succ (B, T) \succ (T, B)$

¿Qué problema se genera? ¿Cómo lo solucionarías?

2.

Supongamos que los pares sobre los que cada individuo es decisivo no tienen elementos en común. Ahora utilizamos esta regla de elección, llamada Dictadura Limitada:

- Si (x, y) no pertenecen a ninguna esfera privada, el Agente 1 (dictador) decide sobre (x, y) .
- Si x pertenece a alguna esfera privada y z pertenece a la esfera “pública”, entonces $x > z$.
- En el orden colectivo, las “esferas privadas” se ordenan según el número de cada votante (primero va la esfera del Agente 1, luego del Agente 2, etc.).
- Cada agente decide sobre su esfera privada (a, b) .

Muestre que la regla genera un orden (completo y transitivo), y satisface Liberalismo y Pareto Relajado.

Parte H: Resultados positivos sobre el voto

La paradoja de Condorcet y el Teorema de Arrow son bastante pesimistas respecto al voto democrático. Sin embargo, el mismo Condorcet, en su *Ensayo* antes citado, probó su *Teorema del Jurado* que muestra que, en algunas circunstancias, el voto puede funcionar. Este resultado funciona de base para las justificaciones *epistémicas* de la democracia (Goodin & Spiekermann, 2019). La idea de las justificaciones epistémicas es que el voto democrático nos acerca a la *verdad*.

Empecemos por imaginar un grupo de expertos sobre un determinado fenómeno, que deben ponerse de acuerdo sobre p o $\neg p$. Este teorema usa probabilidades: ¿cuál es la probabilidad de que el grupo diga lo correcto, si cada uno es bastante confiable? La idea de Condorcet es que, si cada individuo es bastante confiable, y las personas votan de forma independiente, el grupo será *muy* confiable.

Es decir, se asumen dos premisas. En primer lugar, que los agentes son bastante competentes en el asunto:

(Competencia) Hay un $r > 0.5$, tal que para cada individuo i del grupo, la confiabilidad de i sobre p es r . Es decir, $P(i \text{ dice } p \mid p) = r$.

Nótese que r debe ser igual para todos los agentes. También se asume que los agentes tienen un criterio independiente entre sí:

(Independencia) Para dos agentes i y j , $P(i \text{ dice } p \mid j \text{ dice } p) = P(i \text{ dice } p)$. Es decir, los juicios de los agentes son independientes entre sí.

Entonces podemos probar lo siguiente:

(Teorema del jurado de Condorcet)

Asumiendo Competencia e Independencia, podemos probar dos cosas:

Parte finita: A medida que aumenta el número de agentes (siempre que sea impar), y si votan por mayoría, la confiabilidad del grupo irá aumentando (y siempre será mayor a r).

Parte infinita: La confiabilidad del grupo, a medida que agregamos agentes, tiende a 1.

Podríamos entender el teorema como una prueba de que la sabiduría colectiva siempre supera a la individual: no importa cuán confiable sea un agente, un grupo suficientemente grande de agentes bastante confiables votando por mayoría lo va a superar. La prueba general no es tan sencilla, pero mostraremos el primer paso para un caso específico. Supongamos que tenemos un solo agente, con confiabilidad 0.6. Ahora agregamos dos agentes con confiabilidad 0.6. ¿Cuál es la probabilidad de que, votando por mayoría, den con la opción correcta?

Supongamos que sucede p . Ahora hay dos formas de ganar:

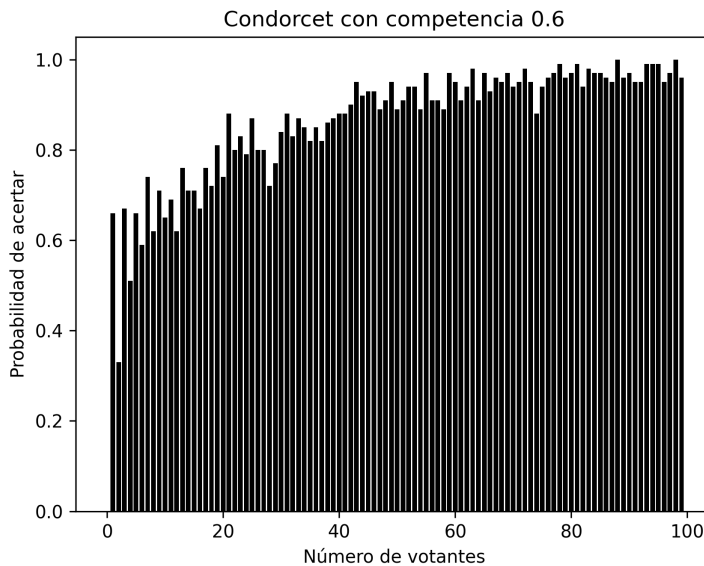
- A. Los tres agentes votan p .
- B. Dos agentes votan p , y un tercer agente vota $\neg p$. Esto puede pasar de tres formas distintas.

La probabilidad de A es $0.6^3 = 0.216$.

La probabilidad de B es $3 \times [0.6 \times 0.6 \times 0.4] = 3 \times 0.144 = 0.432$.

Entonces la probabilidad de que pase A o B es $0.216 + 0.432 = 0.648$. Es decir, es mayor que 0.6.

Para mostrar este resultado es muy común usar simulaciones. Por ejemplo, en esta simulación podemos ver que si la confiabilidad de los agentes es 0.6, mientras más agentes haya, más confiable será el grupo:



Como antes señalé, el Teorema del Jurado de Condorcet se usó muchas veces, especialmente en los últimos años, para justificar la democracia desde un punto de vista epistémico (Estlund 2008, Goodin & Spiekermann 2019).

Sin embargo, el teorema tiene algunos puntos cuestionables:

- a. Supone que *todos* los agentes tienen la misma confiabilidad.
- b. Supone la *independencia* entre opiniones, algo difícil de que funcione en realidad.

La propiedad (a) es la menos problemática, porque la homogeneidad en confiabilidad no es *necesaria* para el teorema: podríamos probar lo mismo (para la versión infinita) si asumimos que los agentes tienen en *promedio* confiabilidad 0.6, ya sea con una distribución uniforme o una distribución normal.

La parte (b) sí es más compleja, porque en el mundo real las opiniones están influidas por “agentes centrales” como el Estado y los medios de comunicación masivos. Grofman y Feld (1988) proponen que el teorema del Jurado de Condorcet representa el ideal de Rousseau de la voluntad general. Recordemos que la voluntad general de Rousseau no admite pactos (esto vendría a representar la Independencia) y está formada de gente moderadamente educada (esto vendría a representar la Competencia).

Hay dos resultados que también suelen utilizarse para cuestionar el teorema. En primer lugar, si bien el teorema muestra que un grupo suficientemente grande de personas competentes e independientes es más confiable que una persona muy competente, no es necesario tener *millones* de personas: con cientos o miles es suficiente. Entonces, algunos autores usan el resultado para cuestionar la democracia (donde votan todos) y defender la *epistocracia*, donde votan las personas que pueden pasar una prueba de conocimiento general (Brennan 2018). Autores como Estlund (2008) sugieren que la democracia no puede justificarse *solo* de forma epistémica, sino que también necesita apelar a la universalidad de los derechos políticos.

Una segunda deriva preocupante del teorema es que, si los agentes son no-confiables (es decir si $r < 0.5$), entonces la democracia los llevará a decisiones equivocadas casi siempre. Es decir, este teorema es simétrico, para lo bueno y para lo malo.

Ejercicio

Mostrar que tres agentes con 0.7 de confiabilidad votan mejor (por mayoría) que un solo agente con 0.7 de confiabilidad.

Parte I: Agregación de juicios

En esta última sección, discutiremos la teoría de agregación de juicios. El problema de la agregación de juicios es similar al problema de la agregación de preferencias, y por eso se lo suele ver como un pariente cercano del problema de la elección social. Similar a la paradoja de Condorcet para la teoría de la elección social, en la teoría de agregación de juicios existe la *Paradoja doctrinal* (Kornhauser & Sager, 1987). En este escenario, tres jurados deben votar sobre si determinado agente es culpable (r) o no ($\neg r$). Todos están de acuerdo en esto: “El acusado es culpable si y sólo si el acusado cometió el acto, y el acto está castigado por ley”. Es decir: $r \leftrightarrow (p \ \& \ q)$. Podemos entonces reemplazar r por $(p \ \& \ q)$. Los agentes votan del siguiente modo:

Juez 1:	p	$\neg q$	$\neg(p \ \& \ q)$
Juez 2:	$\neg p$	q	$\neg(p \ \& \ q)$
Juez 3:	p	q	$p \ \& \ q$

El problema aparece cuando queremos obtener una decisión colectiva a partir de estos votos. Si todas las proposiciones se deciden por voto mayoritario, el resultado es $\{p, q, \neg(p \ \& \ q)\}$. Es decir, el voto mayoritario arroja un resultado inconsistente, aun cuando los votos individuales eran consistentes. Pettit y Rabinowicz (2001) llamaron a esta paradoja “dilema discursivo”.

La teoría de agregación de juicios estudia este tipo de problemas de agregación. Podríamos distinguir entre dos tipos de investigaciones, similar a lo que sucede con la teoría de elección social. Por un lado, podemos desarrollar resultados nuevos de imposibilidad. Otro enfoque, más constructivo, consiste en proponer métodos de agregación que satisfagan determinadas propiedades.

Resultados de imposibilidad

Un resultado de imposibilidad nos sirve para ver de forma más general cuáles son los problemas que causan determinada para-

doja. Así como el Teorema de Arrow puede verse como una generalización de la Paradoja de Condorcet, varios autores intentaron generalizar el Dilema discursivo probando resultados de imposibilidad. El primer resultado fue de List y Pettit (2002); aquí vamos a enunciar (sin prueba) un resultado equivalente de Dietrich y List (2007), que nos ayuda a ver las similitudes entre el Dilema discursivo y el Teorema de Arrow. Para entender este resultado y sus posibles soluciones hace falta un poco más de sofisticación formal.

En el marco de la teoría de la agregación de juicios, una *agenda* es un conjunto X de proposiciones cerrado bajo negación (para simplificar, ignoramos las dobles negaciones). Estas son las proposiciones sobre las cuales votan los agentes. Por ejemplo, una agenda podría ser $\{p, q, \neg p, \neg q, p \& q, \neg(p \& q)\}$.

Por otro lado, un *conjunto de juicios* es un conjunto J de proposiciones de la agenda que acepta determinado agente. Es decir, $J \subseteq X$. Se asume que los conjuntos de juicios J son *consistentes* (tienen modelo) y *completos* (para cada proposición A en la agenda, o eligen A o $\neg A$). Por último, un *perfil* es una n -tupla de conjuntos de juicios, $P = (J_1, \dots, J_n)$.

Ahora podemos presentar el resultado de imposibilidad:

Resultado de imposibilidad (Dietrich & List 2007, Teorema 2): Si la agenda tiene dos letras proposicionales y alguna composición ($p \& q, p \vee q$, etc.), entonces cualquier regla de agregación f que satisface estas condiciones es una **Dictadura**:

- *Dominio universal*: El dominio de f es el conjunto de todos los perfiles de juicios completos y consistentes sobre la agenda X .
- *Racionalidad colectiva*: Siempre el resultado $f(P)$ es un conjunto de juicio (es decir, es consistente y completo).
- *Unanimidad*: Si todos los individuos de un perfil P aceptan A , entonces $A \in f(P)$.

- *Independencia*: Si una proposición A es aceptada por los mismos agentes en distintos perfiles P y P' respectivamente, entonces $A \in f(P)$ si y solo si $A \in f(P')$.

De este modo, podemos ver el paralelo entre este teorema de imposibilidad para la agregación de juicios y el teorema de imposibilidad de Arrow.

Soluciones: Rechazar Independencia

Similar a lo que sucede con el Teorema de Arrow, la respuesta más natural al problema de la agregación de juicios es rechazar la Independencia. La idea es que el voto colectivo sobre una proposición no debería depender solamente de los votos individuales sobre esa misma proposición.

Por ejemplo, el voto sobre una conclusión podría depender de las *razones* a favor de ella. Ese es el espíritu de la Regla de Premisas, que explicaremos a continuación.

En una versión sencilla de la Regla de Premisas, los agentes votan por mayoría sobre las fórmulas atómicas, y los votos sobre las fórmulas compuestas se derivan de ahí. Entonces el voto se decidiría de este modo, manteniendo la consistencia:

Juez 1:	p	$\neg q$	$\neg(p \& q)$
Juez 2:	$\neg p$	q	$\neg(p \& q)$
Juez 3:	p	q	$p \& q$
Premisas:	p	q	$p \& q$

Es fácil ver por qué este método no satisface Independencia. Supongamos que el Juez 2, buscando manipular el resultado sobre $(p \& q)$, decidiera rechazar q . Ahora, la regla de mayoría arrojaría $\neg q$, y la Regla de Premisas nos daría $\neg(p \& q)$. Pero esto viola Independencia: en este nuevo perfil los votos individuales sobre $(p \& q)$ son iguales a los del perfil anterior, pero el voto colectivo sobre esta proposición será distinto.

Hay muchas otras reglas que violan la condición de Independencia, como las reglas secuenciales y reglas de distancia. Por razones de espacio, no indagaremos esas opciones aquí.

Soluciones: Restricciones de dominio

Otra posibilidad, similar a lo que ocurre en el Teorema de Arrow, es restringir el dominio a aquellos perfiles que tienen alguna propiedad deseable. Aquí no usaremos la propiedad de “pico único” sino algo equivalente. Decimos que un grupo está *uniformemente alineado* respecto a una agenda X , cuando podemos ordenar a los agentes de modo tal que, para cada proposición p de la agenda X , los agentes que aceptan la proposición están todos a la derecha (o todos a la izquierda) de los que la rechazan.

Por ejemplo, supongamos que fueran proposiciones sobre el *aborto*, y los votantes van de la ultra-izquierda (UI) a la ultra-derecha (UD):

	UI	I	C	D	UD
p = El feto es persona	0	0	1	1	1
q = Mujeres tienen derecho a abortar.	1	1	1	0	0
r = Estado debe financiar a la Iglesia	0	0	0	1	1

De este modo, p tiene 3 votos, q tiene 3, y r tiene 2. Es decir, el grupo por mayoría va a decidir $\{p, q, \neg r\}$.

Es curioso que el voto mayoritario coincide con lo que vota el agente C. Un resultado muy sencillo nos muestra que en tanto el grupo está uniformemente alineado, podemos aplicar el voto por mayoría, y *va a coincidir con el votante medio* (esto usualmente se llama “Teorema del Votante Medio”). Esto garantiza la consistencia, dado que todos los votantes son individualmente consistentes. Este resultado destaca la solidez de los sistemas políticos con lineamientos claros entre izquierda y derecha. Y, si creemos que la deliberación favorece estos alineamientos, también habla a favor del método deliberativo.

Dryzek y List (2003) usan este resultado para reconciliar la tradición deliberativa con los resultados en agregación de juicios. Ellos proponen que la deliberación ayuda en dos aspectos. Por un lado, luego de discutir, un grupo suele tener posiciones más

homogéneas. Por otro lado, luego de la deliberación, los integrantes del grupo se alinean más, permitiendo luego aplicar métodos de agregación razonables. Es decir, la deliberación ayuda a estructurar a los votantes y homogeneizar los grupos bajo ejes socialmente reconocidos, tales como “derecha” e “izquierda”.

Ejercicio

1.
En el ejemplo de la paradoja discursiva con $(p \ \& \ q)$, el tercer agente acepta esta proposición. Elabore un ejemplo de Paradoja Discursiva donde ningún agente acepte la “conclusión”.
2.
Supongamos que se usa una regla de Dictadura Condicional: si el voto por mayoría genera un resultado inconsistente, el Juez 3 se transforma en un dictador e impone su perfil. Esta regla satisface Unanimidad y evita inconsistencias. Mostrar que viola Independencia.

Soluciones para el capítulo 4

PARTE A

1.
El orden resultante es $b > a > c$.
2.
El problema del método de Lull es que el ganador dependerá del orden de las votaciones. Por ejemplo, en el caso del perfil paradójico de Condorcet, si empezamos votando (a, b) ganará c , y si empezamos votando (a, c) , ganará b .

PARTE C

El método de Dictadura satisface Pareto Débil e Independencia de Alternativas, además de ser completo y transitivo.
Supongamos que el dictador es i .
Pareto Débil: Si todos los agentes votan $a > b$, entonces el agente i también vota $a > b$, entonces el grupo decide $a > b$.

Independencia de Alternativas: Sea P un perfil tal que $a > b$ en $f(P)$. Dado que f es una dictadura, el agente i prefiere $a > b$ en P . Sea P' un perfil donde se mantienen las preferencias sobre $a > b$, pero cambian preferencias sobre otras alternativas. Esto significa que el dictador i aun prefiere $a > b$. Entonces $a > b$ en $f(P')$.

Complejidad y Transitividad: Dado que el ranking individual del dictador i debe ser completo y transitivo, el ranking colectivo también lo será.

PARTE D

1.

Supongamos que i es decisivo sobre $a > b$, j es decisivo sobre $b > c$, y k es decisivo sobre $c > a$. Sea P un perfil donde i prefiere $a > b$, j prefiere $b > c$ y k prefiere $c > a$. Por la definición de “decisivo”, la preferencia colectiva será paradójica: $a > b$, $b > c$, pero $c > a$.

2.

Supongamos que i es decisivo sobre $a > b$, y j es decisivo sobre $b > c$. Supongamos que aplica la propiedad de Pareto Débil. Sea P un perfil donde i prefiere $a > b$ y j prefiere $b > c$, y todos los agentes prefieren $c > a$. Por Pareto Débil y definición de “decisivo”, la preferencia colectiva será paradójica.

PARTE F

1.

El resultado de Borda es $b > c > a > d$, con 9, 8, 7 y 6 votos respectivamente.

2.

Como forma de manipulación, el agente 3 podría mover c hacia arriba y b hacia abajo. Así ganaría c con 9 puntos.

2.

Complejidad: si (x, y) pertenecen al mismo bloque, xR_y si y sólo si xR_iy , donde i es el dictador del orden interno de los bloques. En cambio, si pertenecen a distintos bloques, xR_y si y sólo si xR_jy , donde j es el dictador sobre bloques distintos.

Pareto Débil: Si $x > y$ para todos los agentes, esto significa que $x > y$ para ambos semi-dictadores, entonces en el orden colectivo $x > y$.

3.

Supongamos que $x \geq y$ y $y \geq z$.

Supongamos que (x, y, z) pertenecen al mismo bloque. Entonces se ordenan como disponga el semi-dictador del orden interno de los bloques, cuyo orden de preferencias es transitivo; por lo tanto, $x \geq z$.

Supongamos que (x, y, z) pertenecen a bloques distintos. Entonces se ordenan como disponga el semi-dictador de bloques distintos, cuyo orden de preferencias es transitivo; inferimos $x > z$.

Ahora supongamos que (x, y) pertenecen a un bloque, pero z pertenece a otro. La preferencia $y > z$ la establece el semi-dictador sobre bloques distintos. Entonces también $x > z$, porque las preferencias de ese semi-dictador también se establecen en forma de bloque. (Lo mismo cuando x es de un bloque y (y, z) de otro)

PARTE G

1.

Se genera un ciclo. Alberto genera que $(B, T) > (T, T)$ y $(T, B) > (B, B)$. Y Carolina genera que $(T, T) > (T, B)$ y $(B, B) > (B, T)$. Entonces $(B, T) > (T, T) > (T, B) > (B, B) > (B, T)$.

La explicación de Gibbard (1974) es que nuestras preferencias aquí son condicionales: si Alberto toma tinto, Carolina prefiere tinto sobre blanco, pero si Alberto toma blanco, Carolina prefiere blanco sobre tinto. No podemos ser dictadores sobre nuestra esfera privada si nuestras preferencias son condicionales a lo que hagan los demás.

2.

La satisfacción de Liberalismo es trivial: cada uno decide sobre su esfera privada. La regla satisface Pareto Relajado porque el dictador decide sobre cada par que no pertenece a ninguna esfera privada: entonces cuando todos los agentes prefieren $x > y$, incluyendo al dictador, se infiere que $x > y$.

Respecto a Completitud, hay tres casos a considerar. Primero, x está en una esfera privada, mientras que y es pública. Entonces inferimos $x > y$. Si ambas pertenecen a una esfera privada, decide el agente correspondiente. Si ambas pertenecen a la esfera pública, decide el dictador. Si x es de una esfera privada y y de otra, entonces se decide por la numeración de los agentes.

Transitividad es algo más compleja. Si $x \geq y$ y $y \geq z$, hay varios casos a considerar. Si son de tres esferas privadas distintas, entonces se infiere $x > z$, por la numeración de los agentes. Si (x, y) son de una esfera privada y z de otra, entonces $y > z$ también implica que $x > z$. Si x es de una esfera privada y (y, z) de otra, entonces el hecho de que $x > y$ implica que $x > z$. Si (x, y) pertenecen a una esfera privada y z a la esfera pública, entonces $x > z$ (lo mismo si x es privada pero (y, z) son públicas).

PARTE H

Mostraremos que tres agentes con confiabilidad 0.7 votan mejor que uno.

Si hay tres agentes con confiabilidad 0.7, pueden suceder tres cosas para que el voto mayoritario esté en lo correcto:

Los tres votan bien: $0.7^3 = 0.343$.

Dos votan bien y uno vota mal.

$$3 \times (0.7 \times 0.7 \times 0.3) = 0.147 \times 3 = 0.441$$

$$\text{Entonces } 0.343 + 0.441 = 0.784$$

PARTE I

1. Paradoja donde *ningún* agente acepta la “conclusión”:

	A	B	C	A&B&C
Agente 1	1	1	0	0
Agente 2	0	1	1	0
Agente 3	1	0	1	0
Mayoría	1	1	1	0

2. La regla de Dictadura Condicional viola Independencia. Supongamos que tenemos este perfil:

Juez 1:	$\neg p$	$\neg q$	$\neg(p \& q)$
Juez 2:	$\neg p$	q	$\neg(p \& q)$
Juez 3:	p	q	$p \& q$

Este perfil nos arroja $\{\neg p, q, \neg(p \& q)\}$ por el voto mayoritario. Pero si el Juez 1 cambiara su voto hacia p nos daría lo siguiente

Juez 1:	p	$\neg q$	$\neg(p \& q)$
Juez 2:	$\neg p$	q	$\neg(p \& q)$
Juez 3:	p	q	$p \& q$

Este es el perfil del Dilema Discursivo. Entonces, siguiendo Dictadura Condicional, debería decidir el Juez 3, y se impondrá $\{p, q, (p \& q)\}$. Esto obviamente viola Independencia porque ambos perfiles tienen los mismos votos individuales sobre $(p \& q)$, pero el resultado del voto colectivo es distinto.

EPÍLOGO

Hemos llegado al final del libro. A lo largo de los diferentes capítulos, recorrimos conceptos de probabilidad, teoría de la decisión, teoría de juegos y elección social. Espero haber mostrado a los lectores (especialmente a aquellos que se enfrentan a estos temas por primera vez) que se trata de un campo interesante y con múltiples aplicaciones para la filosofía.

Como resulta obvio, el acercamiento a los temas fue algo general, y quedaron muchos conceptos y discusiones importantes sin explorar. Algunos temas los excluí porque no me interesan tanto, otros por su complejidad matemática y otros porque no son tan fundamentales. Me gustaría, sin embargo, recomendar a los lectores algunos temas para seguir investigando en el área, que no pude recoger en este libro.

Respecto al capítulo 1, algunos temas para explorar que no vimos en el libro son: las condicionalizaciones alternativas como el método de Jeffrey (1992), la teoría de la confirmación bayesiana en ciencia empírica (Earman 1992), y obviamente un análisis más profundo de la filosofía de la probabilidad (Rowbottom 2015).

Respecto al capítulo 2, muchos filósofos se han interesado por la paradoja de Newcomb y temas similares como la teoría de la decisión causal (Weirich 2008). También es interesante explorar más relaciones entre decisión y psicología, como el efecto del “framing” (Bermúdez 2008). Otro tema que interesa a los filósofos son las posibles violaciones de axiomas de la preferencia, como la transitividad (Andreou 2022).

Respecto al capítulo 3, hay muchos temas relacionados con la teoría de juegos que quedaron afuera del libro, por motivos de espacio o dificultad. Esto incluye la selección de equilibrios (Harsanyi & Selten 1988), los métodos de negociación (Vanderschraaf 2023), y los métodos de división justa de bienes (Brams & Taylor 2011).

Finalmente, respecto al capítulo 4, quedaron sin explorar los enfoques más cardinales, como el utilitarismo de Harsanyi (1955), y también otros teoremas fundamentales, como el de Gibbard-Satterthwaite; recomiendo el libro de Gaertner (2006) a los que quieran profundizar en el tema. Respecto a la agregación de juicios, existe también una interesante discusión sobre cómo llegar a acuerdos a partir de opiniones probabilísticas (Elkin & Pettigrew 2025).

Un tema en el que no enfaticé, pero me parece de mucho interés, es la historia intelectual de los temas tratados en el libro. Por ejemplo, de qué modo autores como Arrow, Ellsberg o Schelling participaron en la construcción de la hegemonía militar norteamericana (Erickson 2015). Esto conlleva también una discusión ética sobre el valor de las contribuciones intelectuales más allá de su función original.

Por último, no profundicé en los métodos de simulaciones en filosofía, que a partir de la contribución de autores como Axelrod y Skyrms, se volvieron una práctica usual. Recomiendo a los lectores prestar atención a esta área de investigación y sus recientes exploraciones sobre la desinformación (O'Connor & Weatherall 2019) y el origen de las normas sociales (O'Connor 2019).

Espero que los lectores hayan disfrutado tanto la lectura de este libro como yo disfruté de su escritura.

BIBLIOGRAFÍA

Adams, E. (1998). *A primer in Probability Logic*. CSLI Publications.

Allais, M. (1953a). “Le comportement de l’homme rationnel devant le risque: critique des postulats et axiomes de l’école Américaine”. *Econometrica* 21(4): 503–546.

Allais, M. (1953b). “La psychologie de l’homme rationnel devant le risque : la théorie et l’expérience”, *Journal de la société statistique de Paris* 94:47-73.

Andreoni, J. & Miller, J. (1993). “Rational cooperation in the finitely repeated prisoner dilemma: experimental evidence”. *The Economic Journal* 103(418): 570-585.

Andreou, C. (2022). *Choosing Well*. Oxford University Press.

Arrow, K. (1951). *Social choice and individual values*. Yale University Press. [Traducción al español: *Elección social y valores individuales*. Planeta, 1994]

Axelrod, R. (1984). *The Evolution of Cooperation*, Basic Books. [Traducción al español: *La Evolución de la Cooperación*. Alianza, 1986]

Barenstein, J. (2013). “Los escritos electorales de Ramón Lull: una nueva teoría de la votación en la segunda mitad del s. XIII”, *Revista Española de Filosofía Medieval* 20: 85-99.

Basu, K. (1994). “The Traveler Dilemma: Paradoxes of Rationality in Game Theory”. *The American Economic Review* 84(2): 391-395.

Basu, K., Becchetti, L. & Stanca, L. (2011). “Experiments with the Traveler’s Dilemma: welfare, strategic choice and implicit collusion”. *Social Choice and Welfare* 37: 575–595.

Becker, G. & K. Murphy (1988). “A Theory of Rational Addiction”, *Journal of Political Economy* 96(4):675–700.

- Bentham, J. (1780). *Introduction to the Principles of Morals and Legislation*. [Traducción al español: *Los principios de la moral y la legislación*. Claridad, 2008]
- Bermúdez, J. (2008). *Frame it again: new tools for rational decision-making*. Cambridge University Press.
- Bicchieri, C. (2006). *The Grammar of Society*. Cambridge University Press.
- Binmore, K. (1987). “Modelling Rational Players: Part I”. *Economics and Philosophy* 3(2): 179-214.
- Binmore, K. (1998). *Game Theory and the Social Contract. Vol. II: Just Playing*. The MIT Press.
- Binmore, K. (2007). *Playing for Real: A Text on Game Theory*. Oxford University Press.
- Binmore, K. (2015). “Why all the fuss? The many aspect of the Prisoner’s Dilemma”, en Peterson, M. (ed.) *The Prisoner’s Dilemma*, Cambridge University Press, pp. 16-34.
- Bermúdez, J.L. (2009). *Decision Theory and Philosophy*. Oxford University Press.
- Black, D. (1948). “On the Rationale of Group Decision Making”. *Journal of Political Economy* 56: 23-34.
- Bonanno, G. (2015). *Game Theory*. Disponible en la página web del autor.
- Brams, S. & A. Taylor (1996). *Fair Division*. Cambridge University Press.
- Brennan, J. (2018). *Against Democracy*. Princeton University Press.
- Brocas, I. & J. Carrillo (2025). “Why do children pass in the centipede game? Cognitive limitations v. risk calculations”. *Games and Economic Behavior* 150: 295-311.
- Broome, J. (1999). *Ethics out of economics*. Cambridge University Press.
- Buchak, L. (2013). *Risk and rationality*. Oxford University Press.
- Callard, A. (2018). *Aspiration: the Agency of Becoming*. Oxford University Press.

- Condorcet, M. (1875). *Essai sur l'application de l'analyse à la probabilité des décisions rendues à la pluralité des voix*. Royale.
- Cresto, E. & D. Tajer, (2020). "Confirmational holism and Theory Choice: Arrow meets Duhem", *Mind* 129(513):71-111.
- De Bona, G. & Staffel, J. (2018). "Why be coherent", *Analysis* 78(3): 405-415.
- Dawkins, R. (1989). *The Selfish Gene* (2nd ed). Oxford University Press. [Traducción al español: *El Gen Egoísta (extendido)*. Ed. Salvat, 2017]
- De Finetti, B. (1970). *Theory of Probability: A Critical Introductory Treatment*. New York: John Wiley.
- Dietrich, F. & C. List (2007). "Arrow's theorem in Judgment Aggregation", *Social Choice and Welfare* 29: 19-33.
- Dietrich, F. & C. List (2016). "Mentalism versus behavioralism in economics: a philosophy-of-science perspective", *Economics and Philosophy* 32(2): 249-281.
- Earman, J. (1992). *Bayes or Bust? A Critical Examination of Bayesian Confirmation Theory*. The MIT Press.
- Elkin, L. & R. Pettigrew (2025). *Opinion Pooling*. Cambridge University Press.
- Ellsberg, D. (1961). "Risk, ambiguity, and the Savage Axioms". *The Quarterly Journal of Economics* 75(4): 643-669.
- Erickson, P. (2015). *The World that Game Theorists Made*. University of Chicago Press.
- Estlund, D. (2008). *Democratic Authority: A Philosophical Framework*. Princeton University Press. [Traducción al español: *La autoridad democrática*, Ed. Siglo XXI, 2011]
- Fey, M. (2014). "A straightforward proof of Arrow's theorem", *Economics Bulletin* 34(3): 1792-1797.
- Fitelson, B. (inédito). *Coherence*. Disponible en la página web del autor.
- Gaertner, W. (2006). *A Primer in Social Choice Theory*. Cambridge University Press.

- Gauthier, D. (1986). *Morals by Agreement*. Clarendon Press. [Traducción al español: *La moral por acuerdo*. Gedisa, 2000]
- Gauthier, D. (2015). “How I Learned to Stop Worrying and Love the Prisoner’s Dilemma”, en Peterson, M. (ed.) *The Prisoner’s Dilemma*, Cambridge University Press, pp. 35-53.
- Gibbard, A. (1974). “A Pareto-Consistent Libertarian Claim”. *Journal of Economic Theory* 7: 388-410.
- Goodin, R. & K. Spiekermann (2019). *An epistemic theory of democracy*, Oxford University Press.
- Grofman, B. & S. Feld (1988). “Rousseau’s General Will: a Condorcetian perspective”, *The American Political Science Review* 82(2): 567-576.
- Gustaffson, J. (2022). *Money-pump Arguments*. Cambridge University Press.
- Hacking, I. (2001). *An introduction to probability and inductive logic*. Cambridge University Press.
- Halpern, J. (2003). *Reasoning about Uncertainty*. The MIT Press.
- Harsanyi, J. (1955). “Cardinal Welfare, Individualistic Ethics, and Interpersonal Comparisons of Utility”. *Journal of Political Economy* 63(4): 309-321.
- Harsanyi, J. & R. Selten (1988). *A General Theory of Equilibrium Selection in Games*. The MIT Press.
- Hausman, D. & M. McPherson (2009). “Preference Satisfaction and Welfare Economics”, *Economics and Philosophy* 25: 1-25.
- Hausman, D. (2012). *Preference, Value, Choice, and Welfare*. Cambridge: Cambridge University Press.
- Hintikka, J. (1963). *Knowledge and Belief: a logic of the two notions*. [Traducción al español: *Saber y Creer: una introducción a la lógica de las dos nociones*. Tecnos, 1979.]
- Hobbes, T. (1651). *Leviatán*. [Hay varias traducciones]
- Hume, D. (1739). *Tratado de la Naturaleza Humana*. Traducción de Félix Duque. Tecnos, 1988.

- Jeffrey, R. (1992). *Probability and the Art of Judgment*. Cambridge University Press.
- Kahneman, D. & A. Tversky (1979). "Prospect Theory: an analysis of Decision under Risk", *Econometrica* 47(2): 263-292.
- Keynes, J.M. (1921). *A treatise on Probability*. London: Macmillian.
- Kolodny, N. (2007). "How does coherence matter", *Proceedings of the Aristotelian Society* 107(1):229-263.
- Kornhauser, L.A. & L.G. Sager, "Unpacking the court" (1986) *Yale Law Journal* 96: 82-117.
- Kripke, S. (1963). "Semantical Considerations on Modal Logic", *Acta Philosophica Fennica* 16: 83-94.
- Kripke, S. (1980). *Naming and Necessity*. Harvard University Press. [Traducción al español: *El nombrar y la necesidad*, UNAM, 1995]
- Kyburg, H. (1961). *Probability and the Logic of Rational Belief*. Wesleyan University Press.
- Leitgeb, H. (2014). "The stability theory of belief". *The Philosophical Review* 123(2): 131-171.
- Levi, I. (1974). "On indeterminate probabilities". *Journal of Philosophy* 81(71): 391-418.
- Lewis, D. (1969). *Convention*. Harvard University Press.
- Lewis, D. (1986). "A subjectivists Guide to Objective Chance", en *Philosophical Papers II*, Oxford University Press, 83-132.
- List, C. & J. Dryzek (2003). "Social Choice Theory and Deliberative Democracy: A Reconciliation", *British Journal of Political Science* 33(1): 1-28.
- List, C. & P. Pettit, (2002). "Aggregating sets of judgments: an impossibility result", *Economics and Philosophy* 18(1): 89-110.
- Luce, D. & H. Raiffa (1957). *Games and Decisions*. Wiley.
- MacCrimmon, K. & S. Larsson (1978). "Utility Theory: Axioms Versus 'Paradoxes'", en M. Allais & O. Hagen (Eds.) *Expected Utility Hypotheses and the Allais Paradox*, Springer, pp. 333-409.

- Mackie, G. (2003). *Democracy Defended*. Cambridge University Press.
- Makinson, D. (1965). "The paradox of the Preface". *Analysis* 25(6): 205-207.
- Mill, JS. (1861) *Utilitarianism*. [Hay diversas traducciones al español.]
- Narens, L. & B. Skyrms (2020). *The pursuit of happiness*. Oxford University Press.
- Nash, J. (1950). "Equilibrium Points in n-Person Games," *Proceedings of the National Academy of Sciences*, 36, 48-49.
- Nozick, R. (1974). *Anarchy, State and Utopia*. [Traducción al español: *Anarquía, Estado y Utopía*. Fondo de Cultura Económica, 1988]
- O'Connor, C. (2019). *The origins of unfairness*. Cambridge University Press.
- O'Connor, C. & B. Weatherall, J. (2019). *The Misinformation Age*. Yale University Press.
- Okasha, S. (2011). "Theory choice and social choice: Kuhn versus Arrow", *Mind* 120(477): 83-115.
- Okasha, S. (2016). "On the Interpretation of Decision Theory", *Economics and Philosophy* 32(3): 409-433.
- Paul, L.A. (2014). *Transformative Experience*. Oxford University Press.
- Peterson, M. (2009). *An Introduction to Decision Theory*. Cambridge University Press.
- Peterson, M. (2015). *The Prisoner Dilemma*. Oxford University Press.
- Pettigrew, R. (2019). *Choosing for Changing Selves*. Oxford University Press.
- Pettit, P. & Rabinowicz, W. (2001). "Deliberative Democracy and the Discursive Dilemma", *Philosophical Issues* 11: 268-299.
- Popper, K. (1983). *Realism and the Aim of Science*. Hutchinson. [Traducción al español: *Realismo y el Objetivo de la Ciencia*, Tecnos, 2011]
- Przeworski, A. (2010). *Democracy and the limits of self-government*. Cambridge University Press. [Traducción

- al español: *Qué esperar de la democracia: límites y posibilidades del autogobierno*, Ed. Siglo XXI, 2010].
- Ramsey, F. (1926). "Truth and Probability", en R.B. Braithwaite (ed.) *Foundations of Mathematics and other Logical Essays*, Ed. Kegan, Paul, 1931, pp. 156-198.
- Rapoport, A., Seale, D. & A. Colman (2015). "Is tit-for-tat the answer? On the conclusions drawn from Axelrod's tournaments", *Plos One*.
- Riker, W. (1987). *Liberalism Against Populism*. Waveland Inc.
- Rosenthal, R. (1981). "Games of Perfect Information, Predatory Pricing, and the Chain Store". *Journal of Economic Theory*. **25** (1): 92–100.
- Rousseau, J.J. (1762). *El Contrato Social*. Traducción de Leticia Halperín Donghi. Losada, 2003.
- Rowbottom, P. (2015). *Probability*. Polity.
- Saari, D. (1997). "Are individual rights possible", *Mathematics Magazine* 70(2): 83-92.
- Samuelson, P. (1938). "A Note on the Pure Theory of Consumer's Behaviour". *Economica* 5(17): 61-71.
- Savage, L. (1954). *Foundations of Statistics*. Wiley.
- Schelling, T. (1960). *The Strategy of Conflict*. Harvard University Press. [Traducción al español: *La Estrategia del Conflicto*. Tecnos, 1964]
- Schumpeter, J. (1942). *Socialism, Capitalism, and Democracy*. Harper and Brothers. [Traducción al español: *Socialismo, Capitalismo y Democracia*, Orbis, 1983]
- Selten, R. (1975). "Reexamination of the Perfectness Concept for Equilibrium Points in Extensive Games". *International Journal of Game Theory* 4(1): 25-55.
- Sen, A. (1970a). "The impossibility of a Paretian liberal", *Journal of Political Economy* 78(1): 152-157.
- Sen, A. (1970b). *Collective Choice and Social Welfare*. Harvard University Press.
- Sen, A. (1971). "Choice Functions and Revealed Preference". *The Review of Economic Studies* 38(3): 307-317.
- Sen, A. (1973). "Behaviour and the Concept of Preference". *Economica* 40(159): 241–259.

- Sen, A. (1975). "Is a Paretian Liberal really impossible? A Reply", *Public Choice* 21: 111-113.
- Sen, A. (1977). "Rational Fools: A Critique of the Behavioural Foundations of Economic Theory". *Philosophy and Public Affairs* 6(4): 317-344.
- Sen, A. (1996). "Rights: formulations and consequences", *Analyse & Kritik* 18: 153-170.
- Skyrms, B. (2002). *The Stag Hunt and the Evolution of Social Structure*. Cambridge University Press. [Traducción al español: *La caza del ciervo*, Melusina, 2007]
- Thoma J. (2021). "In defense of revealed preference theory". *Economics and Philosophy* 37: 163-187.
- Ullmann-Margalit, E. & S. Morgenbesser (1977). "Picking and choosing", *Social Research* 44(4): 757-785.
- Vanderschraaf, P. (2019). *Strategic Justice*. Oxford University Press.
- Vanderschraaf, P. (2025). *Bargaining Theory*. Cambridge University Press.
- Von Mises, R. (1928). *Probability, Statistics and Truth*. The MacMillan Company.
- Von Neumann, J., & Morgenstern, O. (1944). *Theory of games and economic behavior*. Princeton University Press.
- Weirich, P. (2008). "Causal Decision Theory". *Stanford Encyclopedia of Philosophy*.

Impreso en noviembre de 2025
en los talleres gráficos de Elías Porter
Buenos Aires, Argentina.