



UNIVERSIDAD  
DE GRANADA

HERRAMIENTAS PARA GARANTIZAR  
JUSTICIA EN APRENDIZAJE AUTOMÁTICO

DANIEL BOLAÑOS MARTÍNEZ

Trabajo Fin de Grado

Doble Grado en Ingeniería Informática y Matemáticas

**Tutores**

Jorge Casillas Barranquero

Pedro González Rodelas

FACULTAD DE CIENCIAS

E.T.S. INGENIERÍAS INFORMÁTICA Y DE TELECOMUNICACIÓN

*Granada, a 23 de octubre de 2021*

# Herramientas para garantizar justicia en aprendizaje automático

Daniel Bolaños Martínez

Daniel Bolaños Martínez. *Herramientas para garantizar justicia en aprendizaje automático.*

Trabajo de fin de Grado. Curso académico 2021-2022.

**Responsables de  
tutorización**

Jorge Casillas Barranquero  
*Ciencias de la Computación  
e Inteligencia Artificial*

Pedro González Rodelas  
*Matemática Aplicada*

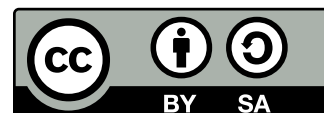
Escuela Técnica Superior  
de Ingeniería Informática  
y Telecomunicación

Facultad de Ciencias

Doble Grado en  
Ingeniería Informática y  
Matemáticas

Universidad de Granada

This work is licensed under a [Creative Commons](#) “[Attribution-ShareAlike 4.0 International](#)” license.



The source code of this text and developed programs are available in the Github repository [danibolanos/TFG-Guarantee\\_Fairness\\_in\\_ML](#)

#### DECLARACIÓN DE ORIGINALIDAD

D. Daniel Bolaños Martínez

Declaro explícitamente que el trabajo presentado como Trabajo de Fin de Grado (TFG), correspondiente al curso académico 2021-2022, es original, entendida esta, en el sentido de que no ha utilizado para la elaboración del trabajo fuentes sin citarlas debidamente.

En Granada a 23 de octubre de 2021

Fdo: Daniel Bolaños Martínez

---

## AGRADECIMIENTOS

---

Me gustaría agradecer a . . .

---

## RESUMEN

---

Los modelos de aprendizaje automático tienen un impacto cada vez mayor en el mundo actual, siendo utilizado para asistir y a veces sustituir a los humanos en muchos entornos. Estos modelos, a menudo funcionan aprendiendo sobre decisiones históricas tomadas por distintos grupos sociales con diferentes tasas de error en su clasificación, por lo que surge la necesidad de vigilar y controlar el sesgo involuntario de los modelos predictivos contra estos grupos de población desfavorecidos.

La mayoría de las herramientas utilizadas para mitigar el sesgo entre los grupos no privilegiados suelen depender del modelo y la métrica utilizada, aumentando así los requerimientos de los procesos en las técnicas de aprendizaje automático empleadas. Nuestro objetivo será encontrar un equilibrio entre un concepto de equidad que nos permita eliminar en mayor medida la desigualdad de rendimiento entre grupos, pero que a su vez no aumente el nivel de complejidad del modelo.

En este trabajo, formalizaremos matemáticamente las definiciones de diferentes medidas de justicia y equidad; presentaremos sus propiedades, limitaciones e incompatibilidades e indagaremos en las posibles opciones para mejorar los resultados obtenidos mediante un proceso de aprendizaje automático en los términos planteados previamente: equidad por desconocimiento, paridad estadística/demográfica, medidas causales, equidad individual, entre otras.

Nos enfocaremos en desarrollar un marco para modelar la equidad usando herramientas de inferencia causal tomando como base la teoría de probabilidad y estadística. Discutiendo el concepto de equidad contrafactual que plantea que una decisión es justa para un individuo si es la misma en el mundo real y en un mundo "contrafactual" en el que el individuo perteneciese a un grupo demográfico diferente.

Finalmente haremos un análisis de las utilidades que ofrece el software Aequitas para garantizar justicia en aprendizaje automático. Realizaremos algunas pruebas para la evaluación de las medidas equidad definidas y lo complementaremos con la resolución de un problema del mundo real utilizando herramientas de equidad contrafactual en *Python*. En este contexto, aportaremos métodos de presentación y visualización de resultados que faciliten una explicación de la causa del sesgo y ayuden a los usuarios a tomar decisiones en el mundo real.

**PALABRAS CLAVE:** medidas de equidad, mitigación del sesgo, impacto dispar, equidad contrafactual, Aequitas.

---

## ABSTRACT

---

Machine learning is having an increasing impact in today's world, being used to assist and sometimes replace humans in many environments. . . .

**KEY WORDS:** fairness metrics, bias mitigation, disparate impact, counterfactual fairness, Aequitas.



---

## ÍNDICE GENERAL

---

1.	INTRODUCCIÓN	12
1.1.	Análisis del problema . . . . .	13
1.2.	Objetivos del trabajo . . . . .	14
1.3.	Contribuciones . . . . .	15
1.4.	Esquema general . . . . .	15
I.	PRINCIPIOS MATEMÁTICOS BÁSICOS	16
2.	TEORÍA DE LA PROBABILIDAD	17
2.1.	Espacio de probabilidad . . . . .	17
2.2.	Variables aleatorias . . . . .	18
2.2.1.	Distribución conjunta . . . . .	20
2.2.2.	Distribución marginal . . . . .	20
2.3.	Probabilidad condicional . . . . .	21
2.3.1.	Distribución condicional . . . . .	22
2.4.	Independencia . . . . .	23
3.	DISTRIBUCIONES DE PROBABILIDAD	24
3.1.	Esperanza de una variable aleatoria . . . . .	24
3.2.	Momentos de una variable aleatoria . . . . .	26
3.3.	Ejemplos de distribuciones . . . . .	27
3.3.1.	Distribución Bernoulli . . . . .	28
3.3.2.	Distribución de Poisson . . . . .	29
3.3.3.	Distribución Normal . . . . .	31
3.3.4.	Distribución Uniforme . . . . .	33
3.3.5.	Distribución Gamma . . . . .	34
4.	ESTADÍSTICA PARAMÉTRICA	35
4.1.	Propiedades del estimador . . . . .	36
4.2.	Criterios de evaluación . . . . .	39
5.	TEORÍA DE GRAFOS	41
5.1.	Grafos, nodos y aristas . . . . .	41
5.2.	Estructura de un grafo . . . . .	42
II.	JUSTICIA EN APRENDIZAJE AUTOMÁTICO	44
6.	CONCEPTOS BÁSICOS DEL APRENDIZAJE AUTOMÁTICO	45
6.1.	¿Qué es el aprendizaje automático? . . . . .	45
6.1.1.	Aprendizaje supervisado . . . . .	45
6.2.	Propiedades del modelo de aprendizaje . . . . .	47
6.3.	Creación de modelos de aprendizaje . . . . .	49
6.3.1.	Ejemplo: Perceptrón . . . . .	51

6.3.2.	Regresión lineal . . . . .	53
6.4.	Evaluación en aprendizaje automático . . . . .	55
7.	FORMALIZACIÓN DE LAS MEDIDAS DE EQUIDAD . . . . .	58
7.1.	¿Qué es la equidad? . . . . .	58
7.1.1.	Principales familias de las medidas de equidad . . . . .	59
7.1.2.	Medición de la parcialidad y la equidad . . . . .	60
7.2.	Equidad por desconocimiento . . . . .	60
7.3.	Equidad individual . . . . .	61
7.4.	Equidad de grupo . . . . .	63
7.4.1.	Paridad demográfica . . . . .	66
7.4.2.	Probabilidades igualadas . . . . .	67
7.4.3.	Tasa de paridad predictiva . . . . .	69
7.4.4.	Medidas basadas en la puntuación . . . . .	70
7.4.5.	Igualdad de las métricas de predicción . . . . .	71
7.4.6.	Impacto desigual . . . . .	71
8.	ALGORITMOS DE MITIGACIÓN DE SESGO . . . . .	73
8.1.	Modelos de aprendizaje justos . . . . .	73
8.1.1.	Selección de los datos del modelo . . . . .	73
8.1.2.	Equilibrio entre equidad y métricas de evaluación . . . . .	74
8.2.	Algoritmos de preprocesamiento . . . . .	74
8.2.1.	Ejemplo: Aprendizaje de la representación justa . . . . .	75
8.3.	Algoritmos de optimización durante el entrenamiento . . . . .	77
8.3.1.	Ejemplo: Aprendizaje en clasificación sin impacto dispar. . . . .	78
8.4.	Algoritmos de posprocesamiento . . . . .	81
8.4.1.	Ejemplo: Aprendizaje en igualdad de oportunidades . . . . .	81
III.	FUNDAMENTOS DE LA EQUIDAD CONTRAFACTUAL . . . . .	83
9.	INFERENCIA CAUSAL . . . . .	84
9.1.	Modelos causales . . . . .	84
9.1.1.	Ejemplo: Construcción de un modelo causal . . . . .	84
9.1.2.	Formalización de los modelos causales estructurales . . . . .	86
9.2.	Grafos causales . . . . .	87
9.2.1.	Forks . . . . .	87
9.2.2.	Colliders . . . . .	88
9.2.3.	Mediador . . . . .	89
9.3.	Intervención y confusión . . . . .	89
9.3.1.	Operadores para realizar actuaciones en el modelo . . . . .	89
9.3.2.	Confusión entre dos variables . . . . .	90
10.	TEOREMA DE IMPOSIBILIDAD DE LA EQUIDAD . . . . .	93
10.1.	Caracterización del teorema . . . . .	93
10.1.1.	Paridad demográfica versus Tasa de paridad predictiva . . . . .	93
10.1.2.	Paridad demográfica versus Probabilidades igualadas . . . . .	95
10.1.3.	Probabilidades igualadas versus Tasa de paridad predictiva . . . . .	97

10.1.4. Enunciado y demostración . . . . .	99
11. MEDIDAS CAUSALES . . . . .	100
11.1. Contrafactuales . . . . .	100
11.1.1. Ejemplo: Modelo de decisión contrafactual . . . . .	100
11.1.2. Formalización del cálculo contrafactual . . . . .	102
11.2. Equidad contrafactual . . . . .	103
11.2.1. Implicaciones de la definición de equidad . . . . .	104
IV. ANÁLISIS EXPERIMENTAL . . . . .	106
12. DESCRIPCIÓN Y DISEÑO . . . . .	107
12.1. Algoritmo de aprendizaje justo . . . . .	107
12.2. Diseño del modelo causal de entrada . . . . .	108
12.3. Aplicación en un problema real . . . . .	109
12.3.1. Descripción del problema . . . . .	109
12.3.2. Escenarios de predicción . . . . .	109
13. IMPLEMENTACIÓN Y RESULTADOS . . . . .	111
13.1. Obtención y tratamiento de los datos . . . . .	111
13.2. Implementación del código . . . . .	111
13.3. Contraste de los resultados . . . . .	112
13.4. Condiciones de la experimentación . . . . .	112
13.4.1. Entorno de ejecución . . . . .	112
13.4.2. Entorno de programación . . . . .	112
13.4.3. Bibliotecas y herramientas auxiliares . . . . .	112
13.5. Manual de ejecución del experimento . . . . .	112
V. CONCLUSIONES Y VÍAS FUTURAS . . . . .	113
14. CONCLUSIÓN . . . . .	114
15. TRABAJOS FUTUROS . . . . .	115
APÉNDICES . . . . .	116
A. HERRAMIENTAS PARA GARANTIZAR JUSTICIA EN AA . . . . .	117
A.1. Aequitas . . . . .	117
A.1.1. Estructura de los datos de entrada y resultados . . . . .	117
A.1.2. Métricas usadas por Aequitas . . . . .	119
B. ESTIMACIÓN DEL COSTE Y PLANIFICACIÓN . . . . .	122
B.1. Estimación del presupuesto del proyecto . . . . .	122
B.2. Planificación del trabajo . . . . .	122
NOTACIÓN . . . . .	126
BIBLIOGRAFÍA . . . . .	130

---

## INTRODUCCIÓN

---

Actualmente, los algoritmos de aprendizaje automático se utilizan en ámbitos diversos con un gran impacto en la sociedad, como son: el proceso de concesión de créditos bancarios (Fuster et al. [2018]), selección de personal para un puesto de trabajo (Miller [2015]) o decisión de una condena en justicia penal (Angwin et al. [2016]). Estos ejemplos de aplicación son propensos a la discriminación, la cual está prohibida por la legislación internacional (Title VII of the Civil Rights Act: Equal Employment Opportunities). Las causas de las disparidades en los sistemas de aprendizaje automático provienen esencialmente del sesgo humano que existe en los conjuntos de datos de entrenamiento debido a razones históricas. Algunas de las posibles causas (Barocas and Selbst [2016]) son:

- Los sistemas de aprendizaje automático mantienen la discriminación existente en los datos antiguos debido al sesgo humano. Por ejemplo, si un sistema de contratación a la hora de seleccionar a los aspirantes al cargo utiliza como etiquetas de predicción las decisiones tomadas por un directivo en lugar de sus capacidades reales, en la mayoría de los casos se podrían rechazar candidatos con un alto nivel de rendimiento.
- Normalmente la cantidad de ejemplos y la información que ofrecen sus características son menores para el grupo minoritario, por lo que es menos probable que se modelen correctamente con respecto a los individuos del grupo mayoritario. Esto tendrá desventajas a la hora de predecir sobre nuevos ejemplos del grupo desfavorecido.
- Si ya existe un sesgo inicial, probablemente se agravará con el tiempo. Por ejemplo, en el registro policial de delitos solo constan delitos observados por la policía. El departamento de policía tiende a enviar más agentes a lugares donde se ha detectado una mayor tasa de delincuencia inicialmente y, por tanto, será más probable que se detecten delitos en esas regiones.
- Aunque los atributos sensibles no se utilicen en el entrenamiento de un sistema de aprendizaje automático, podrán existir atributos derivados de estos. Si se incluyen estas características, el sesgo seguirá presente. A veces, es muy difícil determinar si un atributo relevante está correlacionado con los atributos sensibles y si debemos incluirlo o no en el proceso de entrenamiento.

Estos problemas han llevado al desarrollo de numerosas investigaciones sobre la equidad en el ámbito del aprendizaje automático, enfocadas en cómo surge la discriminación, cómo puede medirse y cómo puede mitigarse. El objetivo de la equidad será por tanto, diseñar algoritmos que hagan predicciones justas, evitando perjudicar a un determinado grupo de la población.

A pesar de estos avances, en algunos escenarios, la discriminación en los modelos sigue siendo difícil de abordar y mucho más de entender debido principalmente a:

- Las aplicaciones generalmente actúan como modelos de caja negra en las cuales tenemos acceso restringido al clasificador por cuestiones de privacidad o derechos de propiedad intelectual de los datos utilizados (Diakopoulos [2014]). Por ejemplo, si estamos usando una API de predicción.
- Los modelos se despliegan en una población donde la distribución de los datos no refleja los patrones contenidos en los datos de entrenamiento (Sugiyama et al. [2017]). Por ejemplo, podrían aparecer diferencias entre los individuos del grupo debido a un cambio en la distribución de la población de interés.

En este trabajo, estudiaremos cómo se formaliza el concepto de equidad en el ámbito del aprendizaje automático y presentaremos estas formalizaciones en las diferentes ramas de la ingeniería informática y las matemáticas. Las medidas de equidad parten de las ideas de justicia distributiva del ámbito de las ciencias sociales y, por ello, pueden aparecer conceptos que entren en conflicto con otros, por lo que las predicciones producidas por los algoritmos y modelos que las utilizan también diferirán enormemente.

Desde el punto de vista práctico, es importante estudiar estos criterios de equidad y sus implicaciones realizando un análisis teórico y empírico a partir de las nociones de la literatura de las ciencias sociales. Este análisis tendrá la intención de ayudar a determinar la bondad de las formalizaciones existentes y poder valorar las ventajas e inconvenientes de cada criterio con el objetivo de mejorar o construir nuevas formalizaciones de equidad en un futuro.

## 1.1 ANÁLISIS DEL PROBLEMA

La equidad en el aprendizaje automático, tiene como objetivo estudiar y mitigar la discriminación en los procesos de toma de decisiones algorítmicas. Actualmente podemos encontrar tres enfoques dentro de la mitigación de los sesgos. Los métodos de preprocesamiento que intentan corregir los sesgos de los datos introducidos en el modelo, los métodos de procesamiento interno que tratan de realizar la corrección de los sesgos producidos durante el proceso de aprendizaje y finalmente, los métodos de posprocesamiento que tienen como objetivo corregir el resultado de un modelo sesgado.

A estos enfoques se le suman la gran cantidad de criterios de justicia que surgen a partir de la literatura de las ciencias sociales aplicadas al aprendizaje automático. Po-

demostramos hacer una primera aproximación a estas familias y por ende a la formalización de los conceptos de equidad, intentando dar respuesta a las siguientes dos preguntas:

- ¿Paridad o preferencia?: si buscamos equidad para lograr una paridad entre los individuos del grupo o en cambio queremos satisfacer unas preferencias dentro del mismo.
- ¿Tratamiento o impacto?: si tratamos de mantener la equidad durante el tratamiento de los datos o por el contrario en los resultados producidos por el modelo (impacto).

Dando respuesta a las preguntas anteriores, surgen los siguientes criterios de equidad que se resumen en el Cuadro 1 (Gajane and Pechenizkiy [2018]).

	Paridad	Preferencia
Tratamiento	Equidad por desconocimiento Medidas causales	Tratamiento preferente
Impacto	Equidad de grupo Equidad individual	Impacto preferente

Cuadro 1: Formalización de los criterios de equidad.

Debido a la gran cantidad de opciones que se nos presentan, nos gustaría saber qué nociones de equidad tienen menos limitaciones en la práctica y qué método de mitigación de sesgo aporta menos complejidad en los algoritmos de predicción. De esta forma, podremos conocer las virtudes de cada familia de equidad y ser capaces de desarrollar nuevos conceptos o adaptar cada uno al caso de estudio que mejores resultados nos proporcione.

## 1.2 OBJETIVOS DEL TRABAJO

Nuestro objetivo principal es hallar modelos en el ámbito del aprendizaje automático que tengan una buena relación equidad-precisión. Queremos encontrar una mitigación del sesgo que no aumente la complejidad de los modelos utilizados, no necesite conocimiento específico del dominio y tenga el menor número de limitaciones en su aplicación sobre problemas del mundo real.

Abordaremos estas tareas realizando un análisis teórico de las nociones de equidad y de los algoritmos construidos en base a ellas. Estudiaremos qué criterios aportan más beneficios sobre un modelo de caja negra del que únicamente podamos conocer los datos antes o después de ser procesados, e indagaremos sobre los métodos de optimización de las nociones de justicia con el propósito de satisfacer la equidad en el modelo minimizando la pérdida de rendimiento del mismo.

Como resultado de los objetivos planteados, sugerimos una revisión bibliográfica exhaustiva de los métodos de mitigación del sesgo y de las medidas de equidad

que se nos presentan en diferentes artículos como [Gajane and Pechenizkiy \[2018\]](#) o [Verma and Rubin \[2018\]](#). Nos proponemos detallar las herramientas matemáticas que nos serán útiles para la formalización de las medidas de equidad, haciendo especial hincapié en la inferencia causal como base sobre la que se construye el modelo de equidad contrafactual que estudiaremos en este trabajo.

Finalmente realizaremos un análisis de la herramienta de software Aequitas ([Saleiro et al. \[2019\]](#)) dedicado al estudio de equidad sobre problemas del mundo real. Haremos pruebas sobre un conjunto de datos y evaluaremos sus carencias respecto a la teoría. Además replicaremos un ejemplo utilizando la noción de equidad contrafactual ([Kusner et al. \[2018\]](#)) en el lenguaje de programación Python y aportaremos herramientas de visualización para facilitar la interpretación de los resultados a otros investigadores o científicos de datos interesados en el tema.

### 1.3 CONTRIBUCIONES

En resumen, las contribuciones principales de este proyecto son:

- Discutir y formalizar las distintas familias de medidas de equidad y algoritmos de mitigación de sesgo de mayor interés en este ámbito.
- Facilitar una demostración alternativa del teorema de imposibilidad para mostrar las incompatibilidades entre los diferentes criterios de equidad de grupo.
- Proporcionar un análisis empírico utilizando el software Aequitas para las evaluaciones de las medidas de justicia estudiadas teóricamente.
- Replicar un ejemplo del mundo real sobre un modelo causal en el lenguaje de programación Python basándonos en el concepto de equidad contrafactual, disponible en: [danibolanos/TFG-Guarantee\\_Fairness\\_in\\_ML](#)
- Implementar gráficas en Python para la interpretación de los resultados al aplicar las técnicas de equidad contrafactual.
- **Comparar nociones de equidad/Esquema prototipo.**

### 1.4 ESQUEMA GENERAL

#### Completar y revisar al final

En la Sección , formalizaremos los conceptos de equidad para cada una de las familias propuestas por la literatura del aprendizaje automático. En la Sección X discutiremos las limitaciones de algunos conceptos vistos en la sección anterior. En la Parte [iv](#) crearemos y discutiremos herramientas para el análisis de los conceptos de equidad y analizaremos empíricamente cada concepto en base a herramientas de visualización de sus resultados. Por último, en la Parte [v](#), discutiremos diferentes vías para la formalización de conceptos de equidad futuros.

## Parte I

# PRINCIPIOS MATEMÁTICOS BÁSICOS

Definiciones y resultados relativos a teoría de probabilidad, estadística y teoría de grafos.



---

## TEORÍA DE LA PROBABILIDAD

---

En este capítulo incluiremos definiciones y resultados previos de la teoría de probabilidad que utilizaremos a lo largo del desarrollo del trabajo. La fuente principal utilizada en este capítulo parte del trabajo contenido en Dembo [2014].

### 2.1 ESPACIO DE PROBABILIDAD

Construiremos la teoría asumiendo que existe un conjunto no vacío  $\Omega$  que representa al conjunto de todos los posibles resultados de un experimento. Llamaremos *suceso* a cualquier subconjunto de  $\Omega$ .

**Definición 1** ( $\sigma$ -álgebra). Sea  $\mathcal{P}(\Omega)$  el conjunto de partes de  $\Omega$ . Llamaremos  $\sigma$ -álgebra a  $\mathcal{A} \subset \mathcal{P}(\Omega)$  que satisfaga:

- $\mathcal{A}$  contiene al conjunto vacío:  $\emptyset \in \mathcal{A}$ .
- $\mathcal{A}$  es cerrado bajo complementarios: si  $A \in \mathcal{A}$ , entonces  $\Omega \setminus A \in \mathcal{A}$ .
- $\mathcal{A}$  es cerrado bajo uniones numerables: si  $A_i \in \mathcal{A}$  para todo  $i \in \mathbb{N}$  y  $B = \bigcup_{i \in \mathbb{N}} A_i$ , entonces  $B \in \mathcal{A}$ .

De las propiedades anteriores deducimos que  $\Omega \in \mathcal{A}$  y que  $\mathcal{A}$  también es cerrado bajo intersecciones numerables.

**Definición 2** (Medida de probabilidad). Una *medida de probabilidad*  $P$  sobre un espacio de medida  $(\Omega, \mathcal{A})$  es una función  $P: \mathcal{A} \rightarrow [0, 1]$  que verifica:

- $P(\Omega) = 1$ .
- Si  $A \subset \Omega$ , entonces  $P(A) \geq 0$ .
- $P$  es  $\sigma$ -aditiva, es decir: dada  $\{A_i\}_{i \in \mathbb{N}}$  una sucesión de conjuntos disjuntos dos a dos en  $\mathcal{A}$ , entonces

$$P\left(\bigcup_{i \in \mathbb{N}} A_i\right) = \sum_{i \in \mathbb{N}} P(A_i).$$

Denotaremos por *suceso seguro* al suceso que siempre va a ocurrir. A partir de la primera condición sabemos que el suceso seguro tiene la máxima probabilidad posible. La segunda condición garantiza la no negatividad de la probabilidad. Por último, la tercera condición implica que dado un conjunto de sucesos disjuntos dos a dos, la probabilidad de que ocurra cualquiera de ellos es igual a la suma de las probabilidades de cada uno.

**Proposición 1.** Toda medida de probabilidad,  $P$ , cumple:

- $P(\emptyset) = 0$ .
- Dados  $A, B \in \mathcal{A}$ , entonces  $P(A \cup B) = P(A) + P(B) - P(A \cap B)$ .

**Definición 3** (Espacio de probabilidad). Definimos como *espacio de medida* a la tupla  $(\Omega, \mathcal{A}, \mu)$  donde  $\mu: \mathcal{A} \rightarrow \mathbb{R}_0^+$  es una medida en  $(\Omega, \mathcal{A})$  y *espacio de probabilidad* a la tupla formada por  $(\Omega, \mathcal{A}, P)$  donde  $P$  una medida de probabilidad en  $(\Omega, \mathcal{A})$ .

## 2.2 VARIABLES ALEATORIAS

Una *variable aleatoria* es una función que asigna un valor, normalmente numérico, al resultado de un experimento aleatorio. Dada una variable aleatoria no es posible saber su valor exacto al ser medida, aunque sí conocemos una distribución de probabilidad para describir la probabilidad de que se den los diferentes valores. A continuación, formalizaremos el concepto de variable aleatoria.

**Definición 4** (Función medible). Sean  $(\Omega_1, \mathcal{A})$  y  $(\Omega_2, \mathcal{S})$  dos espacios de medida. Una *función medible*  $n$ -dimensional es una función  $X: \Omega_1 \rightarrow \Omega_2$  que verifica:

$$X^{-1}(S) = \{\omega \in \Omega : X(\omega) \in S\} \in \mathcal{A}, \text{ para todo } S \in \mathcal{S}.$$

**Definición 5** (Variable aleatoria). Sea  $(\Omega_1, \mathcal{A}, P)$  un espacio de probabilidad y  $(\Omega_2, \mathcal{S})$  un espacio de medida. Una *variable aleatoria*  $\mathbf{X} = (X_1, \dots, X_n)$  es una función medible  $\mathbf{X}: \Omega_1 \rightarrow \Omega_2$  del espacio de probabilidad al espacio de medida.

Diremos que la variable aleatoria es *unidimensional* si  $n = 1$  y *multivariante* cuando  $n > 1$ . Cuando tengamos una variable aleatoria multivariante  $\mathbf{X} = (X_1, \dots, X_n)$ , llamaremos a  $\mathbf{X}$  variable aleatoria conjunta o *vector aleatorio* y a cada  $X_i$  con  $i = 1, \dots, n$  variable aleatoria marginal.

**Definición 6** (Probabilidad inducida). Sea  $(\Omega_1, \mathcal{A}, P)$  un espacio de probabilidad y  $(\Omega_2, \mathcal{S})$  un espacio de medida. La *probabilidad inducida* por una variable aleatoria  $\mathbf{X}$  viene dada por la función:

$$P_{\mathbf{X}}(S) = P(\mathbf{X}^{-1}(S)), \text{ para todo } S \in \mathcal{S}.$$

*Ejemplo 1.* Consideramos el lanzamiento de una moneda. Los posibles resultados del experimento serán cara o cruz, los cuales serán nuestros sucesos aleatorios. Definiremos nuestra variable aleatoria como:

$$X = \begin{cases} 0, & \text{si sale cara,} \\ 1, & \text{si sale cruz.} \end{cases}$$

**Definición 7** (Función de distribución). La *función de distribución* acumulada de una variable aleatoria  $X$  es una función  $F: \mathbb{R} \rightarrow [0, 1]$  definida como:

$$F(x) = P(X \leq x).$$

**Proposición 2.** La función de distribución acumulada  $F$  asociada a la variable aleatoria  $X$  satisface las siguientes propiedades:

- $\lim_{x \rightarrow +\infty} F(x) = 1.$
- $\lim_{x \rightarrow -\infty} F(x) = 0.$
- Es creciente, es decir, si  $x_1 \leq x_2$ , entonces  $F(x_1) \leq F(x_2).$
- Es continua por la derecha, es decir,  $\lim_{x \rightarrow a^+} F(x) = F(a^+).$

Si la imagen de la variable aleatoria  $X$  es numerable, diremos que la variable aleatoria es *discreta* y viene descrita por la función de probabilidad  $p$  que devuelve la probabilidad de  $X$  de ser igual a cierto valor  $x$ .

Si la imagen de la variable aleatoria  $X$  es infinita no numerable, diremos que la variable aleatoria es *continua* y viene descrita por la función de densidad  $f$  que caracteriza la posibilidad relativa de que  $X$  tome un valor cercano a  $x$ .

**Definición 8** (Función de probabilidad). Sea  $X$  una variable aleatoria discreta con posibles valores  $\{x_1, \dots, x_n\}$  su *función de probabilidad* se define como

$$f(x) = \begin{cases} P(X = x), & \text{si } x \in \{x_1, \dots, x_n\}, \\ 0, & \text{en otro caso.} \end{cases}$$

**Definición 9** (Función de densidad). Sea  $X$  una variable aleatoria continua se dice que la función integrable no-negativa  $f$  es su *función de densidad* si para todo  $x \in \mathbb{R}$ ,

$$P(X \leq x) = \int_{-\infty}^x f(u) du.$$

### 2.2.1 Distribución conjunta

Una *distribución conjunta* es la distribución de probabilidad de la intersección de las realizaciones de dos o más variables aleatorias cualesquiera. A continuación, definiremos algunos conceptos que ya discutimos para una única variable aleatoria para el caso multivariante.

**Definición 10** (Función de distribución conjunta). La *función de distribución conjunta* de un vector aleatorio  $\mathbf{X} = (X_1, \dots, X_n)$  es una función  $F_{\mathbf{X}}: \mathbb{R}^n \rightarrow [0, 1]$  definida como:

$$F_{\mathbf{X}}(x_1, \dots, x_n) = P(X_1 \leq x_1, \dots, X_n \leq x_n).$$

**Definición 11** (Función de probabilidad conjunta). Sea un vector aleatorio discreto  $\mathbf{X} = (X_1, \dots, X_n)$  con posibles valores en el conjunto producto

$$\mathcal{X} = \left\{ \{x_{11}, \dots, x_{1n}\} \times \dots \times \{x_{n1}, \dots, x_{nn}\} \right\}$$

su *función de probabilidad conjunta* se define como

$$f_{\mathbf{X}}(x_1, \dots, x_n) = \begin{cases} P(X_1 = x_1, \dots, X_n = x_n), & \text{si } (x_1, \dots, x_n) \in \mathcal{X}, \\ 0, & \text{en otro caso.} \end{cases}$$

**Definición 12** (Función de densidad conjunta). Sea  $\mathbf{X} = (X_1, \dots, X_n)$  un vector aleatorio continuo se dice que la función integrable no-negativa  $f_{\mathbf{X}}$  es su *función de densidad conjunta* si para todo  $(x_1, \dots, x_n) \in \mathbb{R}^n$ ,

$$P(X_1 \leq x_1, \dots, X_n \leq x_n) = \int_{-\infty}^{x_1} \dots \int_{-\infty}^{x_n} f_{\mathbf{X}}(u_1, \dots, u_n) du_1 \dots du_n.$$

### 2.2.2 Distribución marginal

La *distribución marginal* de un subconjunto de un vector aleatorio es la distribución de probabilidad de las variables contenidas en el subconjunto. Procederemos a definir algunos conceptos que serán útiles cuando necesitemos calcular la distribución para alguna componente de  $\mathbf{X}$ .

**Definición 13** (Función de distribución marginal). Sea  $\mathbf{X} = (X_1, \dots, X_n)$  un vector aleatorio, la *función de distribución marginal* de  $X_1$  se define como:

$$F_{X_1}(x_1) = \lim_{x_2 \rightarrow +\infty} \dots \lim_{x_n \rightarrow +\infty} F_{\mathbf{X}}(x_1, \dots, x_n).$$

**Definición 14** (Función de probabilidad marginal). Sea un vector aleatorio discreto  $\mathbf{X} = (X_1, \dots, X_n)$ , la *función de probabilidad marginal* de  $X_1$  se define como:

$$f_{X_1}(x_1) = \sum_{x_2} \cdots \sum_{x_n} f_{\mathbf{X}}(x_1, \dots, x_n),$$

y la *función de probabilidad marginal* del subconjunto  $(X_1, X_2)$  viene dada por:

$$f_{X_1, X_2}(x_1, x_2) = \sum_{x_3} \cdots \sum_{x_n} f_{\mathbf{X}}(x_1, \dots, x_n).$$

**Definición 15** (Función de densidad marginal). Sea  $\mathbf{X} = (X_1, \dots, X_n)$  un vector aleatorio continuo, la *función de densidad marginal* de  $X_1$  se define como:

$$f_{X_1}(x_1) = \int_{x_2} \cdots \int_{x_n} f_{\mathbf{X}}(x_1, \dots, x_n) dx_2 \cdots dx_n,$$

y la *función de densidad marginal* del subconjunto  $(X_1, X_2)$  viene dada por:

$$f_{X_1, X_2}(x_1, x_2) = \int_{x_3} \cdots \int_{x_n} f_{\mathbf{X}}(x_1, \dots, x_n) dx_3 \cdots dx_n.$$

## 2.3 PROBABILIDAD CONDICIONAL

Llamamos *probabilidad condicional* a la probabilidad de que ocurra un suceso  $A$ , sabiendo que también sucede otro suceso  $B$ . Algunos resultados relacionados y que usaremos en el trabajo son el Teorema de probabilidad total y el Teorema de Bayes.

**Definición 16** (Probabilidad condicional). Para cualesquiera dos sucesos  $A, B \in \mathcal{A}$  tales que  $P(B) > 0$ . Definimos la *probabilidad condicional* de  $A$  sobre  $B$  como:

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

**Teorema 1** (Teorema de la probabilidad total). Sea  $\{A_i : i = 1, \dots, n\}$  una partición sobre  $\Omega$  y sea  $B \in \mathcal{A}$  un suceso arbitrario del que se conocen las probabilidades condicionales  $P(B | A_i)$ , entonces la probabilidad del suceso  $B$  viene dada por

$$P(B) = \sum_{i=1}^n P(A_i)P(B | A_i).$$

*Demostración.* Utilizando que  $\{A_i \in \mathcal{A} : i = 1, \dots, n\}$  es una partición del espacio muestral  $\Omega$  y por tanto cumple:

$$\blacksquare \quad \Omega = \bigcup_{i=1}^n A_i.$$

- $A_i \cap A_j = \emptyset$ , para todo  $A_i \neq A_j$ .

Podemos escribir el suceso  $B$  como

$$B = \bigcup_{i=1}^n B \cap A_i.$$

Como los conjuntos  $A_i$  son disjuntos dos a dos, entonces los conjuntos  $B \cap A_i$  también lo son. En consecuencia

$$P(B) = P(B \cap A_1) + \cdots + P(B \cap A_n).$$

Por último, como  $B, A_i \in \mathcal{A}$  usando la Definición 16, tenemos

$$P(B) = P(B | A_1)P(A_1) + \cdots + P(B | A_n)P(A_n) = \sum_{i=1}^n P(B | A_i)P(A_i).$$

□

**Teorema 2** (Teorema de Bayes). Para cualesquiera dos sucesos  $A, B \in \mathcal{A}$  tales que  $P(B) > 0$ , se tiene

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

*Demostración.* Utilizando la Definición 16, tenemos que  $P(B | A)P(A) = P(B \cap A)$  y sabiendo que  $P(B \cap A) = P(A \cap B)$  concluimos:

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

□

**Notación 1.** A partir de este punto, denotaremos  $P(A \cap B)$  como  $P(A, B)$ .

### 2.3.1 Distribución condicional

Una *distribución condicional* se define como la distribución de una de las variables condicionada a cada valor de una o más variables aleatorias. Definiremos el concepto de función de probabilidad y función de densidad condicional, que nos serán útiles en el futuro.

**Definición 17** (Función de probabilidad condicional). Sean  $X, Y$  variables aleatorias discretas con función de probabilidad conjunta  $f_{X,Y}$  y sea  $f_Y(y)$  la función de probabilidad marginal de  $Y$ . Llamaremos *función de probabilidad condicional* de  $X$  dado  $Y = y$ , a la función definida como:

$$f_{X|Y}(x | y) = \frac{f_{X,Y}(x, y)}{f_Y(y)}, \text{ con } f_Y(y) \neq 0.$$

**Comentario 1.** La *función de densidad condicional* se define de la misma forma para  $X, Y$  variables aleatorias continuas, tomando  $f_{X,Y}$  como la función de densidad conjunta y  $f_Y$  como la función de densidad marginal.

## 2.4 INDEPENDENCIA

Dos variables aleatorias son *independientes* entre sí cuando la probabilidad de cada variable no está influida por la ocurrencia de la otra. A continuación formalizaremos esta idea.

**Definición 18** (Variables aleatorias independientes). Sean  $X, Y$  variables aleatorias, diremos que son *independientes* si, y solo si,

$$P(X = x, Y = y) = P(X = x)P(Y = y).$$

Cuando esto ocurra, lo denotaremos como  $X \perp Y$ .

**Definición 19** (Variables aleatorias independientes condicionalmente). Sean  $X, Y, Z$  variables aleatorias, entonces  $X$  e  $Y$  son *condicionalmente independientes* dado  $Z$  si, y solo si,

$$P(X = x, Y = y \mid Z = z) = P(X = x \mid Z = z)P(Y = y \mid Z = z),$$

o equivalentemente,

$$P(X = x \mid Y = y, Z = z) = P(X = x \mid Z = z).$$

En este caso, lo denotaremos como  $X \perp Y \mid Z$ .

**Definición 20** (Variables aleatorias independientes e idénticamente distribuidas). Sea un vector aleatorio  $\mathbf{X} = (X_1, \dots, X_n)$  diremos que sus componentes son *independientes e idénticamente distribuidas* (i.i.d) si, y solo si, son *independientes*:

$$F_{\mathbf{X}}(x) = F_{X_1}(x_1) \cdots F_{X_n}(x_n), \text{ para todo } x_1, \dots, x_n \in \mathbb{R}$$

y son *idénticamente distribuidas*:

$$F_{X_1}(x) = F_{X_i}(x), \text{ para todo } i \in \{2, \dots, n\}, \text{ para todo } x \in \mathbb{R}.$$

---

DISTRIBUCIONES DE PROBABILIDAD

---

En este capítulo recordaremos algunos ejemplos de distribuciones de probabilidad que usaremos durante el desarrollo del trabajo y definiremos algunos conceptos matemáticos estadísticos previos a la construcción de las distribuciones. Utilizaremos como referencias bibliográficas a [Cramer \[2004\]](#) y [Dembo \[2014\]](#).

### 3.1 ESPERANZA DE UNA VARIABLE ALEATORIA

Estamos listos para introducir el concepto de *esperanza matemática* de una variable aleatoria  $X$ . La esperanza representa el valor promedio de los valores que toma la variable.

**Definición 21** (Esperanza de una variable aleatoria). Sea  $X$  una variable aleatoria unidimensional no negativa en un espacio de probabilidad  $(\Omega, \mathcal{A}, P)$ . Definimos su *esperanza* como:

$$\mathbb{E}[X] = \int_{\Omega} X(\omega) dP(\omega).$$

Se denota por  $\mu_X$  o  $\mu$  dependiendo de si queremos destacar o no cual es la variable aleatoria a la que se refiere.

Si  $X$  es una variable aleatoria discreta y toma valores en el conjunto  $\mathcal{X}$ , entonces su esperanza se define como:

$$\mathbb{E}[X] = \sum_{x \in \mathcal{X}} x P_X(x).$$

donde  $x$  es cada posible resultado del experimento y  $P_X(x)$  la probabilidad inducida por  $X$  de obtener el resultado  $x$ .

Si  $X$  es continua y  $f(x)$  es su función de densidad, entonces su esperanza se define como:

$$\mathbb{E}[X] = \int_{-\infty}^{+\infty} x f(x) dx.$$



**Proposición 3.** Dadas  $X, Y$  variables aleatorias y  $a, b \in \mathbb{R}$ .  $\mathbb{E}[X]$  es un operador lineal, es decir:

$$\mathbb{E}[aX + bY] = a\mathbb{E}[X] + b\mathbb{E}[Y].$$

*Demostración.* Consecuencia de la linealidad de la integral de Lebesgue.  $\square$

**Proposición 4.** Sean  $X$  una variable aleatoria con función de densidad  $f_X$  y  $g: \mathbb{R} \rightarrow \mathbb{R}$  una función integrable de Lebesgue, entonces se cumplen las siguientes propiedades:

- $g(X)$  es una variable aleatoria y su esperanza viene dada por:

$$\mathbb{E}[g(X)] = \int_{\Omega} g(x) f_X(x) dx.$$

- Si  $X$  es discreta,  $\mathbb{E}[g(X)] = \sum_{x \in \mathcal{X}} g(x) P_X(x)$ .

- Si  $X$  es continua,  $\mathbb{E}[g(X)] = \int_{-\infty}^{+\infty} g(x) f(x) dx$ .

- Sean  $X_1, \dots, X_n$  variables aleatorias independientes, entonces

$$\mathbb{E} \left[ \prod_{i=1}^n X_i \right] = \prod_{i=1}^n \mathbb{E}[X_i].$$

A continuación, definimos la *esperanza condicional* de una variable aleatoria como el valor esperado de dicha variable respecto a una distribución de probabilidad condicional.

**Definición 22** (Esperanza condicional). Sean  $X, Y$  variables aleatorias discretas y  $f_{X|Y}$  su función de probabilidad condicional, definimos la *esperanza condicional* de  $X$  dado  $Y = y$  como:

$$\mathbb{E}[X | Y = y] = \sum_{x \in \mathcal{X}} x f_{X|Y}(x | y),$$

en el caso continuo, sea  $f_{X|Y}$  la función de densidad condicional, la esperanza condicional se calcula como:

$$\mathbb{E}[X | Y = y] = \int_{-\infty}^{+\infty} x f_{X|Y}(x | y) dy.$$

## 3.2 MOMENTOS DE UNA VARIABLE ALEATORIA

A partir de la definición de esperanza de una variable aleatoria, podemos construir el concepto de *momentos* de una variable aleatoria, que nos permiten extraer información relevante de la distribución desconocida.

**Definición 23** (Momento no centrado de una variable aleatoria). Sea  $X$  una variable aleatoria y  $k \in \mathbb{N}$ . Definimos el *momento no centrado de orden  $k$*  como:

$$\mu_k = \mathbb{E}[X^k],$$

siempre que exista dicha esperanza.

*Comentario 2.* En el caso  $k = 1$ , obtenemos la esperanza matemática de la variable aleatoria  $X$ , a la cual también denominamos *media*  $\mu_X$  o simplemente  $\mu$ .

**Definición 24** (Momento centrado de una variable aleatoria). Sea  $X$  una variable aleatoria,  $k \in \mathbb{N}$  y  $c \in \mathbb{R}$  definimos el *momento centrado en  $c$  de orden  $k$*  como:

$$\mu_k = \mathbb{E}[(X - c)^k].$$

**Definición 25** (Varianza). La *varianza* de una variable aleatoria se define como:

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

A la raíz cuadrada positiva de la varianza la denotaremos como *desviación estándar*  $\sigma_X$  o simplemente  $\sigma$ .

**Proposición 5.** Sean  $X$  una variable aleatoria y  $a, b \in \mathbb{R}$ , entonces se cumplen las siguientes propiedades:

- $\text{Var}(X) = \mathbb{E}[X]^2 - \mathbb{E}[X^2]$ .
- $\text{Var}(b) = 0$ .
- $\text{Var}(aX) = a^2 \text{Var}(X)$ .
- Sea  $X, Y$  variables aleatorias independientes, entonces:

$$\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y).$$

**Proposición 6.** Dadas  $X, Y$  variables aleatorias y  $a, b \in \mathbb{R}$ .  $\text{Var}(X)$  es una operación lineal, es decir:

$$\text{Var}(aX + b) = \mathbb{E}[(aX + b) - \mathbb{E}[aX + b]]^2 = a^2 \mathbb{E}[(X - \mu_X)^2] = a^2 \text{Var}(X).$$

*Demostración.* Consecuencia de la linealidad de la esperanza. □

En el mundo real, cuando aplicamos estos conceptos lo haremos sobre múltiples características observables, es decir, sobre vectores aleatorios como escribiremos a continuación.

**Definición 26** (Esperanza de un vector aleatorio). Sea  $\mathbf{X} = (X_1, \dots, X_n)^T$  un vector aleatorio. Se define la *esperanza* de  $\mu_{\mathbf{X}}$  como:

$$\mu_{\mathbf{X}} = \mathbb{E}[\mathbf{X}] = \begin{bmatrix} \mathbb{E}[X_1] \\ \vdots \\ \mathbb{E}[X_n] \end{bmatrix},$$

siempre que existan las esperanzas unidimensionales.

Para generalizar la varianza de una variable aleatoria, construiremos una matriz multidimensional con la que aparece el concepto de covarianza.

**Definición 27** (Matriz de covarianzas). Sea  $\mathbf{X} = (X_1, \dots, X_n)^T$  un vector aleatorio. Se define, la *matriz de covarianzas* de  $\mathbf{X}$  como:

$$\Sigma_{\mathbf{X}} = \text{Cov}(\mathbf{X}) = \mathbb{E}[(\mathbf{X} - \mu_{\mathbf{X}})(\mathbf{X} - \mu_{\mathbf{X}})^T] = \begin{bmatrix} \sigma_{11} & \dots & \sigma_{1n} \\ \vdots & \ddots & \vdots \\ \sigma_{n1} & \dots & \sigma_{nn} \end{bmatrix},$$

donde  $\sigma_{ij} = \text{Cov}(X_i, X_j) = \mathbb{E}[(X_i - \mu_i)(X_j - \mu_j)] = \sigma_{ji}$  es la covarianza de las variables aleatorias  $X_i, X_j$ . Podremos definirla cuando existan todas las covarianzas.

### 3.3 EJEMPLOS DE DISTRIBUCIONES

La *distribución de probabilidad* de una variable aleatoria es una función que hace corresponder a cada suceso definido sobre la variable, la probabilidad de que dicho suceso ocurra. Describiremos algunas de estas funciones que nos serán de interés a lo largo del desarrollo del trabajo.

**Definición 28** (Moda). La *moda* de una distribución es el valor donde su función de probabilidad alcanza su máximo. Es el valor que aparece con mayor frecuencia en un conjunto de datos.

Las distribuciones pueden ser *unimodales*, *bimodales* o *multimodales* dependiendo de si tienen un solo valor de moda, dos o más, respectivamente.

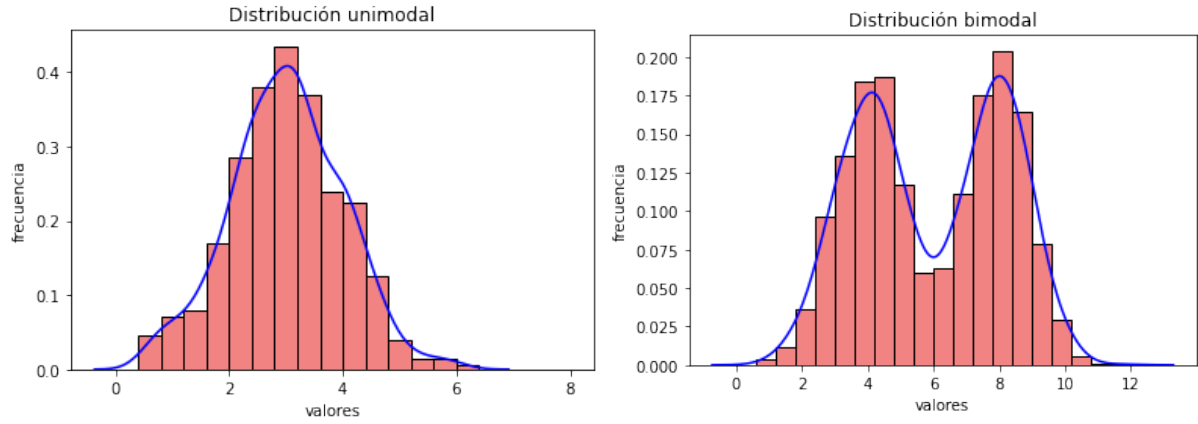


Figura 1: Ejemplos de distribuciones de probabilidad.

### 3.3.1 Distribución Bernoulli

La *distribución Bernoulli* es una distribución de probabilidad aplicada a una variable aleatoria discreta, la cual solo puede tomar dos resultados mutuamente excluyentes (éxito o fracaso).

**Definición 29** (Distribución Bernoulli). Una variable aleatoria unidimensional  $X$  sigue una *distribución Bernoulli* de parámetro  $p$  si su función de probabilidad viene dada por:

$$P(X = x) = \begin{cases} p, & \text{si } x = 1, \\ 1 - p, & \text{si } x = 0. \end{cases}$$

Escribiremos la distribución como  $X \sim \text{Bernoulli}(p)$ , donde el parámetro  $p \in (0, 1)$  indica la probabilidad de éxito y  $(1 - p)$  la probabilidad de fracaso del experimento.

**Proposición 7** (Propiedades). Si  $X \sim \text{Bernoulli}(p)$ , entonces la variable aleatoria  $X$  satisface las siguientes propiedades:

- $\mathbb{E}[X] = p$ .
- $\text{Var}(X) = p(1 - p)$ .

### 3.3.2 Distribución de Poisson

La *distribución de Poisson* es una distribución de probabilidad discreta que modela el número de sucesos raros que ocurren en un determinado periodo de tiempo.

**Definición 30** (Distribución de Poisson). Una variable aleatoria unidimensional  $X$  sigue una *distribución de Poisson* de parámetro  $\lambda$  si su función de probabilidad viene dada por:

$$P(X = x) = \frac{\lambda^x e^{-\lambda}}{x!},$$

donde  $x = 0, 1, \dots$  es el número de ocurrencias del evento o fenómeno.

Escribiremos la distribución como  $X \sim \text{Poisson}(\lambda)$ , donde el parámetro  $\lambda > 0$  representa el número de veces que se espera que ocurra el fenómeno durante un intervalo dado.

**Proposición 8.** Si  $X \sim \text{Poisson}(\lambda)$ , entonces la variable aleatoria  $X$  satisface las siguientes propiedades:

1.  $\mathbb{E}[X] = \lambda$ .
2.  $\text{Var}(X) = \lambda$ .
3. La moda de  $X$  es  $\lfloor \lambda \rfloor$  (el mayor entero menor que  $\lambda$ ).

*Demostración.* Demostraremos la propiedad 1 utilizando la definición de  $\mathbb{E}[X]$ .

$$\begin{aligned} \mathbb{E}[X] &= \sum_{x=0}^{\infty} x \left( \frac{\lambda^x e^{-\lambda}}{x!} \right) \\ &= \lambda e^{-\lambda} \sum_{x=1}^{\infty} \frac{\lambda^{x-1}}{(x-1)!} \\ &\stackrel{(*)}{=} \lambda e^{-\lambda} \sum_{y=0}^{\infty} \frac{\lambda^y}{y!} \\ &= \lambda e^{-\lambda} e^{\lambda} \\ &= \lambda. \end{aligned}$$

donde en  $(*)$  hemos cambiado  $(x-1)$  por  $y$ .

Para la propiedad 2, sabiendo que

$$\text{Var}(X) = \mathbb{E}[X^2] - \mathbb{E}[X]^2 = \mathbb{E}[X(X-1)] + \mathbb{E}[X] - \mathbb{E}[X]^2. \quad (1)$$

Procederemos con el cálculo de  $\mathbb{E}[X(X-1)]$ .

$$\begin{aligned}
 \mathbb{E}[X(X-1)] &= \sum_{x=0}^{\infty} x(x-1) \left( \frac{\lambda^x e^{-\lambda}}{x!} \right) \\
 &= \lambda^2 e^{-\lambda} \sum_{x=2}^{\infty} \frac{\lambda^{x-2}}{(x-2)!} \\
 &\stackrel{(*)}{=} \lambda e^{-\lambda} \sum_{y=0}^{\infty} \frac{\lambda^y}{y!} \\
 &= \lambda^2 e^{-\lambda} e^{\lambda} \\
 &= \lambda^2.
 \end{aligned}$$

donde en  $(*)$  hemos cambiado  $(x-2)$  por  $y$ .

Sustituyendo en la Ecuación (1) concluimos que,  $\text{Var}(X) = \lambda^2 + \lambda - \lambda^2 = \lambda$ . □

**Proposición 9.** Si  $X \sim \text{Poisson}(\lambda)$ , donde  $X$  es una variable aleatoria que representa el número de sucesos raros en una unidad de tiempo e  $Y$  es una variable aleatoria que representa el número de dichos sucesos raros en un tiempo  $t$ , se tiene que

$$Y \sim \text{Poisson}(t\lambda).$$

*Ejemplo 2.* En un instituto, el número medio de suspensos por clase es de 2,4. Es decir, si  $X$  es el número de suspensos por clase, entonces

$$X \sim \text{Poisson}(2,4).$$

¿Cuál es la probabilidad de que en una clase no haya suspensos?

$$P(X=0) = \frac{2,4^0 e^{-2,4}}{0!} = e^{-2,4} = 0,09.$$

¿Cuál es la probabilidad de que en 3 clases haya exactamente 6 suspensos?

Sea  $Y$  el número de suspensos en 3 clases. Sabemos que:

$$Y \sim \text{Poisson}(2,4 \cdot 3) = \text{Poisson}(7,2)$$

$$P(Y=6) = \frac{7,2^6 e^{-7,2}}{6!} = e^{-3,4} = 0,14.$$

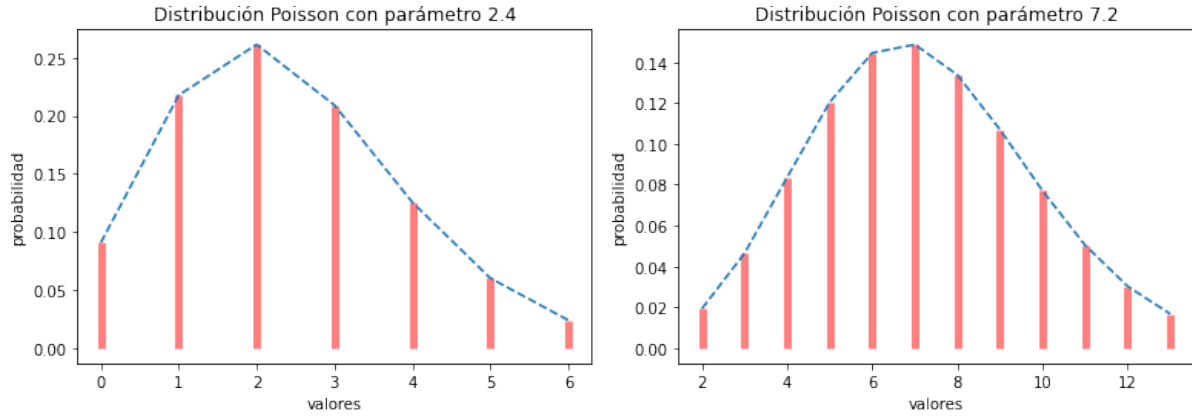


Figura 2: Distribuciones calculadas en el Ejemplo 2.

### 3.3.3 Distribución Normal

La *distribución normal* o *distribución de Gauss* se utiliza para representar variables aleatorias de valor real cuyas distribuciones son desconocidas.

**Definición 31** (Distribución normal). Una variable aleatoria unidimensional  $X$  sigue una *distribución normal* o *gaussiana* de parámetros  $\mu, \sigma$ , si su función de densidad  $f: \mathbb{R} \rightarrow \mathbb{R}$  viene dada por:

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2},$$

donde  $\mu, \sigma \in \mathbb{R}$ .

Donde escribiremos la distribución como  $X \sim \mathcal{N}(\mu, \sigma^2)$ , donde recordemos el parámetro  $\mu$  se refiere a la media y  $\sigma$  a la desviación estándar de la variable aleatoria.

**Proposición 10.** Si  $X \sim \mathcal{N}(\mu, \sigma^2)$ , entonces la variable aleatoria  $X$  satisface las siguientes propiedades :

- La distribución es simétrica respecto a  $\mu$ .
- $P(\mu - \sigma < X < \mu + \sigma) \approx 0,683$ .
- $P(\mu - 2\sigma < X < \mu + 2\sigma) \approx 0,955$ .
- $P(\mu - 3\sigma < X < \mu + 3\sigma) \approx 0,997$ .
- La moda de  $X$  coincide con  $\mu$ .

Como consecuencia de la primera propiedad de la Proposición 10, podemos relacionar todas las variables aleatorias normales con la distribución  $\mathcal{N}(0, 1)$ .

**Proposición 11** (Estandarización de variables aleatorias normales). Sea  $X$  una variable aleatoria tal que  $X \sim \mathcal{N}(\mu, \sigma^2)$ , entonces:

$$Z = \frac{X - \mu}{\sigma} \sim \mathcal{N}(0, 1),$$

donde  $\mathcal{N}(0, 1)$  es la distribución normal estándar.

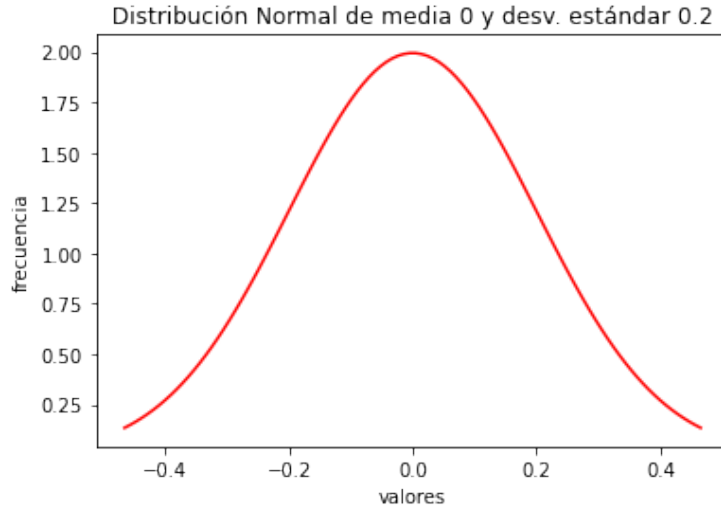


Figura 3: Ejemplo de una distribución normal.

**Definición 32** (Media muestral). Sean  $X_1, \dots, X_n$  variables aleatorias obtenidas a partir de  $X$  y que siguen su misma distribución, se define la *media muestral* como:

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i.$$

La importancia de la distribución normal reside principalmente en el Teorema central del límite que trata sobre la distribución de la media muestral para variables aleatorias independientes e idénticamente distribuidas (i.i.d) y garantiza una distribución aproximadamente normal cuando  $n$  es lo suficientemente grande.

**Teorema 3** (Teorema central del límite). Sean  $X_1, X_2, \dots, X_n$  variables aleatorias i.i.d con media  $\mu$  y desviación estándar  $\sigma^2$  (ambas finitas). Con  $n$  suficientemente grande, se tiene que

$$\frac{\bar{X} - \mu}{\sigma / \sqrt{n}} \sim \mathcal{N}(0, 1).$$



Finalmente, podemos relacionar las dos distribuciones estudiadas aproximando la distribución de Poisson mediante la distribución normal estándar a partir de la siguiente relación:

**Proposición 12.** Sea  $X$  una variable aleatoria tal que  $X \sim \text{Poisson}(\lambda)$  con  $\lambda$  suficientemente grande, entonces:

$$\frac{X - \lambda}{\sqrt{\lambda}} \sim \mathcal{N}(0, 1).$$

Como una extensión del caso unidimensional, aparece el caso de distribución normal multivariante que definiremos a continuación.

**Definición 33** (Distribución normal multivariante). Sea un vector aleatorio  $\mathbf{X} = (X_1, \dots, X_n)^T$  diremos que sigue una *distribución normal multivariante* de parámetros  $\boldsymbol{\mu}, \boldsymbol{\Sigma}$  si su función de densidad  $f: \mathbb{R}^N \rightarrow \mathbb{R}$  viene dada por:

$$f(x) = \frac{1}{\sqrt{\det(2\pi\boldsymbol{\Sigma})}} e^{-\frac{1}{2}(x-\boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(x-\boldsymbol{\mu})},$$

donde  $\boldsymbol{\mu} \in \mathbb{R}^N$  y  $\boldsymbol{\Sigma} \in \mathcal{M}(\mathbb{R})$ .

Escribiremos la distribución como  $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , donde recordemos el parámetro  $\boldsymbol{\mu}$  se refiere al vector de las medias de la distribución y  $\boldsymbol{\Sigma}$  a la matriz de covarianzas.

#### 3.3.4 Distribución Uniforme

**Definición 34** (Distribución uniforme discreta). Una variable aleatoria discreta  $X$  con posibles valores  $\{x_1, \dots, x_n\}$  diremos que sigue una *distribución uniforme* si:

$$P(X = x_i) = \frac{1}{n}, \text{ para todo } i = 1, \dots, n.$$

**Definición 35** (Distribución uniforme continua). Una variable aleatoria continua  $X$  sigue una *distribución uniforme* en el intervalo  $(a, b)$  si su función de densidad viene dada por:

$$f(x) = \begin{cases} \frac{1}{b-a}, & \text{si } x \in (a, b), \\ 0, & \text{si } x \notin (a, b). \end{cases}$$

Escribiremos la distribución como  $X \sim U(a, b)$ , donde la variable aleatoria queda definida por los extremos del intervalo, es decir,  $a$  y  $b$  son sus parámetros.

**Proposición 13** (Propiedades). Si  $X \sim U(a, b)$ , entonces la variable aleatoria  $X$  satisface las siguientes propiedades:

- $\mathbb{E}[X] = \frac{a+b}{2}$ .
- $\text{Var}(X) = \frac{(b-a)^2}{12}$ .
- La moda de  $X$  es cualquier valor en  $(a, b)$ .

### 3.3.5 Distribución Gamma

**Definición 36** (Distribución gamma). Una variable aleatoria continua  $X$  sigue una *distribución gamma* de parámetros  $\alpha$  y  $\lambda$  si su función de densidad viene dada por:

$$f(x) = \frac{\lambda^\alpha x^{\alpha-1} e^{-\lambda x}}{\Gamma(\alpha)}.$$

donde  $\Gamma$  es la *función gamma* definida como  $\Gamma(\alpha) = \int_0^{+\infty} x^{\alpha-1} e^{-x} dx$ .

Escribiremos la distribución como  $X \sim \Gamma(\alpha, \lambda)$ , donde la variable aleatoria queda definida por los parámetros  $\alpha, \lambda > 0$ .

**Proposición 14** (Propiedades). Si  $X \sim \Gamma(\alpha, \lambda)$ , entonces la variable aleatoria  $X$  satisface las siguientes propiedades:

- $\mathbb{E}[X] = \frac{\alpha}{\lambda}$ .
- $\text{Var}(X) = \frac{\alpha}{\lambda^2}$ .

**Definición 37** (Distribución gamma inversa). Una variable aleatoria continua  $X$  sigue una *distribución gamma inversa* de parámetros  $\alpha$  y  $\lambda$  si su función de densidad viene dada por:

$$f(x) = \frac{\lambda^\alpha x^{-(\alpha+1)} e^{-\frac{\lambda}{x}}}{\Gamma(\alpha)}.$$

---

## ESTADÍSTICA PARAMÉTRICA

---

La *estadística paramétrica* es una rama de la inferencia estadística que comprende los procedimientos estadísticos y de decisión basados en distribuciones conocidas que vienen determinadas por un número finito de parámetros. En este capítulo definiremos algunos conceptos de este campo de la estadística que utilizaremos a lo largo del trabajo y que es tratado en mayor profundidad en [Ibarrola and Pérez \[2012\]](#).

**Definición 38** (Muestra aleatoria simple). Sea una variable aleatoria  $X$  que sigue una distribución de probabilidad determinada. Definimos una *muestra aleatoria simple* de tamaño  $n$  como un conjunto de variables aleatorias  $(X_1, \dots, X_n)$ , independientes e idénticamente distribuidas, obtenidas a partir de la distribución de  $X$ .

**Definición 39** (Estadístico muestral). Un *estadístico de una muestra* aleatoria simple  $(X_1, \dots, X_n)$ , es una función medible  $T: \mathbb{R}^n \rightarrow \mathbb{R}$ , que se aplica a la muestra, lo denotaremos como  $T(X_1, \dots, X_n)$ .

*Comentario 3.* Un *estadístico muestral* también es una variable aleatoria.

*Ejemplo 3.* Algunos estadísticos muestrales más comunes son:

- **Media muestral** ( $\bar{X}$ ),

$$T(X_1, \dots, X_n) = \frac{1}{n} \sum_{i=1}^n X_i.$$

- **Varianza muestral** ( $S^2$ ),

$$T(X_1, \dots, X_n) = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2.$$

- **Menor valor muestral**,

$$T(X_1, \dots, X_n) = \min(X_1, \dots, X_n).$$

- **Mayor valor muestral**,

$$T(X_1, \dots, X_n) = \max(X_1, \dots, X_n).$$

Particularizando el uso del término estadístico muestral, aparece el concepto de estimador.

**Definición 40** (Estimador). Un *estimador* es un estadístico cuyos valores se utilizan para obtener información de un parámetro desconocido  $\theta$ , lo denotaremos como  $\hat{\theta}(X_1, \dots, X_n)$ .

#### 4.1 PROPIEDADES DEL ESTIMADOR

A continuación definiremos algunas propiedades que nos permitirán comparar diferentes estimadores de un mismo parámetro y nos informarán de la calidad de su estimación.

**Definición 41** (Estimador insesgado). Si  $\hat{\theta}$  es un estimador del parámetro  $\theta$ , la diferencia

$$\text{sesgo}(\hat{\theta}) = \mathbb{E}[\hat{\theta}] - \theta,$$

se denomina *sesgo* del estimador  $\hat{\theta}$  como estimador de  $\theta$ . Cuando el sesgo es nulo para cualquier valor del parámetro, es decir, si

$$\mathbb{E}[\hat{\theta}] = \theta, \text{ para todo } \theta \in \Theta,$$

diremos que el estimador  $\hat{\theta}$  es *insesgado* para  $\theta$ .

*Ejemplo 4.* Sea  $X_1, \dots, X_n$  una muestra aleatoria de una variable aleatoria  $X$ . Si el parámetro de interés es la media  $\mu = \mathbb{E}[X]$ , podemos utilizar la media muestral para estimarlo. Calculamos su sesgo:

$$\text{sesgo}(\bar{X}) = \mathbb{E}[\bar{X}] - \mu = \frac{1}{n} \sum_{i=1}^n \mathbb{E}[X_i] - \mu = \mu - \mu = 0,$$

demostrando que  $\bar{X}$  es un estimador insesgado para  $\mu$ .

**Definición 42** (Estadístico suficiente). Sea  $T = T(X_1, \dots, X_n)$  un estadístico muestral se dice *suficiente* para el parámetro  $\theta$  si la distribución condicional de la muestra  $(X_1, \dots, X_n)$  dado el valor de  $T(X_1, \dots, X_n)$  no depende del parámetro a estimar  $\theta$ .

**Definición 43** (Estimador suficiente). Un estimador  $\hat{\theta}$  es *suficiente* si es un estadístico suficiente.

A partir de la propiedad de suficiencia aparece el Teorema de factorización de Fisher-Neyman que nos ofrece una caracterización de estadístico suficiente.

**Teorema 4** (Teorema de factorización de Fisher-Neyman). Sea  $T = T(X_1, \dots, X_n)$  un estadístico muestral si la función de densidad es  $f_\theta(x)$ , entonces  $T$  es suficiente para  $\theta$  si, y solo si, podemos encontrar  $g, h: \mathbb{R}^n \rightarrow \mathbb{R}$  funciones no negativas tal que

$$f_\theta(x) = h(x)g_\theta(T(x)).$$

es decir, podemos factorizar la función de densidad como un producto de dos funciones:  $h$  que no depende de  $\theta$  y  $g$  que depende de  $\theta$  y de  $x$  solo a través de  $T(x)$ .

El tamaño de la muestra es un factor determinante en la estimación, a raíz de esto, surge el término de estimador *consistente*. En la definición de consistencia se usan nociones de *convergencia en probabilidad*, que comentaremos a continuación.

**Definición 44** (Convergencia en probabilidad). Una sucesión de variables aleatorias  $\{X_n\}$  converge en probabilidad a la variable aleatoria  $X$  si para todo  $\epsilon > 0$  entonces,

$$\lim_{n \rightarrow +\infty} P(|X_n - X| > \epsilon) = 0.$$

Lo denotaremos como  $X_n \xrightarrow{P} X$ .

**Definición 45** (Estimador consistente). Sea  $\hat{\theta}_n$  un estimador del parámetro  $\theta$  para una muestra de tamaño  $n$ . Diremos que el estimador  $\hat{\theta}_n$  es *consistente* si, cuando  $n \rightarrow +\infty$ , se verifica que

$$\hat{\theta}_n \xrightarrow{P} \theta.$$

**Definición 46** (Estimador asintóticamente insesgado). Sea  $\hat{\theta}_n$  un estimador del parámetro  $\theta$  para una muestra de tamaño  $n$ . Diremos que el estimador  $\hat{\theta}_n$  es *asintóticamente insesgado* si, cuando  $n \rightarrow +\infty$ , se verifica que

$$\text{sesgo}(\hat{\theta}_n) \rightarrow 0.$$

**Corolario 1.** Sea  $\hat{\theta}_n$  un estimador del parámetro  $\theta$  para una muestra de tamaño  $n$ . Si  $\hat{\theta}_n$  es asintóticamente insesgado y además si  $n \rightarrow +\infty$  se cumple que

$$\text{Var}(\hat{\theta}_n) \rightarrow 0.$$

entonces, el estimador  $\hat{\theta}_n$  es consistente.

Finalmente enunciaremos la propiedad de *eficiencia* de un estimador.

**Definición 47** (Estimador eficiente). Sean  $\hat{\theta}_1$  y  $\hat{\theta}_2$  dos estimadores del parámetro  $\theta$ . Se dice que  $\hat{\theta}_1$  es más *eficiente* que  $\hat{\theta}_2$  si verifica que

$$\text{Var}(\hat{\theta}_1) < \text{Var}(\hat{\theta}_2).$$

*Ejemplo 5.* Sea  $X \sim U(0, b)$ . Como  $X$  sigue una distribución uniforme, entonces:

$$\mathbb{E}[X] = \frac{b}{2} \quad \text{y} \quad \text{Var}(X) = \frac{b^2}{12}.$$

Sean dos estimadores insesgados:

$$\hat{b}_1(X_1, \dots, X_5) = \bar{X}.$$

$$\hat{b}_2(X_1, \dots, X_5) = \max\{X_1, \dots, X_5\}.$$

A continuación veremos cual de los dos estimadores es más eficiente, procederemos calculando la varianza de cada uno de ellos:

- Cálculo de la varianza de  $\hat{b}_1$ :

$$\text{Var}(\hat{b}_1) = \text{Var}(\bar{X}) = \frac{\text{Var}(X)}{n} = \frac{b^2/12}{5} = \frac{b^2}{60}.$$

- Cálculo de la varianza de  $\hat{b}_2$ :

$$\mathbb{E}[\hat{b}_2^2] = \int_0^b x^2 f_{\hat{b}_2}(x) dx = \int_0^b x^2 \frac{5x^4}{b^5} dx = \frac{5}{b^5} \int_0^b x^6 dx = \frac{5}{b^5} \left[ \frac{x^7}{7} \right]_0^b = \frac{5}{7} b^2.$$

Como  $\hat{b}_2$  es un estimador insesgado, entonces  $\mathbb{E}[\hat{b}_2] = b$ .

$$\text{Var}(\hat{b}_2) = \mathbb{E}[\hat{b}_2^2] - \mathbb{E}[\hat{b}_2]^2 = \frac{5}{7} b^2 - b^2 = -\frac{2}{7} b^2.$$

Concluimos que  $\hat{b}_2$  es más eficiente que  $\hat{b}_1$  ya que:

$$\text{Var}(\hat{b}_2) = -\frac{2b^2}{7} < \frac{b^2}{60} = \text{Var}(\hat{b}_1).$$

## 4.2 CRITERIOS DE EVALUACIÓN

**Definición 48** (Error cuadrático medio). El *error cuadrático medio* (ECM) de un estimador  $\hat{\theta}$  con respecto al parámetro desconocido  $\theta$  se define como:

$$\text{ECM}(\hat{\theta}) = \mathbb{E}[(\hat{\theta} - \theta)^2].$$

**Proposición 15.** Sea  $\hat{\theta}$  un estimador de  $\theta$ . Se cumple que:

$$\text{ECM}(\hat{\theta}) = \text{Var}(\hat{\theta}) + \text{sesgo}(\hat{\theta})^2.$$

*Demostración.*

$$\text{ECM}(\hat{\theta}) = \mathbb{E}[(\hat{\theta} - \theta)^2] = \mathbb{E}[\hat{\theta}^2 + \theta^2 - 2\hat{\theta}\theta] = \mathbb{E}[\hat{\theta}^2] + \theta^2 - 2\theta\mathbb{E}[\hat{\theta}].$$

Sumando y restando  $\mathbb{E}[\hat{\theta}]^2$  en la expresión anterior, obtenemos:

$$\begin{aligned} \text{ECM}(\hat{\theta}) &= \mathbb{E}[\hat{\theta}^2] + \theta^2 - 2\theta\mathbb{E}[\hat{\theta}] + \mathbb{E}[\hat{\theta}]^2 - \mathbb{E}[\hat{\theta}]^2 \\ &= (\mathbb{E}[\hat{\theta}^2] - \mathbb{E}[\hat{\theta}]^2) + (\theta^2 + \mathbb{E}[\hat{\theta}]^2 - 2\theta\mathbb{E}[\hat{\theta}]) \\ &= (\mathbb{E}[\hat{\theta}^2] - \mathbb{E}[\hat{\theta}]^2) + (\theta - \mathbb{E}[\hat{\theta}])^2 \\ &= \text{Var}(\hat{\theta}) + \text{sesgo}(\hat{\theta})^2. \end{aligned}$$

□

El ECM involucra las propiedades de insesgadez y eficiencia, ya que cuanto más cerca esté la esperanza de un estimador del parámetro y cuanto más pequeña sea su varianza, menor será su error cuadrático medio.

**Proposición 16.** Si  $\hat{\theta}$  es un estimador insesgado, entonces:

$$\text{ECM}(\hat{\theta}) = \text{Var}(\hat{\theta}).$$

A partir del concepto de ECM se define el concepto de *raíz del error cuadrático medio* que es una medida de uso frecuente de las diferencias entre los valores predichos por un estimador y los valores observados y que utilizaremos en la práctica para evaluar modelos de predicción.

**Definición 49** (Raíz del error cuadrático medio). La *raíz del error cuadrático medio* (RMSE) de un estimador  $\hat{\theta}$  con respecto al parámetro desconocido  $\theta$  se define como:

$$\text{RSME}(\hat{\theta}) = \sqrt{\text{ECM}(\hat{\theta})} = \sqrt{\mathbb{E}[(\hat{\theta} - \theta)^2]}.$$

*Comentario 4.* Si  $\hat{\theta}$  es insesgado, por la Proposición 16:

$$\text{RMSE}(\hat{\theta}) = \sqrt{\text{Var}(\hat{\theta})} = \sigma.$$

Estas medidas de evaluación pueden particularizarse para modelos de predicción a partir del siguiente resultado:

**Proposición 17.** Sea  $\hat{\mathbf{y}}$  un vector de  $n$  predicciones e  $\mathbf{y}$  el vector con las etiquetas reales de las mismas, entonces podemos calcular las medidas de de evaluación anteriores como:

$$\text{ECM} = \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2.$$

$$\text{RMSE} = \sqrt{\text{ECM}} = \sqrt{\frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2}.$$



---

## TEORÍA DE GRAFOS

---

En este capítulo incluiremos algunos conceptos fundamentales que nos serán útiles en la construcción de los modelos causales que funcionan como base del concepto de equidad contrafactual que trataremos a lo largo del trabajo. La referencia básica de este capítulo es el libro de [Godsil and Royle \[2001\]](#).

### 5.1 GRAFOS, NODOS Y ARISTAS

**Definición 50** (Par ordenado). Un *par ordenado* es una pareja de elementos, en la que los elementos vienen distinguidos por su orden. El par ordenado donde el primer elemento es  $a$  y el segundo  $b$  se denota como  $(a, b)$ . Si el orden no es relevante, lo llamaremos *par no ordenado* y lo denotaremos  $\{a, b\}$ . En este caso, es claro que  $\{a, b\} = \{b, a\}$ .

**Definición 51** (Grafo). Un *grafo*  $G = (V, E)$  es un conjunto no vacío de vértices o nodos  $V$  y aristas  $E \subset V \times V$  entre ellos. Si  $E$  consta de pares ordenados de vértices lo llamaremos *grafo dirigido*, si en otro caso  $E$  consta de pares no ordenados, lo llamaremos *grafo no dirigido*.

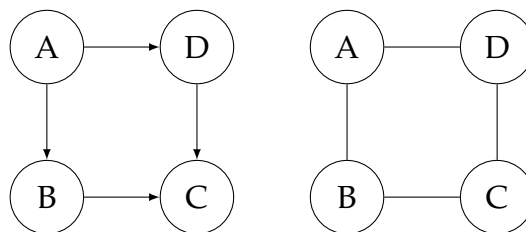


Figura 4: Ejemplo de un grafo dirigido y un grafo no dirigido, respectivamente.

**Definición 52** (Camino). En un grafo dirigido  $G = (V, E)$ , un *camino*  $A \rightarrow B$  es una secuencia de vértices  $\{A = v_0, v_1, \dots, v_{n-1}, v_n = B\}$  donde  $(v_i, v_{i+1}) \in E$  para todo  $i \in \{0, \dots, n-1\}$ . Si  $G$  es un grafo no dirigido,  $A \rightarrow B$  es un *camino* si  $\{v_i, v_{i+1}\} \in E$  para todo  $i \in \{0, \dots, n-1\}$ .

**Definición 53** (Grafo acíclico dirigido). Un *grafo acíclico dirigido* es un grafo dirigido que no tiene ciclos, es decir, que para cada nodo  $v \in V$ , no existe ningún camino que empiece y termine en  $v$ . Si un grafo no es acíclico lo llamaremos *ciclo*.

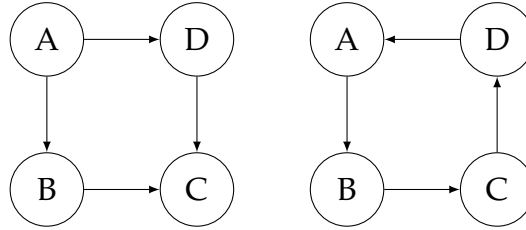


Figura 5: Ejemplo de un grafo acíclico dirigido y un ciclo dirigido, respectivamente.

*Ejemplo 6.* Podemos ver que el grafo de la izquierda de la Figura 5 es acíclico puesto que sea  $v \in \{A, B, C, D\}$  no existe ningún camino que revise  $v$ . Por otro lado, el grafo de la derecha es un ciclo, ya que el camino  $A \rightarrow B \rightarrow C \rightarrow D \rightarrow A$  es directo y empieza y termina en  $A$ .

**Definición 54** (Bucle). Sea  $G = (V, E)$  un grafo y  $v \in V$  un vértice del mismo. Diremos que  $v$  tiene un *bucle* si  $(v, v) \in E$ .

**Definición 55** (Grafo simple). Diremos que  $G$  *grafo simple* si no posee bucles ni ninguna de sus aristas relacionan al mismo par de vértices. Llamaremos a  $G$  *multigrafo* si, y solo si, no es un grafo simple.

## 5.2 ESTRUCTURA DE UN GRAFO

Ahora definiremos algunas relaciones entre nodos en un grafo acíclico dirigido.

**Definición 56** (Ancestros y descendientes de un nodo). Sean  $A, B$  dos vértices de un grafo dirigido  $G$ . Si  $A \rightarrow B$  es un camino dirigido y  $B \not\rightarrow A$  (no existe un camino dirigido de  $B$  a  $A$ ), entonces diremos que  $A$  es el *ancestro* de  $B$  y  $B$  es *descendiente* de  $A$ .

*Ejemplo 7.* En el grafo de la derecha de la Figura 5, el nodo  $A$  es un ancestro de  $B$ ,  $D$  y  $C$ . Por otro lado, el nodo  $C$  es un descendiente de  $A$ ,  $B$  y  $D$ .

**Definición 57** (Padres e hijos de un nodo). Particularizando la Definición 56, definimos a los *padres* de un nodo  $A$  como el conjunto de nodos  $pa(A)$  tal que existe una arista dirigida para cada nodo de  $pa(A)$  hacia  $A$ , la noción inversa de padres define el concepto de *hijos* a los que denotaremos como  $hi(A)$ .

**Definición 58** (Vecinos de un nodo). Sea un grafo  $G$ , definimos los *vecinos* de un nodo como todos los nodos directamente conectados a él. Denotaremos al conjunto de vecinos de un nodo  $A$  como  $ve(A)$ .

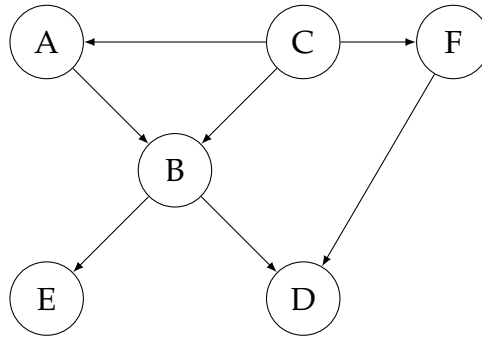


Figura 6: Ejemplo de grafo acíclico dirigido.

*Ejemplo 8.* Para comprender las definiciones anteriores, definiremos los conjuntos de padres, hijos y vecinos del nodo  $B$  para la Figura 6:

$$pa(B) = \{A, C\}, \quad hi(B) = \{E, D\}, \quad ve(B) = \{A, C, E, D\}.$$

## Parte II

### JUSTICIA EN APRENDIZAJE AUTOMÁTICO

Conceptos básicos del aprendizaje automático, formalización de las medidas de equidad y definición de los principales algoritmos de mitigación del sesgo.

---

## CONCEPTOS BÁSICOS DEL APRENDIZAJE AUTOMÁTICO

---

En este capítulo presentaremos algunos conceptos básicos del aprendizaje automático. Definiremos el concepto de aprendizaje supervisado, discutiremos el proceso de creación de modelos predictivos a partir de un conjunto de datos dado y comentaremos las medidas básicas de evaluación de un modelo de clasificación.

### 6.1 ¿QUÉ ES EL APRENDIZAJE AUTOMÁTICO?

El *aprendizaje automático* o *machine learning* se encarga de extraer patrones significativos de un conjunto de datos con el objetivo de inferir en la distribución estadística subyacente (Bishop [2006]). Para ello, durante la fase de entrenamiento, aprendemos un modelo a partir de un conjunto de datos de interés seleccionado previamente. Una vez entrenado el modelo, podremos predecir o tomar decisiones sin que el modelo haya sido específicamente programado para esa tarea. Los patrones generales obtenidos a partir del entrenamiento del modelo podrán aplicarse posteriormente a datos no vistos y seguir obteniendo resultados de utilidad.

A grandes rasgos, consideramos tres tipos de algoritmos de aprendizaje automático: los de aprendizaje supervisado que actúan sobre datos etiquetados, los de aprendizaje no supervisado que actúan sobre datos no etiquetados y los de aprendizaje por refuerzo que actúan en un entorno de ensayo-error. Nos centraremos en el entorno de *aprendizaje supervisado*, ya que nuestro trabajo utilizará algoritmos de este tipo.

#### 6.1.1 Aprendizaje supervisado

Definiremos los componentes del aprendizaje a partir de un ejemplo real: la aprobación de un crédito bancario. En principio, el banco no conoce ninguna fórmula ideal que pueda decirle cuando debe aprobar un crédito. El banco utilizará los registros de los clientes anteriores para aprender sobre ellos y encontrar una buena fórmula para la aprobación de las nuevas peticiones de créditos. Cada registro de clientes tiene información relativa al mismo, como pueden ser salario anual, años de residencia, préstamos pendientes, etc. También registra si la aprobación del crédito para ese cliente fue una buena idea, es decir, si le proporcionó o no beneficios a la entidad

financiera. Estos datos serán los que guíen la construcción de una fórmula de éxito para la aprobación del crédito que podrá utilizarse con futuros solicitantes.

A continuación, formalizaremos los principales componentes de este problema de aprendizaje. Tenemos el vector de entrada  $\mathbf{x}$  que contiene la información del cliente que se utilizará para tomar la decisión del crédito, la *función objetivo desconocida*  $f: \mathcal{X} \rightarrow \mathcal{Y}$  (fórmula ideal para la aprobación del crédito), donde  $\mathcal{X}$  es el conjunto de todas las posibles características e  $\mathcal{Y}$  el conjunto de todos los posibles resultados (en este caso, una decisión binaria si/no). Tenemos un *conjunto de datos* consistente en pares de entrada-salida  $\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$  donde cada  $y_i$  viene dado por una función desconocida  $y_i = f(\mathbf{x}_i)$  con  $i = 1, \dots, n$  (los valores de entrada corresponden a los antiguos clientes y el resultado será la decisión de crédito correcta para ellos en retrospectiva).

Finalmente, contamos con el *algoritmo de aprendizaje*  $\mathcal{A}$  que utiliza el conjunto de datos  $\mathcal{D}$  para elegir una fórmula  $g: \mathcal{X} \rightarrow \mathcal{Y}$  que mejor aproxime a la función ideal  $f$ . El algoritmo elegirá  $g$  de entre un conjunto de fórmulas candidatas consideradas, al que denominamos *conjunto de hipótesis*  $\mathcal{H}$ . Por ejemplo,  $\mathcal{H}$  podría ser el conjunto de todas las fórmulas lineales de las que el algoritmo elegirá la que mejor ajuste linealmente a los datos.

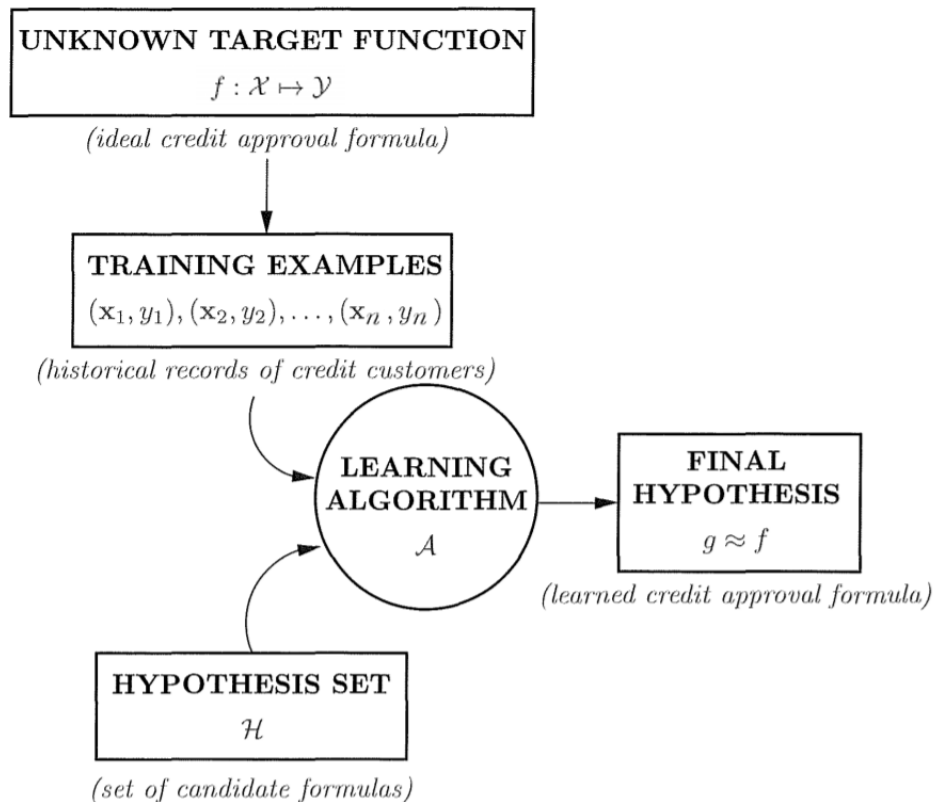


Figura 7: Configuración básica del problema de aprendizaje. (Abu-Mostafa et al. [2012])

Cuando un nuevo cliente solicite un crédito, el banco basará su decisión en  $g$  (la hipótesis que produjo el algoritmo de aprendizaje), no en  $f$  (la función objetivo ideal que sigue siendo desconocida). La decisión será buena solo en la medida en que  $g$  replique  $f$ . Para ello, el algoritmo elige  $g$  que mejor se ajuste a  $f$  en los ejemplos de entrenamiento de clientes anteriores, con la esperanza de que siga coincidiendo con  $f$  en los nuevos clientes.

## 6.2 PROPIEDADES DEL MODELO DE APRENDIZAJE

Definiremos algunos conceptos importantes que surgen junto a algunas cuestiones planteadas una vez que hemos creado nuestro modelo de aprendizaje (Barocas et al. [2019]). Suponemos un conjunto de datos etiquetado  $\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$ . Normalmente los datos  $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$  se extraen de forma independiente e idénticamente distribuida de una población  $(\mathcal{X}, \mathcal{Y})$ .

**Definición 59** (Clasificador arbitrario). Un *clasificador arbitrario* se define como una aplicación  $g: \mathcal{X} \rightarrow \mathcal{Y}$  del conjunto de características  $\mathcal{X}$  al conjunto de resultados  $\mathcal{Y}$ , de forma que  $g(\mathbf{x})$  es el resultado predicho para el individuo  $\mathbf{x}$ .

**Definición 60** (Función de pérdida). Una *función de pérdida* es una función definida como  $\ell: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$  que asigna un valor real no-negativo  $\ell(y', y)$  que denota el coste del valor de la predicción  $y'$  cuando la etiqueta real es  $y$ .

**Definición 61** (Riesgo empírico). El *riesgo empírico* de un clasificador arbitrario  $g$  con respecto a un conjunto de datos  $\mathcal{D}$  dada una función de pérdida  $\ell$  se define como

$$R_{\mathcal{D}}(g) = \frac{1}{n} \sum_{i=1}^n \ell(g(\mathbf{x}_i), y_i).$$

**Definición 62** (Minimización del riesgo empírico). La *minimización del riesgo empírico* es el problema de optimización que busca encontrar un clasificador  $g$  en una familia de funciones  $\mathcal{H}$  tal que minimice el riesgo empírico,

$$\arg \min_{g \in \mathcal{H}} R_{\mathcal{D}}(g).$$

La introducción del concepto de minimización del riesgo empírico genera diversas dudas que intentaremos abordar a lo largo de este capítulo y que pueden resumirse en las siguientes tres preguntas:

1. ¿Cuál es la clase de funciones  $\mathcal{H}$  que deberíamos escoger?
2. ¿Cómo podemos resolver de manera eficiente el problema de optimización resultante?

3. ¿El clasificador encontrado tendrá el mismo rendimiento sobre los ejemplos del conjunto de entrenamiento que sobre el conjunto de datos de prueba?

Estas cuestiones están relacionadas entre sí y dan lugar a los conceptos de *representación*, *optimización* y *generalización* respectivamente.

### *Representación*

Normalmente las aplicaciones que utilizan datos con un número relativamente pequeño de características suelen usar *modelos de predicción lineales*. Por otro lado, si los datos de entrenamiento del modelo incluyen imágenes o audio, se suelen aplicar *modelos no lineales*. Las redes neuronales por ejemplo, aplican una secuencia de transformaciones de cada tipo para obtener mejores resultados sobre el conjunto de datos de entrada.

Lo más relevante de la representación para este trabajo es conocer que la mecánica de entrenamiento de un modelo sigue siendo la misma, y los detalles del tipo de modelo utilizado rara vez importa para las cuestiones relativas a la equidad.

### *Optimización*

Si nuestro objetivo es minimizar la exactitud de un clasificador, sería evidente pensar en resolver directamente el problema de minimización del riesgo empírico con respecto a la siguiente función de pérdida

$$\ell(y', y) = \begin{cases} 1, & \text{si } y \neq y', \\ 0, & \text{si } y = y'. \end{cases}$$

El problema de usar esta función es que es difícil de optimizar. Los gradientes de la pérdida 0-1 toman el valor cero en todo su dominio, por lo que no podemos esperar que los métodos basados en el gradiente optimicen directamente la pérdida 0-1.

A continuación, veremos una serie de métodos de optimización diferentes que, en determinadas circunstancias, encuentran un mínimo global o local del objetivo de riesgo empírico y aproximan en cierta medida la pérdida 0-1.

- **Pérdida al cuadrado** (*squared loss*) dada por  $\frac{1}{2}(y - y')^2$ . La minimización empírica del riesgo con esta función de pérdida equivale a la regresión lineal por mínimos cuadrados.
- **Pérdida de bisagra** (*hinge loss*) se expresa como  $\max\{1 - yy', 0\}$ . Los algoritmos SVM se refieren a la minimización empírica del riesgo con esta función junto con la regularización  $\ell_2$ .



- **Pérdida logística** (*logistic loss*) definida como

$$\begin{cases} -\log(\sigma(y')), & \text{si } y = 1, \\ -\log(1 - \sigma(y')), & \text{si } y = -1. \end{cases}$$

Donde  $\sigma(y') = 1/(1 + e^{-y'})$  es la función logística. La minimización empírica del riesgo con esta función de pérdida equivale a la regresión logística.

La elección de la función de pérdida se realizará comparando los rendimientos de las diferentes aproximaciones mediante prueba y error, eligiendo la que mejor funcione en cada caso.

### Generalización

La generalización en el aprendizaje automático hace referencia a cómo de bueno es un modelo predictivo que etiqueta correctamente datos con los que ha entrenado previamente realizando la misma tarea sobre un nuevo conjunto de datos que sigue la misma distribución de la que se extrajeron los datos de entrenamiento.

No obstante, incluso los modelos más avanzados suelen funcionar peor cuando los datos de prueba se extraen de una distribución que difiere ligeramente de la seguida por los datos de entrenamiento. Un ejemplo de ello fue el caso de la creación de un nuevo conjunto de pruebas para la base de datos ImageNet (Recht et al. [2019]).

## 6.3 CREACIÓN DE MODELOS DE APRENDIZAJE

Dado el esquema de configuración del problema de aprendizaje supervisado de la Figura 7 discutiremos la creación de un modelo simple de aprendizaje (Abu-Mostafa et al. [2012]). Sea  $\mathcal{X} = \mathbb{R}^d$  el espacio de entrada, donde  $\mathbb{R}^d$  es el espacio euclídeo  $d$ -dimensional y sea  $\mathcal{Y} = \{-1, 1\}$  el espacio de salida denotando una decisión binaria (si/no). En el ejemplo de concesión de crédito, las diferentes coordenadas del vector  $\mathbf{x} \in \mathcal{X}$  corresponden a los datos relativos del individuo que solicita el crédito. La salida binaria  $y$  corresponde a la aprobación o denegación del préstamo. Especificamos el conjunto de hipótesis  $\mathcal{H}$  mediante una forma funcional común a todas las hipótesis  $h \in \mathcal{H}$ . La forma funcional  $h(\mathbf{x})$  elegida, asigna pesos diferentes a cada coordenada del vector  $\mathbf{x}$ , reflejando su importancia en la decisión del problema. Las coordenadas ponderadas se combinan para formar una puntuación y el resultado se compara con un valor umbral previamente establecido. Si el solicitante supera el umbral, el crédito es aprobado, si no, es denegado:

$$\text{Aprobar crédito si } \sum_{i=1}^d w_i x_i > \text{umbral},$$

$$\text{Denegar crédito si } \sum_{i=1}^d w_i x_i < \text{umbral}.$$

Esta fórmula se puede escribir de forma más compacta como

$$h(\mathbf{x}) = \text{sign} \left( \left( \sum_{i=1}^d w_i x_i \right) + b \right), \quad (2)$$

donde  $x_1, \dots, x_d$  son los componentes del vector  $\mathbf{x}$ ;  $h(\mathbf{x}) = 1$  significa la aprobación del crédito y  $h(\mathbf{x}) = -1$  significa la denegación del crédito;  $\text{sign}(s) = 1$  si  $s > 0$  y  $\text{sign}(s) = -1$  si  $s < 0$ . Los pesos  $w_1, \dots, w_d$  y el umbral viene determinado en términos del sesgo  $b$ , el crédito se aprueba si  $\sum_{i=1}^d w_i x_i > -b$ .

Uno de los ejemplos más comunes en la literatura del aprendizaje automático es el del *perceptrón* introducido por el neurocientífico e informático teórico Rosenblatt [1957]. Este algoritmo de aprendizaje intenta encontrar  $\mathcal{H}$  buscando los pesos y el sesgo que funcionen bien en el conjunto de datos. Algunos de los pesos  $w_1, \dots, w_d$  podrían acabar siendo negativos, teniendo un efecto adverso en la aprobación del crédito. Por ejemplo, el peso del campo relativo a deudas pendientes debería ser negativo, ya que una mayor deuda no es buena señal para la aprobación de un crédito. El valor del sesgo  $b$  podría terminar siendo muy grande o muy pequeño, reflejando lo indulgente o estricto que debe ser el banco a la hora de conceder créditos. La elección óptima de los pesos y el sesgo define la hipótesis final  $g \in \mathcal{H}$  que produce el algoritmo.

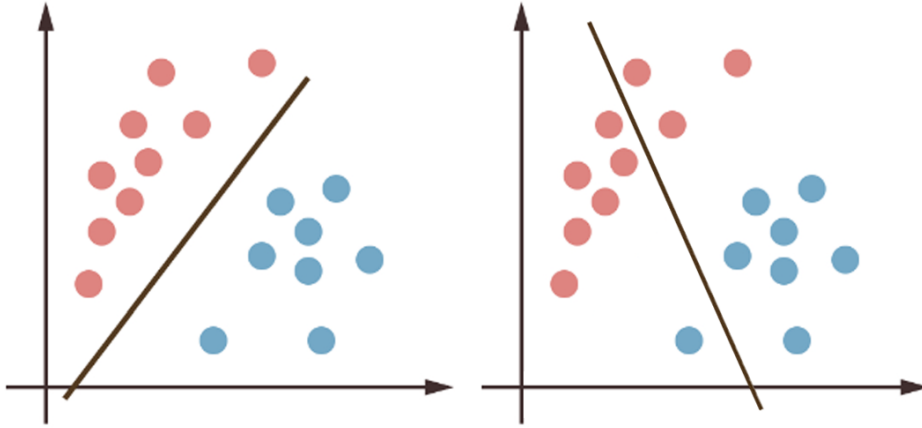


Figura 8: Clasificación de datos linealmente separables en un espacio bidimensional.

En la Figura 8 ilustra dos ejemplos de clasificación de un perceptrón en un espacio de entrada bidimensional ( $d = 2$ ). En el dibujo de la izquierda los ejemplos están perfectamente clasificados mientras que en el de la derecha hay algunos ejemplos mal clasificados. Diferentes valores para los parámetros  $b, w_1, w_2$  dan lugar a diferentes rectas  $w_1 x_1 + w_2 x_2 + b = 0$ . Si el conjunto de datos es linealmente separable existe una elección de parámetros que clasifica todos los ejemplos correctamente.

### 6.3.1 Ejemplo: Perceptrón

Introduciremos el algoritmo de del perceptrón como un ejemplo simple de un modelo de aprendizaje. Para simplificar la notación de la fórmula del perceptrón, trataremos el sesgo  $b$  como un peso  $w_0 = b$  y lo añadiremos como una coordenada más al vector de pesos  $\mathbf{w} = (w_0, w_1, \dots, w_d)^T$ . Además añadimos una coordenada  $x_0 = 1$  al vector  $\mathbf{x} = (x_0, x_1, \dots, x_d)^T$ . Observar que tratamos a  $\mathbf{x}$  y  $\mathbf{w}$  como vectores columna. Formalmente denotaremos el espacio de entrada como

$$\mathcal{X} = \{1\} \times \mathbb{R}^d = \{\mathbf{x} = (x_0, x_1, \dots, x_d)^T : x_0 = 1, x_1, \dots, x_d \in \mathbb{R}\}.$$

Denotando  $\mathbf{w}^T \mathbf{x} = \sum_{i=0}^d w_i x_i$ , podemos reescribir la Ecuación 2 como

$$h(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x}).$$

El algoritmo determinará el valor de  $\mathbf{w}$  basado en los datos. Supondremos que el conjunto de datos es linealmente separable, lo que significa que podemos encontrar un vector  $\mathbf{w}$  tal que  $h(\mathbf{x})$  consigue una decisión correcta  $h(\mathbf{x}_n) = y_n$  para todos los ejemplos del conjunto de datos de entrenamiento.

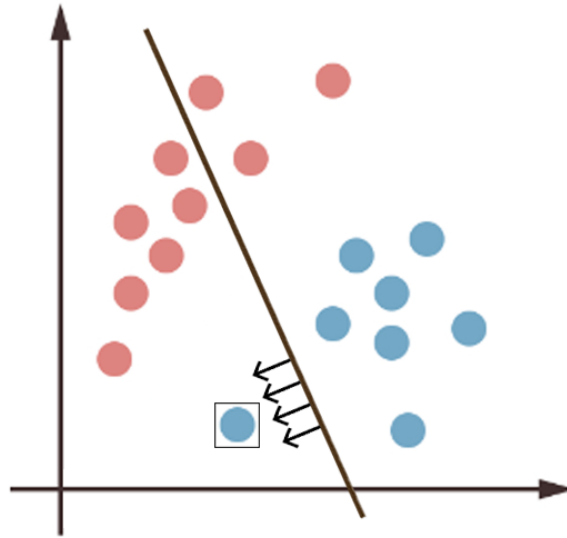


Figura 9: Esquema de actualización del algoritmo del perceptrón.

Nuestro algoritmo de aprendizaje encontrará este  $w$  usando un simple método iterativo, que funciona de la siguiente forma:

En cada paso  $t = 0, 1, \dots$ , hay un valor actual del vector de pesos  $\mathbf{w}_t$ , seleccionamos aleatoriamente un índice  $i \in \{1, \dots, n\}$  correspondiente a un ejemplo actualmente mal clasificado lo denota como  $(\mathbf{x}_t, y_t)$  y lo usa para actualizar el valor de  $\mathbf{w}_t$ . Como el

ejemplo está mal clasificado, tenemos que  $y_t \neq \text{sign}(\mathbf{w}_t^T \mathbf{x}_t)$ . La regla de actualización es

$$\mathbf{w}_{t+1} = \mathbf{w}_t + y_t \mathbf{x}_t.$$

Esta regla mueve la frontera con el objetivo de cambiar la dirección en la clasificación correcta de  $\mathbf{x}_t$ , como se puede ver en la Figura 9. El algoritmo continua con las sucesivas iteraciones hasta que no haya ejemplos mal clasificados en el conjunto de datos.

Aunque la regla de actualización solo considera un único ejemplo del conjunto de entrenamiento y podría desordenar la clasificación para otros ejemplos no implicados en la iteración actual, el algoritmo garantiza una solución óptima (véase Teorema 1 en Collins [2012]). El resultado se mantiene de forma independiente al ejemplo que elijamos y a la inicialización del vector de pesos al comienzo del algoritmo. Por simplicidad, elegiremos un ejemplo mal clasificado aleatorio e inicializaremos  $w_0$  a un vector de ceros.

#### *Otra definición del algoritmo*

Podemos definir el algoritmo de perceptrón como una instancia de la minimización del riesgo empírico (Barocas et al. [2019]). Por la descripción del algoritmo, buscamos un separador lineal y por tanto, nuestra clase de funciones corresponde con el conjunto de funciones lineales

$$\mathcal{H} = \{f(\mathbf{x}) = \langle \mathbf{w}, \mathbf{x} \rangle : \mathbf{w} \in \mathbb{R}^d\}.$$

Utilizaremos el método de gradiente estocástico como base del algoritmo, ya que es un método de optimización que elige un ejemplo aleatorio en cada paso y realiza una actualización de los parámetros del modelo. Se define con la regla siguiente:

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \nabla_{\mathbf{w}_t} \ell(f(\mathbf{x}_t), y_t).$$

Donde  $\nabla \ell(f(\mathbf{x}_t), y_t)$  es el gradiente de la función de pérdida con respecto a los parámetros del modelo  $\mathbf{w}_t$  para un ejemplo  $(\mathbf{x}_t, y_t)$ . El escalar  $\eta > 0$ , es un parámetro denominado *tamaño de paso*, pensaremos en él como una constante pequeña.

Consideramos ahora la función de pérdida

$$\ell(y, \langle \mathbf{w}, \mathbf{x} \rangle) = \max(1 - y \langle \mathbf{w}, \mathbf{x} \rangle, 0),$$

donde su gradiente se define como:

$$\nabla \ell(y, \langle \mathbf{w}, \mathbf{x} \rangle) = \begin{cases} -y\mathbf{x}, & \text{si } y \langle \mathbf{w}, \mathbf{x} \rangle < 1, \\ 0, & \text{si } y \langle \mathbf{w}, \mathbf{x} \rangle > 1. \end{cases}$$

La expresión anterior define una parte de la regla de actualización del perceptrón, la otra parte la deduciremos al añadir la penalización  $\frac{\alpha}{2} \|\mathbf{w}\|^2$  (donde  $\|\cdot\|$  denota la norma euclídea) a la función de pérdida para impedir a los pesos tomar valores por encima de un umbral establecido. Esta penalización se denomina *regularización  $\ell_2$*  o *regularización de Tíjonov* y su propósito es fomentar la generalización.

Sumando las dos funciones de pérdida, obtenemos una expresión del riesgo empírico regularizado por  $\ell_2$  para la función de pérdida de bisagra:

$$R_{\mathcal{D}}(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \max(1 - y_i \langle \mathbf{w}, \mathbf{x}_i \rangle, 0) + \frac{\alpha}{2} \|\mathbf{w}\|^2.$$

Definiremos el algoritmo perceptrón como la resolución a este problema de minimización del riesgo empírico utilizando el método del gradiente estocástico.

### 6.3.2 Regresión lineal

Volviendo al ejemplo del apartado 6.1.1, recordemos que el banco tiene un registro de clientes con variables que pueden ser usadas para aprender un clasificador lineal de decisión para la aprobación del crédito. En este caso, en lugar de limitarnos a tomar una decisión binaria (aprobar o no el crédito), en el caso de aprobación, podríamos querer establecer un umbral en el valor del crédito concedido. Esta tarea, puede ser automatizada haciendo uso del aprendizaje por *regresión* (Abu-Mostafa et al. [2012]).

El banco parte de un conjunto de datos  $\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$ , donde  $\mathbf{x}_n$  es la información del cliente e  $y_n$  es el límite de crédito establecido por uno de los expertos del banco. Ahora  $y_n$  será un número real en lugar de un valor binario. El banco querrá usar un modelo de aprendizaje para encontrar una hipótesis  $g$  que replique la actuación humana en los límites de crédito. En este caso, no buscaremos una función determinista  $y = f(x)$  y en su lugar, asumiremos que la etiqueta  $y_n$  proviene de una distribución  $P(y | \mathbf{x})$ . No obstante, la naturaleza del problema sigue siendo la misma. Tenemos una distribución desconocida  $P(\mathbf{x}, y)$  que genera cada  $(\mathbf{x}_n, y_n)$  y queremos encontrar una hipótesis  $g$  que minimice el error entre  $g(\mathbf{x})$  e  $y$  con respecto a esa distribución.

#### Algoritmo de regresión lineal

El algoritmo se basa en la minimización del riesgo empírico para la pérdida al cuadrado entre  $h(\mathbf{x})$  e  $y$ . El problema es equivalente a minimizar la siguiente función:

$$R_{\mathcal{D}}(h) = \frac{1}{n} \sum_{i=1}^n (h(\mathbf{x}_i) - y_i)^2.$$

En regresión lineal,  $h$  se expresa como combinación lineal de las componentes de  $\mathbf{x}$ :

$$h(\mathbf{x}) = \sum_{i=0}^d w_i x_i = \mathbf{w}^T \mathbf{x},$$

donde  $x_0 = 1$  y  $\mathbf{x} \in \{1\} \times \mathbb{R}^d$  y  $\mathbf{w} \in \mathbb{R}^{d+1}$ . Para el caso lineal se suele definir una matriz de representación para  $R_{\mathcal{D}}(h)$ . La matriz de datos  $X \in \mathbb{R}^{N \times (d+1)}$  tiene como filas los vectores  $\mathbf{x}_n$  e  $\mathbf{y} \in \mathbb{R}^N$  como vector columna cuyas componentes son los valores objetivo  $y_n$ . Podemos expresar el error en función de  $\mathbf{w}$ ,  $X$  e  $\mathbf{y}$  como:

$$\begin{aligned} R_{\mathcal{D}}(\mathbf{w}) &= \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}_i - y_i)^2 \\ &= \frac{1}{n} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 \\ &= \frac{1}{n} (\mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} - 2\mathbf{w}^T \mathbf{X}^T \mathbf{y} + \mathbf{y}^T \mathbf{y}). \end{aligned} \quad (3)$$

Nuestro problema se reduce a minimizar la función dada por:

$$\arg \min_{\mathbf{w} \in \mathbb{R}^{d+1}} R_{\mathcal{D}}(\mathbf{w}). \quad (4)$$

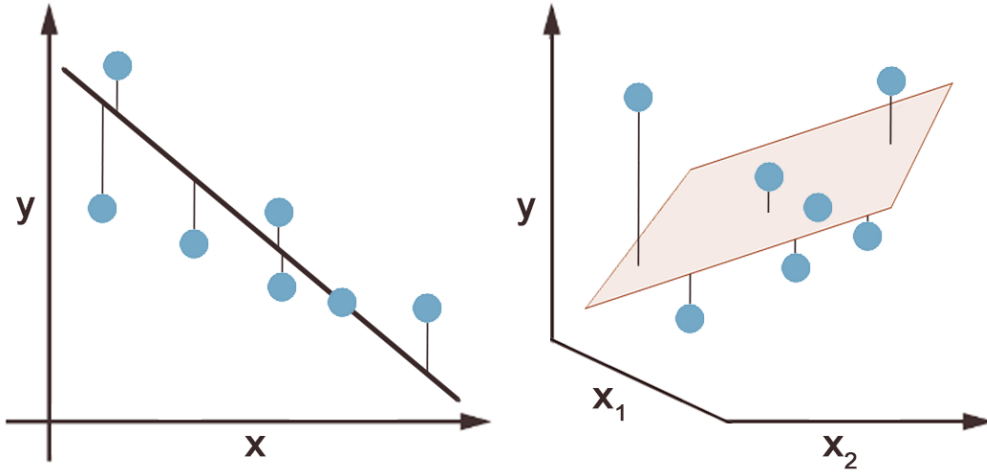


Figura 10: Problema de regresión lineal en una y dos dimensiones respectivamente.

La Figura 10 ilustra la solución para los casos unidimensional y bidimensional respectivamente. La Ecuación 3 implica que  $R_{\mathcal{D}}(\mathbf{w})$  es diferenciable, por lo que podemos encontrar el mínimo de esta función igualando su gradiente a cero.

$$\nabla R_{\mathcal{D}}(\mathbf{w}) = \frac{2}{n} (\mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{X}^T \mathbf{y}) = 0.$$

Para encontrar la solución se debe cumplir que  $\mathbf{X}^T \mathbf{X} \mathbf{w} = \mathbf{X}^T \mathbf{y}$ . Si  $\mathbf{X}^T \mathbf{X}$  es regular,  $\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$  es la única solución óptima para la Ecuación 4. En otro caso, existirá una solución óptima, pero no será única.

## 6.4 EVALUACIÓN EN APRENDIZAJE AUTOMÁTICO

Un modelo en aprendizaje automático funciona bien, cuando puede predecir resultados correctos a partir de un conjunto de datos de entrada no conocido previamente. A este proceso lo conocemos como generalización. La manera de medir este fenómeno, no es única, por lo que existen una gran cantidad de *métricas de rendimiento* en el ámbito del *machine learning*. La elección de la métrica dependerá del problema específico, su dominio y de las restricciones del mismo al mundo real. Sin embargo, la mayoría de las métricas vienen descritas como funciones de una *matriz de confusión*.

Una matriz de confusión define un modelo de predicción a partir de las dimensiones de los *valores reales de las etiquetas* y de sus posibles *predicciones*, resumiendo el número de predicciones correctas e incorrectas por clase.

		Etiqueta real	
		$y = 1$	$y = -1$
Predicción	$\hat{y} = 1$	Verdadero Positivo (TP)	Falso Positivo (FP) (Error tipo I)
	$\hat{y} = -1$	Falso Negativo (FN) (Error tipo II)	Verdadero Negativo (TN)

Cuadro 2: Matriz de confusión, ilustra la relación entre la etiqueta real y la predicción.

En el caso de una matriz de confusión binaria como la del Cuadro 2, tenemos una matriz  $2 \times 2$  que nos ofrece información sobre el número de *verdaderos positivos* (TP) ( $y = 1 \wedge \hat{y} = 1$ ), *falsos positivos* (FP) ( $y = -1 \wedge \hat{y} = 1$ ), *falsos negativos* (FN) ( $y = 1 \wedge \hat{y} = -1$ ) y *verdaderos negativos* (TN) ( $y = -1 \wedge \hat{y} = -1$ ). A partir de ellos, podemos obtener el número total de positivos predichos (TP+FP), el número total de negativos predichos (TN+FN), el número total de etiquetados positivos (TP+FN), y el número total de etiquetados negativos (TN+FP).

Además, podemos definir las siguientes métricas de clasificación avanzadas construidas a partir de combinaciones lineales de las básicas:

- **Tasa de verdaderos positivos** (*Recall*),

$$\text{TPR} = \frac{\text{TP}}{\text{TP} + \text{FN}}.$$

- **Tasa de falsos negativos**,

$$\text{FNR} = 1 - \text{TPR} = \frac{\text{FN}}{\text{TP} + \text{FN}}.$$

- Tasa de verdaderos negativos (*Specificity*),

$$\text{TNR} = \frac{\text{TN}}{\text{TN} + \text{FP}}.$$

- Tasa de falsos positivos,

$$\text{FPR} = 1 - \text{TNR} = \frac{\text{FP}}{\text{TN} + \text{FP}}.$$

- Valor positivo predictivo (*Precision*),

$$\text{PPV} = \frac{\text{TP}}{\text{TP} + \text{FP}}.$$

- Tasa de falso descubrimiento,

$$\text{FDR} = 1 - \text{PPV} = \frac{\text{FP}}{\text{TP} + \text{FP}}.$$

- Valor negativo predictivo,

$$\text{NPV} = \frac{\text{TN}}{\text{TN} + \text{FN}}.$$

- Tasa de falsa omisión,

$$\text{FOR} = 1 - \text{NPV} = \frac{\text{FN}}{\text{TN} + \text{FN}}.$$

- Exactitud (*Accuracy*),

$$\text{Accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}}.$$

		Etiqueta real			
		$y = 1$	$y = -1$		
Predicción	$\hat{y} = 1$	Verdadero Positivo (TP)	Falso Positivo (FP)	<b>Precision</b> $\frac{TP}{TP + FP}$	<b>Tasa de Falso Descubrimiento</b> $\frac{FP}{TP + FP}$
	$\hat{y} = -1$	Falso Negativo (FN)	Verdadero Negativo (TN)	<b>Tasa de Falsa Omisión</b> $\frac{FN}{TN + FN}$	<b>Valor Negativo Predictivo</b> $\frac{TN}{TN + FN}$
		<b>Recall</b> $\frac{TP}{TP + FN}$	<b>Tasa de Falsos Positivos</b> $\frac{FP}{TN + FP}$	<b>Accuracy</b> $\frac{TP + TN}{TP + TN + FP + FN}$	
		<b>Tasa de Falsos Negativos</b> $\frac{FN}{TP + FN}$	<b>Specificity</b> $\frac{TN}{TN + FP}$		

Cuadro 3: Matriz de confusión con métricas de evaluación avanzadas.



Normalmente mediremos el rendimiento de un modelo con una de las métricas mencionadas o con una combinación de ambas. Una de las medidas de rendimiento más usadas es el valor positivo predictivo o precisión, que mide el porcentaje de predicciones correctas realiza. No obstante, existen problemas derivados de utilizar la precisión a la hora de medir el rendimiento cuando abordamos problemas que contienen desequilibrio entre las clases del conjunto de datos. Por ejemplo, si solo 5 de cada 100 muestras es positiva, el modelo trivial que siempre prediga la clase negativa tendrá una precisión del 95 %.

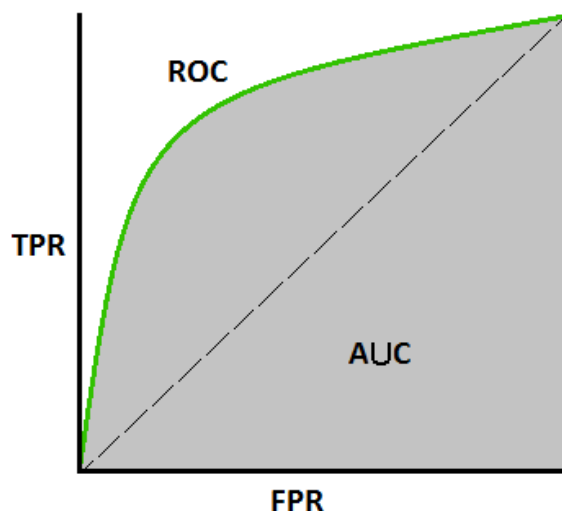


Figura 11: Ejemplo de un gráfico de curva ROC.

Otras métricas comúnmente utilizadas son la  $F_1$ -score la cual podemos calcularla como  $F_1\text{-score} = 2 \cdot \frac{PPV \cdot TPR}{PPV + TPR}$  y el *área bajo la curva ROC* (AUC). El espacio ROC se representa en un gráfico bidimensional, en el que la tasa de verdaderos positivos (TPR) se representa en el eje vertical y la tasa de falsos positivos (FPR) en el eje horizontal. Como podemos ver en la Figura 11, diferentes umbrales de clasificación corresponden a diferentes puntos en el espacio ROC.

El área bajo la curva ROC (AUC) se interpreta como la probabilidad de que el modelo clasifique un ejemplo positivo aleatorio más alto que un ejemplo negativo aleatorio. Al ser una probabilidad, el AUC oscilará entre los valores 0 y 1. Un modelo cuyas predicciones sean un 100 % incorrectas tendrá un AUC de 0, mientras que otro cuyas predicciones sean un 100 % correctas tendrá un AUC de 1. También es importante saber que el AUC es invariable con respecto a la escala y con respecto al umbral de clasificación.

---

## FORMALIZACIÓN DE LAS MEDIDAS DE EQUIDAD

---

En este capítulo presentaremos el concepto de equidad, realizaremos un análisis de las distintas nociones de equidad conocidas y formalizaremos sus definiciones.

### 7.1 ¿QUÉ ES LA EQUIDAD?

Con el aumento de los métodos de toma de decisiones automatizadas en la actualidad, la necesidad de satisfacer *equidad* en los modelos de *machine learning* ha cobrado importancia. Por ello, cabe hacerse las siguientes preguntas: ¿Qué es la equidad? ¿Cómo podemos medirla? ¿Y, cómo podemos fomentarla en nuestros algoritmos? En esta sección se intentará dar respuesta a estas preguntas.

En la Sección 6.1, formalizamos el proceso de aprendizaje automático sobre un ejemplo de aprobación de créditos bancarios. Hemos visto que existe un proceso de aprendizaje sobre registros históricos de datos que nos aportan información a la hora de predecir nuevos ejemplos, pero: ¿Existirán atributos que discriminen a un grupo determinado de la población?, ¿Dos clientes con características similares recibirán la misma predicción? Estas y otras preguntas son las motivaron el concepto de equidad.

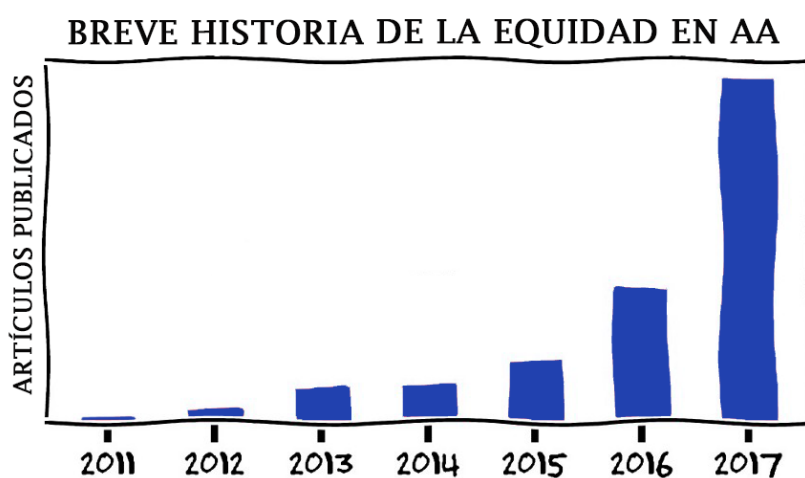


Figura 12: Incremento de las publicaciones sobre equidad entre 2011 y 2017.

Podemos dividir los trabajos sobre equidad en *machine learning* en, detectar el sesgo y discriminación en los modelos (Dwork et al. [2011]) y mitigar el sesgo algorítmico (Corbett-Davies et al. [2017]). Para estas tareas, al igual que en cualquier ciencia experimental, debemos ser capaces de medir el concepto partiendo de una definición teórica. La equidad es un concepto inherentemente subjetivo y que depende en gran medida del ámbito en el que lo apliquemos. Por lo tanto, a partir de conceptos de la literatura de la ciencias sociales, se han ido proponiendo diferentes medidas y formalizando estos conceptos para que puedan ser aplicadas al aprendizaje automático.

La primera idea es buscar apoyo legal y comprobar si existe alguna definición que pueda utilizarse para formular la equidad matemáticamente. Las leyes antidiscriminatorias de muchos países prohíben el trato desigual entre personas en función de atributos sensibles, tales como el sexo o la raza (Title VII of the Civil Rights Act: Equal Employment Opportunities). Estas leyes suelen evaluar la imparcialidad de un proceso de toma de decisiones utilizando dos nociones distintas (Barocas and Selbst [2016]): el tratamiento dispar y el impacto dispar. Un proceso de toma de decisiones sufrirá un trato dispar si basamos su juicio en el atributo sensible del sujeto, y tendrá un impacto dispar si sus resultados perjudican (o benefician) de forma desigual a personas con valores de atributos sensibles diferentes (por ejemplo, mujeres o afroamericanos).

#### 7.1.1 Principales familias de las medidas de equidad

Los conceptos anteriores son demasiado abstractos como para tener una formulación cuantitativa directa por lo que, poco a poco, se han ido añadiendo numerosas definiciones de equidad a la literatura del aprendizaje automático. Sin embargo, el Cuadro 1 recoge los criterios más importantes, los cuales han sido previamente recopilados en trabajos como Gajane and Pechenizkiy [2018] o Verma and Rubin [2018].

En este capítulo, nos centraremos en las medidas de equidad relativas a *impacto y tratamiento dispar*, estableciendo a su vez una división más específica en algunas de ellas. Los conceptos que trataremos a lo largo de este capítulo son:

- **Equidad por desconocimiento.**
- **Equidad individual.**
- **Equidad de grupo:** formalizaremos la paridad demográfica, el criterio de probabilidades igualadas y la tasa de paridad predictiva.
- **Medidas causales:** analizaremos su base matemática y desarrollaremos un ejemplo práctico de su aplicación en la Parte IV de este trabajo.

La equidad basada en *preferencias* está fuera del marco de este trabajo, por lo que no será explicado en profundidad. Para el lector interesado en este concepto, puede consultar el artículo Zafar et al. [2017c].

### 7.1.2 Medición de la parcialidad y la equidad

Consideramos una tarea de la clasificación estándar en la que el objetivo es predecir una variable de resultado binaria  $y \in \mathcal{Y}$  utilizando un vector de variables de entrada  $\mathbf{x} \in \mathcal{X}$  que sigue una distribución de probabilidad  $P_{\mathbf{x}}$ .

Sea un clasificador arbitrario  $g : \mathcal{X} \rightarrow \mathcal{Y}$ , donde  $\mathcal{X} = \mathbb{R}^d$  y donde  $\mathcal{Y} = [0, 1]$  si produce una probabilidad predicha (por ejemplo, regresión logística) y  $\mathcal{Y} = \{-1, 1\}$  si el clasificador produce un resultado predicho (por ejemplo, SVM). A lo largo del proyecto, normalmente asumiremos que el espacio de salida trabaja sobre decisiones binarias (si/no), es decir,  $\mathcal{Y} = \{-1, 1\}$ .

Muchos de los problemas en los que aparece el concepto de equidad pueden formularse como problemas estadísticos de evaluación de riesgos en los que asignamos una puntuación de valor real  $s \in [0, 1]$  a cada individuo del conjunto de datos y se toma una decisión  $\hat{y}$  basada en la puntuación, normalmente seleccionando un número predefinido ( $k$ ) de entidades que deben clasificarse como positivas.

Las principales definiciones que usaremos en este capítulo son las siguientes:

- **Vector de características** -  $\mathbf{x} \in \mathcal{X}$  es un vector de características reales que identifican a un individuo.
- **Puntuación** -  $s \in [0, 1]$  es una puntuación de valor real asignada a cada entidad por el clasificador.
- **Predicción** -  $\hat{y} \in \mathcal{Y}$  es una predicción binaria asignada a un individuo determinado, basada en el umbral de la puntuación (por ejemplo, el máximo  $k$ ).
- **Etiqueta real** -  $y \in \{-1, 1\}$  es la etiqueta binaria que representa el valor real de un individuo en concreto.

## 7.2 EQUIDAD POR DESCONOCIMIENTO

La *equidad por desconocimiento* se basa en eliminar los atributos sensibles para todos los individuos en el proceso de predicción. Algunos clasificadores propuestos en la literatura de aprendizaje automático satisfacen esta medida (Dwork et al. [2011]) debido a que es intuitiva y muy fácil de aplicar.

**Notación 2.** Sea  $\Delta = \{\pi(\mathbf{x}) : \mathbb{R}^d \rightarrow \mathbb{R}^m : \mathbf{x} \in \mathcal{X}, 1 \leq m \leq d\}$  el conjunto formado por las proyecciones del conjunto de individuos a todas las posibles dimensiones menores o iguales que la dimension de  $\mathcal{X}$ . Sea  $\mathcal{A} \subset \Delta$  un subconjunto que contiene todos los individuos de  $\mathcal{X}$  a los que se ha aplicado una proyección  $\pi$  sobre ciertas características. Notaremos como  $\mathcal{A}_i$  al elemento de  $\mathcal{A}$  que contiene las características del individuo  $\mathbf{x}_i \in \mathcal{X}$  que cierta proyección  $\pi$  ha seleccionado.

**Definición 63** (Equidad por desconocimiento). Sea  $\mathcal{A} \subset \Delta$  donde  $\mathcal{A}_i$  contiene las proyecciones sobre los atributos sensibles para el individuo  $\mathbf{x}_i$  y sean  $g: \mathcal{X} \rightarrow \mathcal{Y}$  y  $h: \mathcal{X} \setminus \mathcal{A} \rightarrow \mathcal{Y}$  dos clasificadores arbitrarios. Diremos que  $g$  logra *equidad por desconocimiento* si, y solo si,

$$g(\mathbf{x}_i) = h(\mathbf{x}_i \setminus \mathcal{A}_i), \text{ para todo } \mathbf{x}_i \in \mathcal{X}.$$

Uno de los principales problemas de la equidad por desconocimiento es que no da una condición suficiente para evitar la discriminación ya que puede haber muchas características altamente correlacionadas (por ejemplo, la zona de residencia) que funcionen como sustitutos del atributo sensible (por ejemplo, la raza). Por lo tanto, no bastaría con eliminar el atributo sensible para eliminar las disparidades. Además, se han documentado diversos ejemplos de equidad por desconocimiento para la raza en ámbitos como educación, concesión de préstamos o justicia penal y se ha demostrado que, a largo plazo, el enfoque ciego de la raza es menos eficaz que el enfoque consciente de la misma (Fryer et al. [2008]).

Las críticas anteriores cuestionan la idoneidad de la equidad por desconocimiento en los dominios en los que, los atributos sensibles pueden deducirse a partir de los atributos no sensibles disponibles y tenemos conocimiento de la existencia de barreras estructurales, que obstaculizan a los grupos desfavorecidos, a partir de encuestas verosímiles sobre los grupos demográficos.

*Ejemplo 9.* Supongamos un modelo utilizado para aprobar o denegar créditos bancarios. Por sesgos históricos, sabemos que uno de los atributos sensibles en la concesión de préstamos, es la raza. Procedemos entonces a eliminar esta información de todos los individuos en el modelo de predicción. El problema surge cuando notamos que el código postal, otro atributo presente en el vector de características, está altamente correlacionado con la raza y por tanto, las decisiones basadas en este serán racialmente discriminatorias. En consecuencia, el criterio de equidad por desconocimiento en este caso concreto, sería insuficiente.

Esta práctica se conoce como *redlining*, cuyo término fue acuñado en la década de 1960 debido a la práctica de negar bienes y servicios a las minorías mediante la *redlining* de barrios específicos en un mapa (Custers et al. [2012]).

### 7.3 EQUIDAD INDIVIDUAL

La *equidad individual* se basa en métricas de similitud sobre los atributos y establece que individuos similares deben recibir predicciones similares independientemente del atributo sensible (Dwork et al. [2011]). Además, la equidad individual es más precisa que la equidad de grupo, ya que impone restricciones en el tratamiento de cada par de individuos.

**Definición 64** (Equidad individual). Sea  $g: \mathcal{X} \rightarrow \mathcal{Y}$  un clasificador arbitrario,  $D: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$  una medida de distancia sobre el espacio de clasificación resultante y  $d: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  una medida de distancia sobre los individuos, se dice que  $g$  cumple con la *equidad individual* si, y solo si,

$$D(g(x_i), g(x_j)) \leq d(x_i, x_j), \text{ para todo } x_i, x_j \in \mathcal{X}.$$

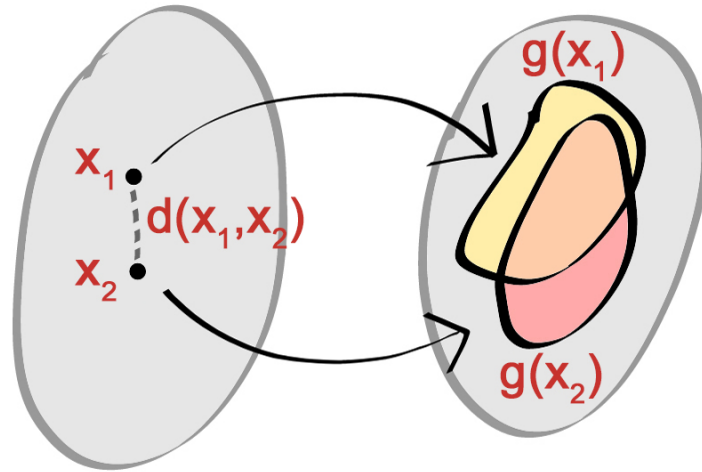


Figura 13: Ilustración de la noción de equidad individual.

La definición de equidad individual, también es conocida como *propiedad  $(D, d)$ -Lipschitz*. Cualquier clasificador que satisfaga esta propiedad también verificará la paridad demográfica con un cierto sesgo (véase Lema 3.1 en [Dwork et al. \[2011\]](#)).

En la literatura de las ciencias sociales, esta formalización equivale al individualismo igualitario, conocido por ser el principio formal de justicia. Esta noción responsabiliza a la *métrica de la distancia* de garantizar la justicia del clasificador. Si la métrica de la distancia utiliza los atributos sensibles directa o indirectamente para calcular la distancia entre dos individuos, un clasificador que satisfaga la Definición 64 podría seguir causando impacto dispar. Por tanto, la equidad individual, no podría considerarse adecuada para dominios en los que no se dispone de una métrica de distancia fiable y no discriminatoria.

*Ejemplo 10.* Imaginemos tres candidatos a un puesto de trabajo,  $A$ ,  $B$  y  $C$ .  $A$  tiene únicamente el título de graduado y un año de experiencia laboral relacionada.  $B$  tiene un máster y un año de experiencia laboral relacionada.  $C$  tiene un doble grado pero no tiene experiencia laboral. En principio, no podemos disponer del rendimiento de los tres individuos ya que no podemos contratar a todos. Entonces: ¿Está  $A$  más cerca de  $B$  que de  $C$ ? Si es así, ¿por cuánto? La cosa se complica aún más cuando entran en juego los atributos sensibles. ¿Cómo deberíamos cuantificar la diferencia de pertenencia a un grupo en nuestra función métrica?

En este ejemplo podemos observar los problemas de la equidad individual comentados anteriormente y cómo dependen directamente de la construcción de la métrica de distancia entre individuos.

#### 7.4 EQUIDAD DE GRUPO

La *equidad de grupo* mide el impacto dispar entre grupos desfavorecidos y privilegiados, como podrían ser grupos de diferentes razas, edad o género.

Supongamos que en  $\mathcal{X}$  se definen todas los posibles valores para las características {raza, género, salario, trabajo, edad}, denotaremos por  $\mathcal{A}$  el conjunto de atributos sensibles sobre  $\mathcal{X}$  considerando por ejemplo, en este caso,

$$\mathcal{A}=\{\text{raza, género, edad}\}.$$

Consideremos un *atributo multivaluado*  $a = \{a_1, \dots, a_n\} \in \mathcal{A}$ , por ejemplo,

$$\text{raza}=\{\text{hispanico, caucasico, afroamericano, otro}\}.$$

Definimos un *grupo*  $G(a_i)$  como un conjunto de entidades que tienen en común un valor específico del atributo  $a$ , por ejemplo raza=hispanico corresponde a todos los individuos de raza hispanica del conjunto de datos.

Teniendo en cuenta todos los grupos definidos por el atributo  $a$ , las predicciones  $\hat{y}$  y la etiqueta real  $y$  para cada entidad de cada grupo, podemos hablar ahora de las métricas de grupo. Las principales definiciones sobre grupos para evaluar el sesgo y la equidad son las siguientes:

- **Atributo** -  $a = \{a_1, \dots, a_n\}$  es un atributo multivaluado, por ejemplo,  
 $\text{raza}=\{\text{hispanico, caucasico, afroamericano, otro}\}.$
- **Grupo** -  $G(a_i)$  es un grupo de todas las entidades que comparten el mismo valor de atributo  $a = a_i$ , por ejemplo, raza=hispanico.
- **Grupo de objetivo** -  $G(a_o)$  es un grupo utilizado como objetivo del cálculo de las medidas de sesgo.
- **Grupo de referencia** -  $G(a_r)$  es el grupo que se utiliza como referencia para calcular las medidas de sesgo. Suele fijarse siguiendo un criterio determinado.
- **Etiquetado positivo** -  $LP_G$  número de entidades etiquetadas como positivas dentro de un grupo.
- **Etiquetado negativo** -  $LN_G$  número de entidades etiquetadas como negativas dentro de un grupo.

**Notación 3.** De aquí en adelante, utilizaremos la siguiente notación:

- $Y$  es una variable aleatoria binaria que representa la etiqueta real de un individuo de  $\mathcal{X}$ .

- $\hat{Y}$  es una variable aleatoria binaria que representa el resultado de la predicción de un clasificador  $g: \mathcal{X} \rightarrow \mathcal{Y}$  para un individuo de  $\mathcal{X}$ .
- $A$  es una variable aleatoria binaria que representa si un individuo de  $\mathcal{X}$  pertenece al grupo objetivo ( $a_o$ ) o de referencia ( $a_r$ ).

### *Métricas de grupos de distribución*

Definiremos *métricas de decisión* a nivel de grupo centradas en la distribución de los individuos entre los grupos del conjunto seleccionado para la intervención (máximo  $k$ ) y por tanto, no requieren del valor de la etiqueta  $Y$ . Definimos las métricas de distribución de los grupos como sigue:

- **Positivos predichos** -  $PP_G$  número de entidades dentro de un grupo donde la decisión es positiva, es decir,  $\hat{Y} = 1$ .
- **Negativos predichos** -  $PN_G$  número de entidades dentro de un grupo cuya decisión es negativa, es decir,  $\hat{Y} = -1$ .
- **Total de predicciones positivas** - número total de entidades predichas como positivas en los grupos definidos por  $a$ ,

$$K = \sum_{i=1}^n PP_{G(a_i)}.$$

- **Prevalencia predicha** - fracción de entidades dentro de un grupo que se predijo como positiva,

$$PP_{\text{prev}_G} = \frac{PP_G}{|G|} = P(\hat{Y} = 1 \mid A = a_i).$$

- **Tasa de positivos predichos** - fracción de las entidades predichas como positivas que pertenecen a un determinado grupo,

$$PPR_G = \frac{PP_G}{K} = P(A = a_i \mid \hat{Y} = 1).$$

### *Métricas de grupo basadas en la etiqueta real*

A continuación, discutiremos diferentes métricas que surgen dependiendo de la coincidencia o no entre los valores de predicción  $\hat{Y}$  y, en este caso, de la etiqueta real  $Y$ . La mayoría de ellas, ya fueron presentadas en la Sección 6.4. Las métricas de grupo basadas en los errores y aciertos son las siguientes:

- **Falso positivo** -  $FP_G$  es el número de entidades del grupo con

$$\hat{Y} = 1 \wedge Y = -1.$$

- **Falso negativo** -  $FN_G$  es el número de entidades del grupo con

$$\hat{Y} = -1 \wedge Y = 1.$$



- **Verdadero positivo** -  $TP_G$  es el número de entidades del grupo con

$$\hat{Y} = 1 \wedge Y = 1.$$

- **Verdadero negativo** -  $TN_G$  es el número de entidades del grupo con

$$\hat{Y} = -1 \wedge Y = -1.$$

- **Prevalencia** - fracción de entidades dentro de un grupo cuyo resultado verdadero fue positivo,

$$Prev_G = \frac{LP_G}{|G|} = P(Y = 1 \mid A = a_i).$$

- **Tasa de falso descubrimiento** - fracción de falsos positivos de un grupo dentro de los positivos predichos del grupo,

$$FDR_G = \frac{FP_G}{PP_G} = P(Y = -1 \mid A = a_i, \hat{Y} = 1).$$

- **Tasa de falsa omisión** - fracción de falsos negativos de un grupo dentro de los negativos predichos del grupo,

$$FOR_G = \frac{FN_G}{PN_G} = P(Y = 1 \mid A = a_i, \hat{Y} = -1).$$

- **Tasa de falsos positivos** - fracción de falsos positivos de un grupo dentro de los negativos etiquetados del grupo,

$$FPR_G = \frac{FP_G}{LN_G} = P(\hat{Y} = 1 \mid A = a_i, Y = -1).$$

- **Tasa de falsos negativos** - fracción de falsos negativos de un grupo dentro de los positivos etiquetados del grupo,

$$FNR_G = \frac{FN_G}{LP_G} = P(\hat{Y} = -1 \mid A = a_i, Y = 1).$$

- **Valor negativo predictivo** - fracción de verdaderos negativos de un grupo dentro de los negativos predichos del grupo,

$$NPV_G = \frac{TN_G}{PN_G} = P(Y = -1 \mid A = a_i, \hat{Y} = -1).$$

- **Valor positivo predictivo (Precision)** - fracción de verdaderos positivos de un grupo dentro de los positivos predichos del grupo,

$$PPV_G = \frac{TP_G}{PP_G} = P(Y = 1 \mid A = a_i, \hat{Y} = 1).$$

- **Tasa de verdaderos positivos (Recall)** - fracción de verdaderos positivos de un grupo dentro de los positivos etiquetados del grupo,

$$TPR_G = \frac{TP_G}{LP_G} = P(\hat{Y} = 1 \mid A = a_i, Y = 1).$$

- **Tasa de verdaderos negativos (Specificity)** - fracción de verdaderos negativos de un grupo dentro de los negativos etiquetados del grupo,

$$\text{TNR}_G = \frac{\text{TN}_G}{\text{LN}_G} = P(\hat{Y} = -1 \mid A = a_i, Y = -1).$$

- **Exactitud (Accuracy)** - fracción de resultados verdaderos de un grupo dentro del total de casos examinados del grupo,

$$\text{Accuracy}_G = \frac{\text{TP}_G + \text{TN}_G}{\text{LP}_G + \text{LN}_G} = P(\hat{Y} = Y \mid A = a_i).$$

- **Tasa global de clasificación errónea.** - fracción de resultados falsos de un grupo dentro del total de casos examinados del grupo,

$$\text{OMR}_G = \frac{\text{FP}_G + \text{FN}_G}{\text{LP}_G + \text{LN}_G} = P(\hat{Y} \neq Y \mid A = a_i).$$

En los apartados siguientes, formalizaremos algunas de las nociones populares de equidad de grupo que podemos encontrar en la literatura. Sea un atributo sensible multivaluado  $a \in \mathcal{A}$ , las métricas relativas a equidad de grupo se definen como una igualdad entre las probabilidades de un *grupo objetivo* ( $a_o$ ) en comparación con un *grupo de referencia* ( $a_r$ ).

El grupo de referencia se suele seleccionar en base a diferentes criterios. Por ejemplo, se podría utilizar el grupo mayoritario entre los grupos definidos por  $A$ , o el enfoque tradicional de fijar un grupo históricamente favorecido, por ejemplo, en el caso de la raza, los individuos de raza caucásica.

#### 7.4.1 Paridad demográfica

La *paridad demográfica*, también conocida como *paridad estadística* o *independencia*, es uno de los criterios de equidad de grupo más conocidos. Esta noción de equidad afirma que la probabilidad de ser clasificado con el resultado positivo (o negativo) debe ser independiente de que el individuo pertenezca al grupo protegido, es decir, que los datos demográficos de los individuos clasificados positivamente son idénticos a los de la población en su conjunto (Dwork et al. [2011]).

**Definición 65** (Independencia en clasificación binaria). Sean  $C, A$  variables aleatorias. La *independencia* entre  $C$  y  $A$  equivale a que se cumpla la siguiente restricción:

$$P(C = c \mid A = a_r) = P(C = c \mid A = a_o).$$

Lo denotaremos como  $C \perp A$ .

**Definición 66** (Paridad demográfica). Sea  $A \in \mathcal{A}$  un atributo sensible multivaluado y  $g: \mathcal{X} \rightarrow \mathcal{Y}$  un clasificador arbitrario. Se dice que  $g$  cumple con la *paridad demográfica* si, y solo si,  $\hat{Y} \perp A$ .

### Relajaciones y aproximaciones

Podemos *relajar* el concepto de paridad demográfica suponiendo que  $\hat{Y} = 1$  y aproximando su definición con una acotación en valor absoluto de las probabilidades a partir una constante fijada  $\epsilon \in (0, 1]$ , donde  $\epsilon = \frac{p}{100}$  para que se satisfaga la "regla  $p$ " (Zafar et al. [2017b]). De esta manera, aproximamos este criterio de equidad como:

$$|P(\hat{Y} = 1 \mid A = a_r) - P(\hat{Y} = 1 \mid A = a_o)| \leq \epsilon.$$

Tomando un  $\epsilon \in (0, 1]$  también podemos aproximar el concepto de paridad demográfica de la siguiente manera:

$$1 - \epsilon \leq \frac{P(\hat{Y} = 1 \mid A = a_o)}{P(\hat{Y} = 1 \mid A = a_r)}. \quad (5)$$

En algunos trabajos se suele obviar la acotación por una constante y simplemente se define la paridad demográfica utilizando la Definición 66 y asumiendo que  $\hat{Y} = 1$ , lo que sería equivalente a igualar las métricas PPrev entre los subgrupos.

$$\text{PPrev}_{a_r} = \text{PPrev}_{a_o} \Rightarrow P(\hat{Y} = 1 \mid A = a_r) = P(\hat{Y} = 1 \mid A = a_o).$$

*Ejemplo 11.* Consideremos un sistema de detección de delitos y dos grupos de igual tamaño,  $A$  y  $B$ . Suponemos que los miembros del grupo  $B$  tienen el doble probabilidades reales de cometer un delito que los individuos del grupo  $A$ . Al igualar la probabilidad de un resultado positivo, el mismo número de predicciones positivas se distribuiría entre un grupo mucho mayor de delincuentes para  $B$  que para  $A$ . Así, un delincuente del grupo  $B$  tendría menos probabilidades de serlo que un delincuente del grupo  $A$  ( $\text{FNR}_A < \text{FNR}_B$ ). De hecho, para la misma precisión, la tasa de verdaderos positivos del grupo  $B$  sería la mitad de la del grupo  $A$ ,  $\frac{1}{2}\text{TPR}_A = \text{TPR}_B$ , cumpliendo la paridad demográfica.

### 7.4.2 Probabilidades igualadas

Uno de los problemas de la paridad demográfica es que ignora una posible correlación entre  $Y$  y  $A$ . El criterio de las *probabilidades igualadas*, también conocido como *ratio de paridad positiva* o *separación*, tiene en cuenta la etiqueta real de cada grupo y su condicionamiento al resto de variables. Además, proporciona un incentivo para reducir los errores de manera uniforme en todos los grupos sin descartar el clasificador

perfecto (que obtenga  $\hat{Y} = Y$ ) a diferencia de la paridad estadística.

**Definición 67** (Independencia condicional en clasificación binaria). Sean  $C, Y, A$  variables aleatorias. La *independencia condicional* entre  $C$  y  $A$  dado  $Y$  equivale a que se cumpla la siguiente restricción:

$$P(C = c \mid A = a_r, Y = y) = P(C = c \mid A = a_o, Y = y).$$

Lo denotaremos como  $C \perp A \mid Y$ .

El criterio de las probabilidades igualadas establece que  $\hat{Y}$  debe ser condicionalmente independiente de  $A$  dado  $Y$ , permitiendo que el clasificador dependa de  $A$  a través de la variable objetivo (Hardt et al. [2016]). Para utilizar este criterio, es necesario conocer las etiquetas reales de cada individuo, por lo que esta medida restringe su uso para determinadas tareas en las que no conozcamos previamente el resultado de la acción sobre el individuo.

**Definición 68** (Probabilidades igualadas). Sea  $A \in \mathcal{A}$  un atributo sensible multivaluado y  $g: \mathcal{X} \rightarrow \mathcal{Y}$  un clasificador arbitrario. Se dice que  $g$  cumple con el criterio de las *probabilidades igualadas* si, y solo si,  $\hat{Y} \perp A \mid Y$ .

#### Relajaciones y aproximaciones

El concepto previo depende de varias variables por lo que normalmente en la práctica se tiende a relajar el criterio fijando algunos valores en la definición. A partir de estas relajaciones surgen otros criterios de equidad que también son ampliamente utilizados y conocidos en la literatura.

La relajación más común surge al suponer que  $\hat{Y} = 1$ , en este caso, el criterio de las probabilidades igualadas equivale a igualar las métricas FPR y TPR entre los subgrupos. Esta aproximación beneficia al individuo, y equilibra la probabilidad de tener un resultado beneficioso, en todos los subgrupos de los individuos etiquetados tanto positiva como negativamente.

$$\text{FPR}_{a_r} = \text{FPR}_{a_o} \Rightarrow P(\hat{Y} = 1 \mid A = a_r, Y = -1) = P(\hat{Y} = 1 \mid A = a_o, Y = -1).$$

$$\text{TPR}_{a_r} = \text{TPR}_{a_o} \Rightarrow P(\hat{Y} = 1 \mid A = a_r, Y = 1) = P(\hat{Y} = 1 \mid A = a_o, Y = 1).$$

Definiremos el concepto de igualdad de oportunidades como la igualdad de las tasas de verdaderos positivos entre los subgrupos (Hardt et al. [2016]). En algunos trabajos, también se define este criterio de forma equivalente para las tasas de verdaderos negativos.

$$\text{TPR}_{a_r} = \text{TPR}_{a_o} \Rightarrow P(\hat{Y} = 1 \mid A = a_r, Y = 1) = P(\hat{Y} = 1 \mid A = a_o, Y = 1).$$

**Definición 69** (Igualdad de oportunidades). Sea  $A \in \mathcal{A}$  un atributo sensible multi-valuado y  $g: \mathcal{X} \rightarrow \mathcal{Y}$  un clasificador arbitrario. Se dice que  $g$  cumple con la *igualdad de oportunidades* si, y solo si,

$$P(\hat{Y} = 1 \mid A = a_r, Y = 1) = P(\hat{Y} = 1 \mid A = a_o, Y = 1).$$

La igualdad de oportunidades es naturalmente más débil que la Definición 68, ya que asumimos los valores de  $\hat{Y} = 1$  e  $Y = 1$ . Este concepto, iguala la probabilidad de que los individuos etiquetados positivamente sean correctamente clasificados con el resultado positivo (beneficioso). Por ejemplo, dos individuos, un hombre y una mujer, que están cualificados para un trabajo ( $Y = 1$ ), deberían tener la misma probabilidad de conseguir el trabajo ( $\hat{Y} = 1$ ).

*Ejemplo 12.* Consideremos un sistema de contratación y dos grupos de igual tamaño,  $A$  y  $B$ . Imaginemos que en el grupo  $A$  de los 100 aspirantes al cargo, 58 están cualificados, mientras que el grupo  $B$  solo 2 de ellos son aptos para el cargo. Si la empresa decide aceptar a 30 solicitantes y satisfacer la igualdad de oportunidades, se concederán 29 ofertas al grupo  $A$  mientras que solo se concederá 1 oferta al grupo  $B$ . Si el trabajo es bien remunerado, el grupo  $A$  mejorará sus condiciones de vida y a la larga, podrá permitir una mejor educación para sus hijos, y en consecuencia una mejor cualificación de los mismos en el futuro.

En este ejemplo podemos observar que la igualdad de oportunidades no ayuda a cerrar la brecha entre los dos grupos, es más la brecha entre el grupo  $A$  y el grupo  $B$ , tenderá a ampliarse con el tiempo.

#### 7.4.3 Tasa de paridad predictiva

La *tasa de paridad predictiva*, también denominada *suficiencia* surge de una motivación equivalente a la del criterio de las probabilidades igualadas. El concepto se define de igual manera haciendo uso de la independencia condicional, pero cambiando los papeles de  $\hat{Y}$  e  $Y$ .

**Definición 70** (Tasa de paridad predictiva). Sea  $A \in \mathcal{A}$  un atributo sensible multi-valuado y  $g: \mathcal{X} \rightarrow \mathcal{Y}$  un clasificador arbitrario. Se dice que  $g$  cumple con la *tasa de paridad predictiva* si, y solo si,  $Y \perp A \mid \hat{Y}$ .

### Relajaciones y aproximaciones

La relajación más común surge al suponer que  $\hat{Y} = Y$ , en este caso, el criterio de las probabilidades igualadas equivale a igualar las métricas PPV y NPV entre los subgrupos.

$$\begin{aligned} \text{PPV}_{a_r} = \text{PPV}_{a_o} &\Rightarrow P(Y = 1 \mid A = a_r, \hat{Y} = 1) = P(Y = 1 \mid A = a_o, \hat{Y} = 1). \\ \text{NPV}_{a_r} = \text{NPV}_{a_o} &\Rightarrow P(Y = -1 \mid A = a_r, \hat{Y} = -1) = P(Y = -1 \mid A = a_o, \hat{Y} = -1). \end{aligned}$$

Las limitaciones de este concepto de equidad también son similares a las de las probabilidades igualadas, pudiendo acrecentar las diferencias entre los grupos privilegiado y desfavorecido.

#### 7.4.4 Medidas basadas en la puntuación

A diferencia de las definiciones anteriores, que se basan en los índices de clasificación binaria, algunas nociones de equidad se basan en la *puntuación* de la probabilidad predicha  $S$  y la etiqueta real  $Y$  (Mitchell et al. [2021]).

El *balance para la clase positiva* (o *negativa*) se cumple cuando la puntuación esperada para un individuo clasificado positivamente (o negativamente) es igual en todos los grupos.

**Definición 71** (Balance para la clase positiva). Sea  $A \in \mathcal{A}$  un atributo sensible multivaluado y  $S \in [0, 1]$  la puntuación asignada por el clasificador arbitrario  $g: \mathcal{X} \rightarrow \mathcal{Y}$ . Se dice que  $g$  cumple con el *balance para la clase positiva* si, y solo si,

$$\mathbb{E}[S \mid A = a_r, Y = 1] = \mathbb{E}[S \mid A = a_o, Y = 1].$$

**Definición 72** (Balance para la clase negativa). Sea  $A \in \mathcal{A}$  un atributo sensible multivaluado y  $S \in [0, 1]$  la puntuación asignada por el clasificador arbitrario  $g: \mathcal{X} \rightarrow \mathcal{Y}$ . Se dice que  $g$  cumple con el *balance para la clase negativa* si, y solo si,

$$\mathbb{E}[S \mid A = a_r, Y = -1] = \mathbb{E}[S \mid A = a_o, Y = -1].$$

Sin embargo, hay que tener en cuenta que en los problemas del mundo real es casi imposible cumplir el equilibrio simultáneo para la clase negativa y el equilibrio para la clase positiva simultáneamente.

### 7.4.5 Igualdad de las métricas de predicción

Como hemos podido observar, en general, la mayoría de los criterios de equidad definidos surgen como una igualdad de las métricas de grupo presentadas en el Apartado 7.4. De esta forma podemos definir, un nuevo criterio para cada métrica de la siguiente forma.

**Definición 73** (Paridad métrica). Sea  $A \in \mathcal{A}$  un atributo sensible multivaluado y  $g: \mathcal{X} \rightarrow \mathcal{Y}$  un clasificador arbitrario. Se dice que  $g$  cumple con la *paridad métrica* si, y solo si,

$$\text{Métrica}_{a_r} = \text{Métrica}_{a_o}.$$

Donde Métrica podrá ser cualquier métrica de grupo definida previamente.

Aunque esta definición nos permite crear una gran variedad de criterios de equidad, la formulación a partir de una igualdad sigue siendo difícil a la hora de aplicarla en la práctica.

### 7.4.6 Impacto desigual

Para facilitar la implementación práctica de las medidas anteriores, aparece el término *impacto desigual* (no confundir con el término impacto dispar definido en el Cuadro 1) como una aproximación más de la paridad demográfica. Esta noción está directamente relacionada con la "regla  $p$ ", antes mencionada, y que podemos encontrar en la literatura jurídica (Binns [2021]), según la cual una decisión es discriminatoria si el coeficiente del impacto desigual es inferior a un valor  $p$ .

**Definición 74.** (Impacto desigual) Sea  $A \in \mathcal{A}$  un atributo sensible multivaluado,  $g: \mathcal{X} \rightarrow \mathcal{Y}$  un clasificador arbitrario y  $p \in (0, 1]$ . Se dice que  $g$  satisface el *impacto desigual* si, y solo si,

$$\frac{P(\hat{Y} = 1 \mid A = a_o)}{P(\hat{Y} = 1 \mid A = a_r)} \leq p. \quad (6)$$

En el mundo real, si un clasificador aplicado a una tarea en una empresa satisface el impacto desigual, debe justificarse que su existencia es esencial para el funcionamiento seguro y eficiente del negocio, y no existen procedimientos alternativos que sean sustancialmente igual de válidos y tengan un impacto menos adverso. La Comisión para la Igualdad de Oportunidades en el Empleo (EEOC) de Estados Unidos adopta la regla del 80 % ( $p = 0,80$ ) para considerar si una decisión tiene impacto desigual (Adverse Impact Analysis / Four-Fifths Rule).

### Construcción de otras medidas en la práctica

Combinando las acotaciones de las Ecuaciones (6) y (5), obtenemos un criterio de equidad popularmente utilizado, en concreto por el software de Aequitas (Saleiro et al. [2019]). Además esta formalización también puede extenderse a cualquier métrica de grupo de las explicadas en el Apartado 7.4.

Podemos, por ejemplo, definir la tasa de falsas omisiones (FOR) como:

$$\text{Métrica FOR}_G = \frac{\text{FOR}_{a_o}}{\text{FOR}_{a_r}} = \frac{P(Y = 1 \mid A = a_o, \hat{Y} = -1)}{P(Y = 1 \mid A = a_r, \hat{Y} = -1)}.$$

Aequitas utiliza la tasa de positivos predichos (PPR) para aproximar el concepto de paridad demográfica. Además, usando la prevalencia predicha (PPrev), observamos que la métrica generada es equivalente a la Definición 74 y por tanto al concepto de paridad de impacto.

$$\text{Métrica PPrev}_G = \frac{\text{PPrev}_{a_o}}{\text{PPrev}_{a_r}} = \frac{P(\hat{Y} = 1 \mid A = a_o)}{P(\hat{Y} = 1 \mid A = a_r)}.$$

La formulación aportada por Aequitas se basa en calcular la fracción de la métrica de grupo elegida y una vez calculada, comprobar si se encuentra dentro de un rango definido. Tomando  $p = \frac{1}{1-\epsilon}$ , obtenemos la acotación que define el rango. La pertenencia o no al rango es una aproximación a la igualdad de las métricas en la Definición 73.

$$1 - \epsilon \leq \text{Métrica de disparidad}_G \leq \frac{1}{1 - \epsilon}. \quad (7)$$

Usamos  $\epsilon \in (0, 1]$  para controlar el rango de valores de disparidad que pueden considerarse justos. Para aplicar la regla del 80 %, simplemente bastaría con tomar  $1 - \epsilon = 0,8$ . Diremos que un clasificador será tan justo como lo permita el valor máximo del sesgo entre los grupos definidos por  $A$ .

Usaremos unas métricas u otras en función del impacto y el objetivo que quiera intervenir el usuario. Si las intervenciones pueden perjudicar a los individuos (punitivas), entonces queremos minimizar los falsos positivos (centrándonos en la Tasa de Falsos Descubrimientos (FDR) o la Tasa de Falsos Positivos (FPR)). Si por otro lado, tienen como objetivo beneficiar a los individuos (asistenciales), deberíamos preocuparnos más por los falsos negativos (centrándonos en la Tasa de Falsa Omisión (FOR) o la Tasa de Falsos Negativos (FNR)).



---

## ALGORITMOS DE MITIGACIÓN DE SESGO

---

En este capítulo discutiremos los diferentes algoritmos existentes para la mitigación del sesgo, daremos algunos ejemplos específicos de los aportados en la bibliografía y comentaremos sus características más relevantes.

### 8.1 MODELOS DE APRENDIZAJE JUSTOS

En la literatura, podemos encontrar una gran variedad de métodos y algoritmos que nos pueden ayudar a mejorar la equidad en un modelo de aprendizaje. Los enfoques de *mitigación del sesgo* pueden subdividirse en tres categorías: algoritmos de preprocesamiento, que intentan aprender representaciones justas de los datos; algoritmos de optimización durante el entrenamiento, que ajustan el proceso de aprendizaje para cumplir los criterios de justicia; y algoritmos de posprocesamiento, que adaptan las predicciones del modelo en función de sus resultados. Estas categorías no tienen por qué ser mutuamente excluyentes y a veces pueden tener varios métodos de actuación sobre los datos.

A lo largo de este capítulo, presentaremos las características comunes entre los algoritmos de cada uno de los tres grupos presentados y discutiremos algunos ejemplos que podemos encontrar en la literatura, profundizando en algún caso específico de cada tipo. Además, contextualizaremos en su categoría correspondiente el algoritmo de optimización de equidad contrafactual presentado por [Kusner et al. \[2018\]](#) sobre el que basaremos la parte práctica de nuestro trabajo.

#### 8.1.1 Selección de los datos del modelo

En la Sección 6.1, observamos que para construir un modelo de aprendizaje es necesario tener un conjunto de datos  $\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$  a partir del cual poder entrenar el modelo y extraer la información que se quiere aprender. Normalmente para entrenar los modelos, no utilizaremos el conjunto de datos al completo, sino que tomaremos un conjunto de entrenamiento  $X$ , que contiene  $m$  individuos extraídos aleatoriamente del conjunto de datos total, donde  $m < n$ . Cada individuo  $\mathbf{x} \in X$  es un vector de longitud  $d$  donde cada componente del vector describe una caracte-

rística de la persona. Además, cada individuo  $x$  tiene asociado un atributo sensible  $a \in \{0,1\}$ , donde el valor 0 denota la pertenencia al grupo de referencia ( $a_r$ ) y 1 al grupo objetivo ( $a_o$ ). Denotaremos por  $Y$  al conjunto con las etiquetas reales de los individuos que están en  $X$ . De esta forma, tenemos que  $X \times Y \subset \mathcal{D}$ .

En la práctica, se suele utilizar el conjunto de los individuos restantes contenidos en  $\mathcal{D}$  como conjunto de prueba del modelo con el objetivo de poder comprobar su rendimiento sobre un conjunto de datos "nuevo" con el que no ha sido entrenado.

### 8.1.2 Equilibrio entre equidad y métricas de evaluación

Algunas medidas de evaluación como la precisión o exactitud dependen directamente del conjunto de datos, la definición de equidad utilizada y los algoritmos empleados. La equidad en la práctica, perjudica a métricas como la exactitud. Si queremos mitigar el sesgo entre grupos, debemos hacer una compensación entre la equidad y la exactitud, sacrificando en parte esta última como se puede observar en la Figura 14.

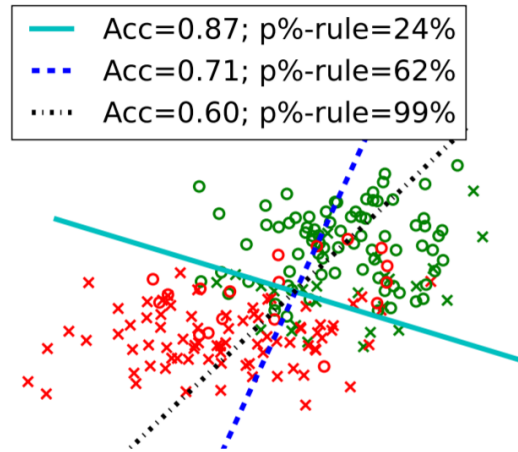


Figura 14: Exactitud vs. Independencia en un problema de clasificación (Zafar et al. [2017b]).

## 8.2 ALGORITMOS DE PREPROCESAMIENTO

Los *algoritmos de preprocesamiento* buscan mejorar la equidad antes de entrenar el modelo, modificando los datos de entrenamiento de forma que no presenten sesgos antes de ser procesados. El propósito de estos algoritmos es aprender una nueva representación  $Z$  que elimine la información correlacionada con el atributo sensible  $A$  y preserve, en la medida de lo posible, la información del conjunto de individuos  $X$  sin necesidad de conocer el valor de sus etiquetas  $Y$ . La tarea posterior (por ejemplo, regresión o clasificación) desempeñada por  $g$ , será independiente del método usado y podrá producir resultados que preserven diversos criterios de equidad.

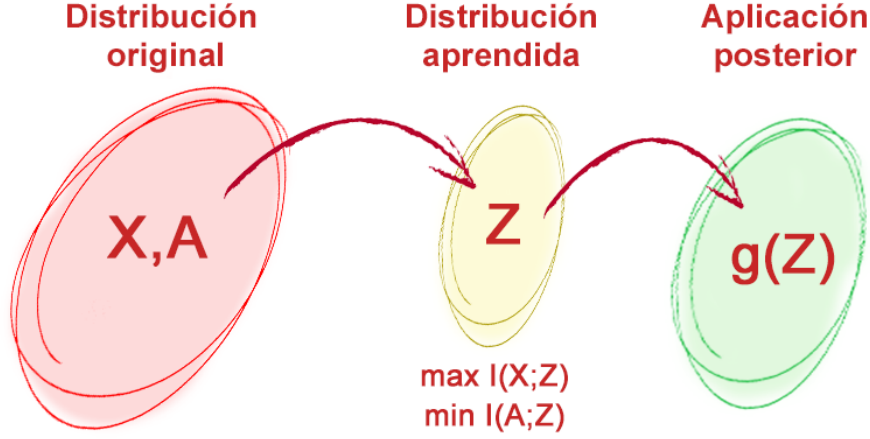


Figura 15: Etapas del proceso de preprocesamiento.

### 8.2.1 Ejemplo: Aprendizaje de la representación justa

Zemel et al. [2013] presenta este problema como el aprendizaje de una función que lleva muestras individuales a representaciones intermedias. Este trabajo tiene dos objetivos: minimizar la pérdida de información (los datos originales deberían preservar la misma cantidad de información) y maximizar la equidad (la pertenencia a un grupo protegido no debería afectar negativamente a los individuos del mismo).

Este ejemplo, podría confundirse como una aplicación de equidad por desconocimiento, pero no es así, ya que la pertenencia a grupos protegidos no es simplemente ignorada, sino que se trata activamente (junto con la información redundante de estos atributos). El modelo propuesto, tiene como objetivo mantener la equidad de grupo e individual, al mismo tiempo que maximiza la exactitud.

**Notación 4.** Introducimos la siguiente notación que utilizaremos en este ejemplo:

- $X^+ = \{\mathbf{x}_n \in X : A = 1\} \subset X$  es el conjunto de datos de entrenamiento cuyos miembros pertenecen al grupo protegido,  $X^- = \{\mathbf{x}_n \in X : A = 0\}$  denota al conjunto cuyos miembros pertenecen al grupo no protegido.
- $Z$  es una variable aleatoria multivariante, donde cada uno de sus  $r$  valores representa un "prototipo". Asociado a cada prototipo existe un vector  $\mathbf{v}_k$  en el mismo espacio que los individuos  $\mathbf{x}_n$ .

La idea es representar cada individuo  $\mathbf{x}_n \in X$  como una combinación lineal ponderada de  $r$  prototipos para satisfacer la paridad demográfica, minimizando la pérdida de información original y maximizando la exactitud en la medida de lo posible.

Se define la probabilidad *softmax* de que un elemento sea un prototipo concreto como:

$$M_{n,k} := P(Z = k \mid \mathbf{x}_n) = \frac{\exp(-d(\mathbf{x}_n, \mathbf{v}_k))}{\sum_{j=1}^r \exp(-d(\mathbf{x}_n, \mathbf{v}_j))} \text{ para todo } n, k.$$

donde  $d$  es una función de medida de distancia (por ejemplo, la distancia  $\ell_2$ ). Los autores definen la regularización  $\ell_2$  como una función de distancia ponderada que viene dada por

$$d(\mathbf{x}_n, \mathbf{v}_k, \alpha) = \sum_{i=1}^d \alpha_i (x_{ni} - v_{ki})^2.$$

Se define  $\hat{y}_n$  como la predicción para  $y_n$  calculada partiendo de la marginalización sobre el valor de  $Y$  para la predicción de cada prototipo

$$\hat{y}_n = \sum_{k=1}^r M_{n,k} w_k.$$

El modelo de aprendizaje, minimiza la siguiente función de pérdida:

$$L = A_z L_z + A_x L_x + A_y L_y.$$

Donde  $L_z$  regulariza la paridad demográfica,  $L_x$  es el error de reconstrucción de la distribución y  $L_y$  cuantifica la pérdida de predicción. Los factores  $A_z, A_x, A_y$  se definen como hiperparámetros para equilibrar estas pérdidas.

$$L_z = \sum_{k=1}^r |M_k^+ - M_k^-|.$$

expresando  $M_k^+ = \frac{1}{|X^+|} \sum_{n \in X^+} M_{n,k}$  y  $M_k^-$  se formula de forma similar a partir de  $X^-$ .

$$L_x = \sum_{n=1}^m (\mathbf{x}_n - \hat{\mathbf{x}}_n)^2.$$

donde  $\hat{\mathbf{x}}_n$  son los nuevos valores de  $\mathbf{x}_n$  para  $Z$ :

$$\hat{\mathbf{x}}_n = \sum_{k=1}^r M_{n,k} \mathbf{v}_k.$$

$$L_y = \sum_{n=1}^m -y_n \log \hat{y}_n - (1 - y_n) \log(1 - \hat{y}_n).$$

En la fase de entrenamiento, los valores de  $\mathbf{v}, \mathbf{w}, \alpha$  se optimizan conjuntamente a través del método L-BFGS (Zhu et al. [1997]) para minimizar la función objetivo  $L$ . Los valores de  $A_x, A_y, A_z$  se seleccionan a través de la afinación de los hiperparámetros, optando por los que producen mejores resultados (búsqueda en malla). Tenemos que tener en cuenta que la función objetivo no es convexa y, por tanto, no garantiza la optimización.

### *Ventajas e inconvenientes*

Algunas de las principales ventajas de los algoritmos de preprocesamiento ya han sido mencionadas previamente. En general, estos métodos, son muy útiles cuando el clasificador que utilizaremos para la fase de entrenamiento es un modelo de caja negra y no conocemos su actuación sobre los datos. En estos casos, al devolvernos una nueva distribución que no contiene correlación con los atributos sensibles, los datos podrán ser usados de forma independiente para cualquier tarea posterior sin necesidad de preocuparnos por la existencia de sesgo entre grupos.

Por otro lado, solo podemos utilizar los métodos de preprocesamiento para optimizar los criterios de equidad que no requieran información sobre el valor de las etiquetas  $Y$  (por ejemplo, la paridad demográfica o la equidad individual). Como desconocemos el uso que se le va a dar a los nuevos datos, en algunos casos, podría no garantizarse la equidad en el modelo final aprendido. Además, este grupo de algoritmos suele ser inferior respecto los otros dos en términos de rendimiento entre exactitud y equidad.

### *Otros ejemplos en la literatura*

El concepto de preprocesamiento optimizado es introducido por [Calmon et al. \[2017\]](#). En este se plantea la reducción de la discriminación como una tarea de optimización convexa con el objetivo de minimizar la pérdida de exactitud (preservando la utilidad) mientras se limita por ciertas medidas de equidad de grupo e individual.

[Wang et al. \[2019\]](#) ofrece métodos que usan distribuciones contrafactuales para resolver el trato dispar de un clasificador de caja negra en una población de despliegue sin la necesidad de reentrenar el modelo. El método propuesto se basa en la construcción de una nueva distribución para los individuos del grupo objetivo de manera que mejoren sus resultados en promedio.

## 8.3 ALGORITMOS DE OPTIMIZACIÓN DURANTE EL ENTRENAMIENTO

Los *algoritmos de optimización* alteran el entrenamiento del propio modelo. En este contexto, la mitigación del sesgo se plantea como un aprendizaje al que se le añade una restricción o un término de regularización al objetivo de optimización existente. Es decir, un modelo aprende a optimizar una función de pérdida en los datos de entrenamiento, sujeta a restricciones de equidad (por ejemplo, la distancia máxima a la paridad demográfica).

Otro enfoque diferente, sería el de la optimización de métricas de rendimiento complejas que incluyan alguna noción de equidad. Por ejemplo, introduciendo una penalización relacionada con la equidad en la función objetivo.

### 8.3.1 Ejemplo: Aprendizaje en clasificación sin impacto dispar.

Uno de los enfoques más populares es el de la optimización con restricciones, donde el objetivo es encontrar un conjunto de parámetros,  $\theta \in \Theta$  que minimicen una función objetivo  $l_0(\theta)$ , sujeta a  $m$  restricciones funcionales,  $l_i(\theta)$  para todo  $i \in \{1, \dots, m\}$ :

$$\theta^* = \arg \min_{\theta \in \Theta} l_0(\theta) \quad \text{donde } l_i(\theta) \leq 0 \quad \text{para todo } i \in \{1, \dots, m\}.$$

Zafar et al. [2017b] enmarca la tarea de clasificación justa imponiendo restricciones lineales a la covarianza entre los atributos sensibles y las predicciones (o la distancia con signo entre el vector de características del individuo y el límite de decisión del clasificador). Este método es adecuado para múltiples atributos sensibles y para cualquier clasificador basado en contornos convexos (por ejemplo, SVM o regresión logística). Además, se propone otra formulación similar, destinada a satisfacer necesidades del mundo real, maximizando la equidad sujeta a restricciones de exactitud.

En el trabajo anterior, se asume que el conjunto de datos original presenta un sesgo histórico, por lo que en Zafar et al. [2017a] se amplía este enfoque a los casos en los que tenemos acceso a la verdad histórica no sesgada en la fase de entrenamiento y podemos saber si una decisión histórica fue correcta o incorrecta.

**Notación 5.** Introducimos la siguiente notación para el ejemplo propuesto:

- $\mathcal{D}'$  es el conjunto de datos de entrenamiento que se define como

$$\mathcal{D}' = \{(\mathbf{x}_i, y_i) \in X \times Y : i = 1, \dots, m\}.$$

- $\theta$  son los parámetros que debemos aprender.
- $L(\theta)$  es la función de pérdida convexa original.
- $d_\theta(\mathbf{x})$  es la función de distancia con signo del vector de características  $\mathbf{x}$  al límite de decisión del clasificador.
- $f_\theta(\mathbf{x})$  es la función de clasificación, definida por

$$f_\theta(\mathbf{x}) = \begin{cases} 1, & \text{si } d_\theta(\mathbf{x}) \geq 0, \\ -1, & \text{en otro caso.} \end{cases}$$

Podemos añadir como restricciones al problema de optimización original, la *paridad OMR* (definida como una relajación para  $\hat{Y} \neq Y$  del criterio de probabilidades igualadas) o la *paridad FNR*. Aunque, para el ejemplo, utilizaremos la *paridad FPR* construida a partir de la Definición 73 y que se define como:

$$P(\hat{Y} = 1 \mid A = 0, Y = -1) = P(\hat{Y} = 1 \mid A = 1, Y = -1).$$

Cabe señalar que la paridad FNR y FPR implican la paridad TPR y TNR respectivamente, y por tanto, la igualdad de oportunidades. En este ejemplo, usaremos como restricción la paridad FPR a partir de la cual surge la siguiente formulación de optimización:

$$\begin{aligned} \text{minimizar: } & L(\theta) \\ \text{sujeto a: } & P(\hat{Y} = 1 \mid A = 0, Y = -1) - P(\hat{Y} = 1 \mid A = 1, Y = -1) \leq \epsilon, \\ & P(\hat{Y} = 1 \mid A = 0, Y = -1) - P(\hat{Y} = 1 \mid A = 1, Y = -1) \geq -\epsilon. \end{aligned} \quad (8)$$

La complejidad de las restricciones, hacen que el problema que plantea minimizar la función  $L(\theta)$  (a priori no convexa) sea intratable a nivel computacional al no poder utilizar los algoritmos tradicionales como el de descenso de gradiente estocástico u otros resolutores para encontrar la solución óptima al problema que se plantea.

Para subsanar esta cuestión, se presentan algunas relajaciones de las restricciones que utilizan la covarianza entre los atributos sensibles de los individuos y  $d_\theta(\mathbf{x})$  para detectar la relación entre el atributo sensible y las predicciones (condicionales) a nivel de grupo:

$$\begin{aligned} \text{Cov}(a, g_\theta(y, \mathbf{x})) &= \mathbb{E}[(a - \bar{a})(g_\theta(y, \mathbf{x}) - \bar{g}_\theta(y, \mathbf{x}))] \\ &\approx \frac{1}{m} \sum_{(\mathbf{x}, y, a) \in \mathcal{D}'} (a - \bar{a})g_\theta(y, \mathbf{x}), \end{aligned}$$

donde el término  $\mathbb{E}[(z - \bar{z})]\bar{g}_\theta(\mathbf{x})$  se anula, ya que  $\mathbb{E}[(z - \bar{z})] = 0$  y la función  $g_\theta(y, \mathbf{x})$  se puede definir como:

$$g_\theta(y, \mathbf{x}) = \min(0, \frac{1-y}{2} y d_\theta(\mathbf{x})).$$

Tras la relajación, podemos reescribir (8) como:

$$\begin{aligned} \text{minimizar: } & L(\theta) \\ \text{sujeto a: } & \frac{1}{m} \sum_{(\mathbf{x}, y, a) \in \mathcal{D}'} (a - \bar{a})g_\theta(y, \mathbf{x}) \leq c, \\ & \frac{1}{m} \sum_{(\mathbf{x}, y, a) \in \mathcal{D}'} (a - \bar{a})g_\theta(y, \mathbf{x}) \geq -c, \end{aligned} \quad (9)$$

donde el umbral de covarianza  $c \in \mathbb{R}^+$  controla el grado de desempeño del criterio de igualdad de oportunidades.

Esta formulación sigue siendo no convexa, por lo que a continuación convertiremos estas restricciones en un *programa convexo-cóncavo disciplinado* (DCCP), que puede resolverse de manera eficiente aprovechando los recientes avances en la programación convexa-cóncava (Shen et al. [2016]).

En primer lugar, consideramos la restricción descrita en (9), es decir:

$$\sum_{(\mathbf{x}, y, a) \in \mathcal{D}'} (a - \bar{a}) g_{\theta}(y, \mathbf{x}) \sim c,$$

donde  $\sim$ , podría denotar ' $\leq$ ' o ' $\geq$ '. Además, dejamos de lado la constante  $\frac{1}{n}$  para simplificar. Como  $a \in \{0, 1\}$ , dividimos la suma en la expresión anterior en dos términos:

$$\sum_{(\mathbf{x}, y) \in \mathcal{D}'_0} (0 - \bar{a}) g_{\theta}(y, \mathbf{x}) + \sum_{(\mathbf{x}, y) \in \mathcal{D}'_1} (1 - \bar{a}) g_{\theta}(y, \mathbf{x}) \sim c, \quad (10)$$

donde  $\mathcal{D}'_0$  y  $\mathcal{D}'_1$  son subconjuntos del conjunto de datos  $\mathcal{D}'$  que toman valores  $a = 0$  y  $a = 1$ , respectivamente. Definimos  $m_0 = |\mathcal{D}'_0|$  y  $m_1 = |\mathcal{D}'_1|$ , entonces  $\bar{z} = \frac{m_1}{m}$  y podemos reescribir (10) como:

$$-\frac{m_1}{m} \sum_{(\mathbf{x}, y) \in \mathcal{D}'_0} g_{\theta}(y, \mathbf{x}) + \frac{m_0}{m} \sum_{(\mathbf{x}, y) \in \mathcal{D}'_1} g_{\theta}(y, \mathbf{x}) \sim c,$$

que, dado que  $g_{\theta}$  es convexa en  $\theta$  (por suposición), resulta en una función convexa-cóncava.

Por lo tanto, podemos reescribir el problema definido por (9) como:

$$\begin{aligned} &\text{minimizar: } L(\theta) \\ &\text{sujeto a: } -\frac{m_1}{m} \sum_{(\mathbf{x}, y) \in \mathcal{D}'_0} g_{\theta}(y, \mathbf{x}) + \frac{m_0}{m} \sum_{(\mathbf{x}, y) \in \mathcal{D}'_1} g_{\theta}(y, \mathbf{x}) \leq c, \\ &\quad -\frac{m_1}{m} \sum_{(\mathbf{x}, y) \in \mathcal{D}'_0} g_{\theta}(y, \mathbf{x}) + \frac{m_0}{m} \sum_{(\mathbf{x}, y) \in \mathcal{D}'_1} g_{\theta}(y, \mathbf{x}) \geq -c, \end{aligned}$$

que es un DCCP para cualquier función de pérdida convexa  $L(\theta)$ , y puede ser resuelto eficientemente usando heurísticas bien conocidas como la propuestas en [Shen et al. \[2016\]](#).

#### *Ventajas e inconvenientes*

Entre las principales ventajas que tiene este grupo de algoritmos, se encuentra que al poder modificar el modelo y optimizar las métricas utilizadas, consiguen un mejor rendimiento en términos de las medidas de equidad y exactitud, variando en función del algoritmo utilizado.

En cambio, como estos métodos modifican directamente el proceso de aprendizaje, a menudo son difíciles de generalizar a diferentes modelos o métricas. Además, en el ámbito de los problemas de *machine learning* en el mundo real, no siempre tendremos acceso al modelo de clasificación, por lo que esto supone un gran inconveniente en la optimización del mismo.



### *Otros ejemplos en la literatura*

La equidad contrafactual presentada por [Kusner et al. \[2018\]](#), se basa en la noción de que un resultado debería ser el mismo independientemente del grupo demográfico del individuo. Definiendo una decisión como justa si se mantiene igual cuando se cambia el valor del atributo protegido. En el artículo, se plantea como objetivo la mitigación del sesgo como un problema de optimización de equidad usando como base la inferencia causal. Este enfoque funciona bien para capturar los sesgos sociales e identificar el equilibrio entre equidad y utilidad.

[Russell et al. \[2017\]](#) presenta un artículo donde extiende el trabajo de [Kusner et al. \[2018\]](#) aportando un método para ofrecer predicciones justas con respecto a varios modelos causales posibles simultáneamente.

## 8.4 ALGORITMOS DE POSPROCESAMIENTO

Los *algoritmos de posprocesamiento* tienen como objetivo ajustar un clasificador ya entrenado para que cumpla con unas restricciones de equidad específicas. Esto se suele hacer mediante la calibración de umbrales, cuya idea principal es encontrar un umbral adecuado utilizando una función de puntuación para cada grupo.

### 8.4.1 Ejemplo: Aprendizaje en igualdad de oportunidades

[Hardt et al. \[2016\]](#) desarrolla un marco para eliminar la discriminación de forma óptima para cualquier modelo de clasificación aprendido. Los autores definen la equidad como una restricción del concepto de probabilidades igualadas. Esta técnica, conocida como *predictor derivado*, utiliza la calibración del umbral para navegar por la curva ROC hasta que se cumplan los criterios de equidad establecidos.

Se utilizan diferentes valores de umbral para los distintos subgrupos, y se buscan soluciones factibles a lo largo de la intersección de cada intervalo convexo de la curva ROC correspondiente. Al definir el rendimiento del modelo como el rendimiento mínimo en cualquier subgrupo, el objetivo de equidad coincidirá con el objetivo del modelo. Dado un clasificador y las correspondientes curvas ROC para ambos grupos. Podemos encontrar el umbral basado en las curvas, como podemos observar en la Figura 16.

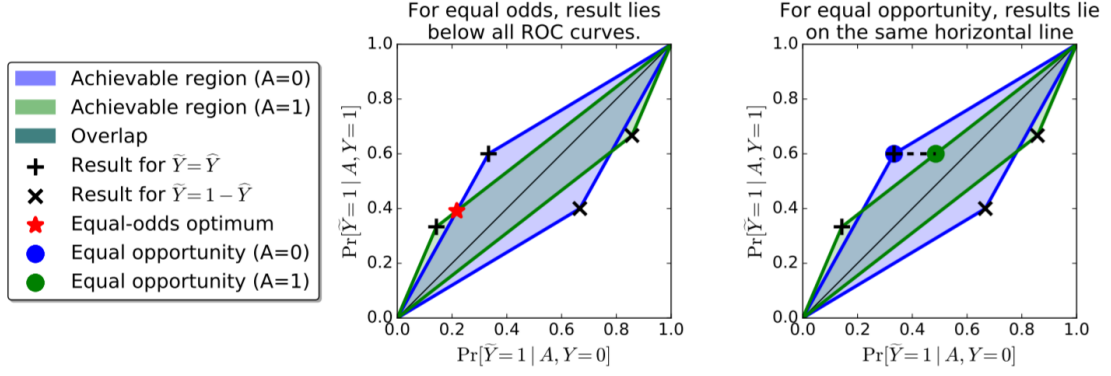


Figura 16: Búsqueda del clasificador óptimo para los criterios de probabilidades igualadas e igualdad de oportunidad, respectivamente. (Hardt et al. [2016])

El criterio de probabilidades igualadas se satisface únicamente cuando las curvas ROC de los dos grupos se cruzan, como se muestra en la imagen izquierda de la Figura 16; la igualdad de oportunidades, como relajación de la noción anterior, puede satisfacerse tomando un umbral tal que las tasas de verdaderos positivos de los dos grupos sean iguales, como se puede ver en la imagen de la derecha.

#### *Ventajas e inconvenientes*

Los algoritmos de posprocesamiento comparten algunas ventajas con los de preprocesamiento, pudiéndose aplicar de forma independiente al modelo de clasificación sin necesidad de modificar su método de actuación. Además, consigue un rendimiento relativamente bueno en la optimización de la mayoría de definiciones de equidad (excepto la equidad contrafactual).

Sin embargo, se puede argumentar que al actuar sobre el modelo después de haberlo aprendido, este proceso es intrínsecamente subóptimo. Siendo equivalente a aprender a sabiendas un modelo sesgado y luego corregirlo, en lugar de aprender un modelo insesgado desde el principio.

#### *Otros ejemplos en la literatura*

En el artículo presentado por Woodworth et al. [2017] se amplía el trabajo de Hardt et al. [2016] demostrando que su método de posprocesamiento podría ser subóptimo en algunos casos. Se ofrece una demostración de que el problema es intratable y se presenta un método nuevo que aproxima el criterio de equidad proporcionando un resultado estadísticamente cercano al óptimo.

## Parte III

### FUNDAMENTOS DE LA EQUIDAD CONTRAFACTUAL

Discusión sobre la inferencia causal, el cálculo de contrafactuales y el teorema de incompatibilidad como fundamentos matemáticos de la equidad contrafactual.

---

## INFERENCIA CAUSAL

---

En este capítulo formalizaremos algunos conceptos básicos para desarrollar la teoría relativa a la causalidad en el ámbito de la equidad. Además, introduciremos los grafos como herramienta para describir los modelos causales explicados, así como los efectos que tienen estos modelos en las poblaciones donde se aplican.

### 9.1 MODELOS CAUSALES

Elegiremos los *modelos causales estructurales* aprovechando que pueden ofrecernos una base sólida para las diferentes nociones causales utilizadas en este trabajo. La forma más sencilla de conceptualizar un modelo causal estructural es como un programa que genera una distribución a partir de *variables de ruido* independientes mediante una secuencia de instrucciones formales.

Imaginemos que en lugar de muestras de una distribución, tenemos un programa informático que genera muestras a partir de una semilla aleatoria. El código de este programa, partiría de una semilla aleatoria simple e iría construyendo muestras cada vez más complejas. Esta idea es la misma que utiliza un modelo causal estructural cambiando la sintaxis de programación por lenguaje matemático.

#### 9.1.1 Ejemplo: Construcción de un modelo causal

Supongamos una población en la que un individuo hace ejercicio regularmente con una probabilidad de  $\frac{1}{2}$ . Con una probabilidad de  $\frac{1}{3}$ , el individuo tiene predisposición a desarrollar sobrepeso en ausencia de ejercicio regular. Del mismo modo, en ausencia de ejercicio, la aparición de una enfermedad cardíaca puede aparecer con una probabilidad de  $\frac{1}{3}$ . Denotaremos por  $X$  el indicador de ejercicio regular, por  $Y$  el de exceso de peso, y por  $Z$  el indicador de la enfermedad cardíaca. A continuación, construiremos un modelo causal estructural para generar muestras de esta población hipotética (Barocas et al. [2019]).

---

**Algoritmo 1:** Programa distribución causal 1.

---

Muestras de variables aleatorias independientes de Bernoulli:

 $U_1 \sim \text{Bernoulli}\left(\frac{1}{2}\right), U_2, U_3 \sim \text{Bernoulli}\left(\frac{1}{3}\right);$  $X := U_1;$  $Y := \text{if } X = 1 \text{ then } 0 \text{ else } U_2;$  $Z := \text{if } X = 1 \text{ then } 0 \text{ else } U_3;$ 

---

A partir de la descripción anterior, observamos que en nuestra población el ejercicio evita tanto el sobrepeso como las enfermedades cardíacas, pero en ausencia de ejercicio, ambos son independientes. Nuestro programa genera una distribución conjunta sobre las variables aleatorias  $(X, Y, Z)$ . Podemos calcular las probabilidades bajo esta distribución. Por ejemplo, la probabilidad de sufrir una enfermedad cardíaca bajo la distribución especificada por nuestro modelo es de  $\frac{1}{2} \cdot \frac{1}{3} = \frac{1}{6}$ . También podemos calcular la probabilidad condicional de padecer enfermedades cardíacas dado el sobrepeso. Dado el suceso  $Y = 1$  podemos inferir que el individuo no hace ejercicio  $X = 0$ , por lo que la probabilidad de sufrir una enfermedad cardíaca debido al sobrepeso es  $\frac{1}{3}$ .

Formalmente, tener un programa que genere una distribución es más potente que el simple acceso al muestreo. Una de las razones es que podemos manipular el programa de la manera que queramos, mientras resulte en un programa funcional. Podríamos, por ejemplo, establecer  $Y := 1$ , dando lugar a una nueva distribución. El programa resultante tiene el siguiente aspecto:

---

**Algoritmo 2:** Programa distribución causal 2.

---

Muestras de variables aleatorias independientes de Bernoulli:

 $U_1 \sim \text{Bernoulli}\left(\frac{1}{2}\right), U_2, U_3 \sim \text{Bernoulli}\left(\frac{1}{3}\right);$  $X := U_1;$  $Y := 1;$  $Z := \text{if } X = 1 \text{ then } 0 \text{ else } U_3;$ 

---

Calculando de nuevo la probabilidad de sufrir una enfermedad cardíaca sobre la nueva distribución, de nuevo obtenemos  $\frac{1}{6}$ . Este cálculo revela una idea importante, la sustitución  $Y := 1$  no corresponde a un condicionamiento de  $Y = 1$ . Una se trata de una acción y la otra es una observación de la que podemos extraer conclusiones. En este ejemplo, si observamos que un individuo tiene sobrepeso, podemos inferir que tiene un mayor riesgo de enfermedad cardíaca. Sin embargo, esto no significa que la reducción del peso corporal evite las enfermedades cardíacas. En cambio, la intervención  $Y := 1$  crea un nuevo modelo en el que todos los individuos de la población tienen sobrepeso con todo lo que ello conlleva.

A continuación, profundizaremos un poco más en este punto considerando otra población hipotética, especificada por el siguiente programa:

**Algoritmo 3:** Programa distribución causal 3.

---

Muestras de variables aleatorias independientes de Bernoulli:

 $U_1 \sim \text{Bernoulli}\left(\frac{1}{2}\right), U_2, U_3 \sim \text{Bernoulli}\left(\frac{1}{3}\right);$ 
 $Y := U_2;$ 
 $X := \text{if } Y = 0 \text{ then } 0 \text{ else } U_1;$ 
 $Z := \text{if } X = 1 \text{ then } 0 \text{ else } U_3;$ 


---

En esta población, la única razón por la que los individuos eligen hacer ejercicio con cierta probabilidad es el sobrepeso. Por otro lado, las enfermedades del corazón se desarrollan en ausencia de ejercicio. La sustitución  $Y := 1$  en este modelo conduce a un aumento de la probabilidad de hacer ejercicio y por tanto, a una disminución de la probabilidad de sufrir una enfermedad cardíaca. El condicionamiento de  $Y = 1$  también tiene el mismo efecto y en ambos casos, la probabilidad de sufrir un problema cardíaco es de  $\frac{1}{6}$ .

### 9.1.2 Formalización de los modelos causales estructurales

Los modelos causales estructurales nos proporcionan un cálculo preciso para razonar sobre el efecto de las acciones hipotéticas. Formalmente, un modelo causal estructural es una secuencia de asignaciones que generan una distribución conjunta a partir de variables de ruido independientes. A continuación ofreceremos la definición de modelo causal estructural presentada por Pearl [2000].

**Definición 75** (Modelo causal estructural). Un *modelo causal estructural*  $M$  se define como una tupla  $(U, V, F)$  de conjuntos tales que:

- $U$  es un conjunto de variables aleatorias de ruido, las cuales deben ser conjuntamente independientes. Corresponden a factores no causados por ninguna variable del conjunto  $V$  de variables observadas.
- $F$  es un conjunto de funciones  $\{f_1, \dots, f_n\}$ , una para cada  $V_i \in V$ , tal que,

$$V_i = f_i(pa_i, U_i), \text{ para todo } i = 1, \dots, n.$$

donde  $pa_i \subseteq V \setminus \{V_i\}$  y  $U_i \subseteq U$ . Estas ecuaciones también son conocidas como ecuaciones estructurales.

El modelo es causal en el sentido de que, dada una distribución de probabilidad  $P(U)$  sobre las variables de ruido  $U$ , podemos derivar la distribución de un subconjunto  $W \subseteq V$  tras una intervención en  $V \setminus W$ . Cuando  $M$  denota un modelo causal estructural, escribiremos la probabilidad de un evento  $E$  bajo la distribución conjunta vinculada como  $P_M(E)$ .

Para familiarizarnos con la notación, supongamos que  $M$  denota el modelo causal estructural del Apartado 9.1.1, entonces la probabilidad de sufrir una enfermedad cardíaca en este modelo será  $P_M(Z) = \frac{1}{6}$ .

## 9.2 GRAFOS CAUSALES

La notación  $pa_i$  utilizada en la Definición 75 se refiere al subconjunto de variables  $pa(V_i)$  que contiene los padres del nodo  $V_i$ . Esta notación viene motivada por la suposición de que el modelo se factoriza como un grafo dirigido, el cual en este trabajo, restringiremos al caso acíclico (DAG). A este grafo lo llamaremos: el *grafo o diagrama causal* correspondiente al modelo causal estructural especificado.

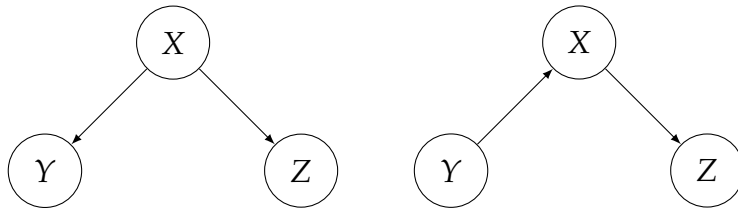


Figura 17: Grafos causales de los modelos descritos por los Programas 1 y 3, respectivamente.

Los diagramas causales se utilizan cuando las asignaciones exactas en un modelo causal estructural son secundarias y lo que es realmente relevante son los caminos presentes y ausentes entre nodos.

Los grafos también nos permiten aprovechar el lenguaje establecido de la teoría de grafos para discutir nociones causales. En particular, los grafos causales nos permiten distinguir la causa y el efecto (de tipo directo o indirecto) en función de si un nodo es ancestro o descendiente de otro.

### 9.2.1 Forks

**Definición 76 (Fork).** Sea  $G$  un grafo acíclico dirigido,  $U$  el camino entre dos nodos  $y$  y  $A \in U$ . Llamaremos *fork* al nodo  $A$  si  $(A, B) \in E$ , para todo  $B \in ve(A) \cap U$ .

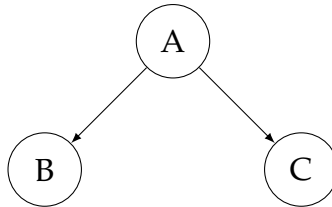


Figura 18: Ejemplo de *fork*.

En la Figura 18, el nodo  $A$  es un ejemplo de *fork* o dicho de otro modo, el nodo  $A$  es una causa común de los nodos  $B$  y  $C$ . En el grafo causal resultante de la distribución del Programa 1, el nodo  $X$  también es un ejemplo de *fork* ( $Y \leftarrow X \rightarrow Z$ ). En ese caso, la variable indicadora de ejercicio regular  $X$  influía tanto en el aumento de peso  $Y$  como en el riesgo de enfermedad  $Z$ , pero como ya discutimos en el Apartado 9.1.1, las variables  $Y$  y  $Z$  no están correlacionadas positivamente. Llegamos a la conclusión de que el nodo *fork*, tiene un efecto de confusión que conduce a un desacuerdo entre el cálculo de las probabilidades condicionales y las intervenciones.

*Ejemplo 13.* En un conocido estudio médico, un presunto efecto beneficioso de la terapia de sustitución hormonal para reducir las enfermedades cardiovasculares desapareció tras identificar el estatus socioeconómico como variable de confusión (Humphrey et al. [2002]). Los ejemplos de confusión suponen una amenaza para la validez de las conclusiones extraídas de los datos en problemas del mundo real.

### 9.2.2 Colliders

**Definición 77 (Collider).** Sea  $G$  un grafo acíclico dirigido,  $U$  el camino entre dos nodos y  $A \in U$ . Llamaremos *collider* al nodo  $A$  si  $(B, A) \in E$ , para todo  $B \in ve(A) \cap U$ .

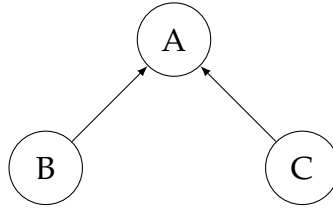


Figura 19: Ejemplo de *collider*.

En la Figura 19, el nodo  $A$  es un ejemplo de *collider*. Cabe destacar que los *colliders* no dan lugar a situaciones en las que se pueda dar confusión. De hecho, en la figura anterior, la relación entre  $B$  y  $C$  no es confusa, lo que significa que podemos sustituir las intervenciones por probabilidades condicionales. Sin embargo, condicionar un *collider*, podría crear una correlación entre  $B$  y  $C$ , un fenómeno al que denominaremos sesgo de *collider*.

*Ejemplo 14.* En el ámbito sanitario, dos enfermedades independientes pueden correlacionarse negativamente cuando se analizan pacientes hospitalizados. La razón es que cuando cualquiera de las dos enfermedades ( $B$  o  $C$ ) es suficiente para el ingreso en el hospital (indicado por la variable  $A$ ), observar que un paciente tiene una enfermedad hace que la otra sea estadísticamente menos probable. A esto es lo que se le conoce como paradoja de Berkson (Berkson [2014]).



### 9.2.3 Mediador

En la definición de *fork*, no tenemos una relación directa entre los nodos  $B$  y  $C$ . Si queremos un efecto total de  $B$  sobre  $C$  estableceremos esta relación causal a través de  $A$ . En este caso,  $A$  no será un factor confusión y recibirá el nombre de *mediador*.

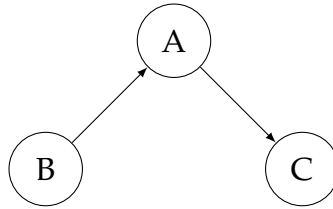


Figura 20: Ejemplo de mediador.

En el grafo causal asociado a la distribución generada por el Programa 3, el nodo  $X$  es un ejemplo de mediador ( $Y \rightarrow X \rightarrow Z$ ). La noción de mediador es especialmente relevante para el tema del análisis de la discriminación ya que establece una relación directa entre las diferentes variables que definen a un individuo. Esta relación causal nos servirá como herramienta para extraer conclusiones sobre las causas de segregación entre grupos.

## 9.3 INTERVENCIÓN Y CONFUSIÓN

Los modelos causales estructurales nos proporcionan una herramienta de formalización del efecto de acciones e intervenciones sobre la población donde se aplican. Como hemos visto previamente, para modelar estos efectos simplemente necesitamos la capacidad de realizar sustituciones.

### 9.3.1 Operadores para realizar actuaciones en el modelo

A partir de los ejemplos propuestos en el Apartado 9.1.1, hemos observado que fijar una variable por sustitución puede corresponder o no a una probabilidad condicional. Esto refuerza nuestra intuición de que una observación no es una acción. En cambio, una sustitución sí es una acción, ya que al sustituir un valor estamos rompiendo el curso natural de la acción captada por nuestro modelo.

**Definición 78** (Intervención). Dado un modelo causal estructural  $M$ , se define una *intervención* sobre una variable observada  $X$  como la sustitución de la ecuación  $X := f(pa, U)$  por la ecuación  $X := x$  para un valor  $x$  constante.

Denotaremos el modelo resultante por  $M' = M[X := x]$  para indicar la modificación que realizamos sobre el modelo original  $M$ . Bajo esta asignación mantenemos  $X$  constante eliminando la influencia de sus nodos padres y por tanto, de cualquier otra variable del modelo. Por otra parte, los nodos hijos de  $X$  recibirán un valor constante  $x$  cuando consulten el valor de su padre.

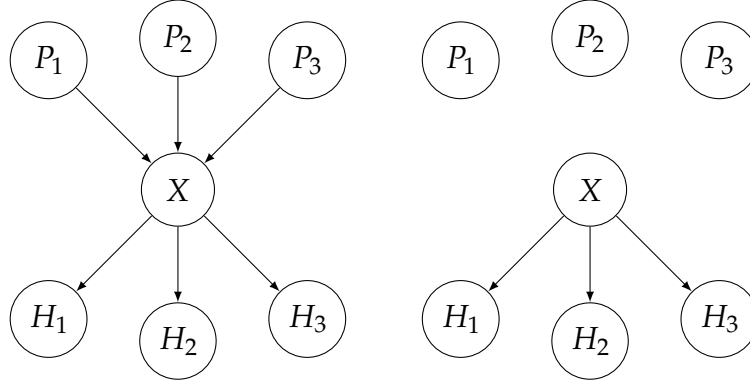


Figura 21: Grafo causal antes y después de la sustitución.

El operador de asignación también se denomina *operador do* para destacar que corresponde a la realización de una acción o intervención. La notación que usaremos para calcular las probabilidades para un evento cualquiera  $E$  después de aplicar el operador *do*, será  $P_{M[X:=x]}(E)$ . También podemos utilizar otra notación que aproxima el concepto de probabilidad condicional y es equivalente al anterior y que se define como:

$$P(E \mid \text{do}(X := x)) = P_{M[X:=x]}(E).$$

### 9.3.2 Confusión entre dos variables

Las cuestiones importantes en inferencia causal están relacionadas con cuándo podemos reescribir una operación *do* en términos de probabilidades condicionales. Cuando esto sea posible, podremos estimar el efecto de la operación *do* a partir de las probabilidades condicionales estimadas a partir de los datos.

Sea una variable  $Y$  sobre la que actúa una variable  $X$ , nos interesará que exista una equivalencia entre el efecto causal de la acción y la probabilidad condicional correspondiente, es decir, que se cumpla la siguiente igualdad:

$$P(Y = y \mid \text{do}(X := x)) = P(Y = y \mid X = x).$$

En general, esto no es cierto. Al fin y al cabo, la diferencia entre la observación (probabilidad condicional) y la acción (intervención) es la principal motivación de la inferencia causal.

**Definición 79** (Confusión). Sean  $Y$  una variable aleatoria sobre la que actúa otra variable  $X$ , diremos que son *confusas* si, y solo si,

$$P(Y = y \mid \text{do}(X := x)) \neq P(Y = y \mid X = x).$$

Cuando tenemos dos variables aleatorias confusas, podemos estimar el efecto de una intervención en términos de probabilidades condicionales a partir de la denominada fórmula de ajuste.

**Proposición 18.** Sean  $X, Y$  dos variables confusas, podemos aproximar el efecto causal de una intervención dada a partir de probabilidades condicionales como:

$$P(Y = y \mid \text{do}(X := x)) = \sum_z P(Y = y \mid X = x, PA = z)P(PA = z),$$

donde  $PA$  indica el conjunto  $pa(X)$ .

Dependiendo de la estructura del grafo podremos eliminar o no la confusión entre dos variables utilizando la fórmula de ajuste sobre un nodo u otro. Si el grafo tiene una estructura de *fork* (por ejemplo,  $B \leftarrow A \rightarrow C$ ), eliminaremos la confusión entre los nodos  $B$  y  $C$ , condicionando  $A$ . En cualquier otro caso (mediador o *collider*), ajustar una variable tendría consecuencias opuestas a las que buscamos.

#### *Criterio de backdoor*

El tratamiento de la confusión a partir de la fórmula de ajuste, puede ser una tarea complicada cuando la cantidad de nodos en el grafo aumente considerablemente. Para detectar las variables sobre las que deberemos condicionar, aparece el *criterio de backdoor* (Pearl [2000]). Este método parte de la idea de seleccionar un conjunto de variables que "bloqueen" todos los *caminos de backdoor* entre los dos nodos sobre los que queremos eliminar la confusión.

**Definición 80** (Camino de *backdoor*). Un *camino de backdoor* entre dos nodos  $A$  y  $B$  es cualquier camino que empiece con una arista de la forma " $\leftarrow$ " hacia  $A$ .

**Definición 81** (Conjunto de *backdoor*). Un *conjunto de backdoor* es una secuencia de variables o nodos contenida en un camino de *backdoor*.

Para aplicar el criterio de *backdoor*, primero seleccionaremos un conjunto de *backdoor* de nodos del grafo. Si el conjunto está formado por una secuencia de nodos relacionados únicamente por aristas de tipo " $\rightarrow$ ", podremos eliminar la confusión entre las variables aplicando la fórmula de ajuste sobre un nodo central de la cadena. Por otro lado, si el camino contiene un *collider* o un descendiente de este, la confusión es inevi-

table, ya que "bloqueando" el camino podríamos impedir que la información fluyera a través de los nodos.

*Ejemplo 15.* Sea un grafo causal dado por la secuencia  $A \leftarrow C \rightarrow D \rightarrow E \rightarrow B$ , nuestro objetivo será eliminar la confusión entre las variables  $A$  y  $B$ . Es evidente que la cadena  $A \leftarrow C \rightarrow D \rightarrow E \rightarrow B$  es un camino de *backdoor*. A continuación, seleccionamos un conjunto de *backdoor* entre ambos nodos, por ejemplo  $C \rightarrow D \rightarrow E$ . En vista de la forma de la secuencia anterior, podemos eliminar la confusión entre  $A$  y  $B$  aplicando la fórmula de ajuste sobre el nodo  $D$ .

### Confusión no observada

La fórmula de ajuste presentada en la Proposición 18, podría sugerir que siempre podemos eliminar el sesgo de confusión condicionando a los nodos padres. Sin embargo, esto no se cumple cuando aparecen *factores de confusión no observados*. En la práctica, a menudo hay variables que son difíciles de medir o que no fueron registradas. Podemos incluir estos nodos no observados en un grafo indicando su influencia con líneas discontinuas.

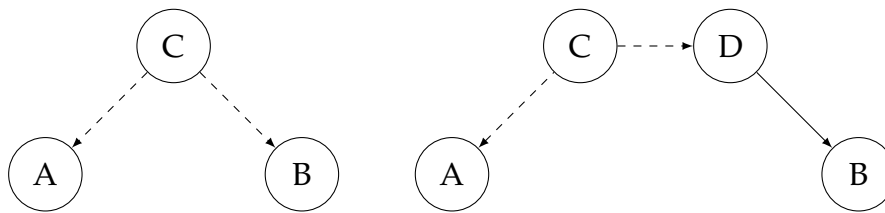


Figura 22: Ejemplos de confusión no observada.

La Figura 22 muestra dos casos de confusión no observada. En el primer ejemplo, el efecto causal de  $A$  sobre  $B$  no es identificable. En el segundo caso, podemos eliminar la confusión entre  $A$  y  $B$  a partir del criterio de *backdoor*. Sea  $C \rightarrow D \rightarrow B$  un conjunto de *backdoor*, podemos eliminar la confusión entre las variables  $A$  y  $B$  ajustando la fórmula sobre la variable  $D$  aunque  $C$  no se observe.

Cabe destacar que podemos combatir la confusión no observada aumentando el número de variables consideradas, pero esto aumentaría progresivamente la complejidad de nuestro modelo causal. En la práctica, es habitual controlar el mayor número posible de variables con el objetivo de eliminar el sesgo de confusión. Sin embargo, como hemos visto, el control de mediadores y *colliders* podría ser problemático en la resolución de nuestro problema.

---

## TEOREMA DE IMPOSIBILIDAD DE LA EQUITAD

---

En este capítulo se propondrá una demostración del teorema de imposibilidad de equidad. Este enunciado surge como una formalización de la *incompatibilidad* entre los criterios de paridad demográfica, probabilidades igualadas y tasa de paridad predictiva.

### 10.1 CARACTERIZACIÓN DEL TEOREMA

La mayoría de los criterios de equidad definidos en el Capítulo 7 se construyen a partir de restricciones no triviales de la distribución de probabilidad conjunta. Por ello, es lógico pensar que la imposición de varios de ellos de forma simultánea, restringirían el espacio de búsqueda hasta el punto de que solo obtendríamos soluciones degeneradas.

El Teorema de la Imposibilidad, cuya primera aproximación fue ofrecida por [Kleinberg et al. \[2016\]](#), establece que no se puede satisfacer más de una medida de equidad al mismo tiempo para un clasificador bien entrenado y un atributo sensible que sea capaz de introducir un sesgo en el modelo. En nuestro caso, presentaremos una versión del Teorema de la Imposibilidad para tres de los criterios de equidad de grupo estudiados: paridad demográfica, tasa de paridad predictiva y probabilidades igualadas. Los enunciados de los lemas demostrados a lo largo de este capítulo, han sido definidos en base al trabajo de [Barocas et al. \[2019\]](#).

#### 10.1.1 Paridad demográfica versus Tasa de paridad predictiva

Comenzamos con un lema que muestra cómo, en general, la paridad demográfica y la paridad predictiva se excluyen mutuamente. La única suposición necesaria es que el atributo sensible  $A$  y la variable  $Y$  no son independientes, es decir, dependen una de la otra. Esto es una forma diferente de decir que un grupo tiene mayor tasa de resultados positivos que otro, lo que es cierto en la mayoría de casos.

**Lema 1.** Supongamos que  $A$  e  $Y$  son variables dependientes. Entonces la paridad demográfica ( $\hat{Y} \perp A$ ) y la tasa de paridad predictiva ( $Y \perp A \mid \hat{Y}$ ) no pueden verificarse simultáneamente.

*Demostración.* El enunciado del lema es análogo a la siguiente expresión,

$$Y \not\perp A \implies \neg(\hat{Y} \perp A \wedge Y \perp A \mid \hat{Y}).$$

Procederemos por contrarrecíproco, lo que equivale a demostrar que,

$$\hat{Y} \perp A \wedge Y \perp A \mid \hat{Y} \implies Y \perp A.$$

Si se da la independencia entre las variables  $A$  e  $\hat{Y}$ , entonces se cumple,

$$P(A = a, \hat{Y} = \hat{y}) = P(A = a)P(\hat{Y} = \hat{y}). \quad (11)$$

Por otro lado, sabemos que la independencia condicional dada por  $Y \perp A \mid \hat{Y}$  satisface que,

$$P(Y = y, A = a \mid \hat{Y} = \hat{y}) = P(Y = y \mid \hat{Y} = \hat{y})P(A = a \mid \hat{Y} = \hat{y}). \quad (12)$$

Finalmente, aplicando las hipótesis del enunciado y usando el teorema de probabilidad total (T.P.T) sobre  $P(Y = y, A = a)$ , llegamos a la siguiente expresión:

$$\begin{aligned} P(Y = y, A = a) &\stackrel{\text{T.P.T}}{=} \sum_{\hat{y}} P(\hat{Y} = \hat{y})P(Y = y, A = a \mid \hat{Y} = \hat{y}) \\ &\stackrel{(12)}{=} \sum_{\hat{y}} P(\hat{Y} = \hat{y})P(Y = y \mid \hat{Y} = \hat{y})P(A = a \mid \hat{Y} = \hat{y}) \\ &= \sum_{\hat{y}} P(\hat{Y} = \hat{y})P(Y = y \mid \hat{Y} = \hat{y}) \frac{P(A = a, \hat{Y} = \hat{y})}{P(\hat{Y} = \hat{y})} \\ &\stackrel{(11)}{=} \sum_{\hat{y}} P(\hat{Y} = \hat{y})P(Y = y \mid \hat{Y} = \hat{y}) \frac{P(A = a)P(\hat{Y} = \hat{y})}{P(\hat{Y} = \hat{y})} \\ &= \sum_{\hat{y}} P(\hat{Y} = \hat{y})P(Y = y \mid \hat{Y} = \hat{y})P(A = a) \\ &= P(A = a) \sum_{\hat{y}} P(\hat{Y} = \hat{y})P(Y = y \mid \hat{Y} = \hat{y}) \\ &\stackrel{\text{T.P.T}}{=} P(A = a)P(Y = y). \end{aligned}$$

La última igualdad nos da que las variables  $A$  e  $Y$  son independientes y por tanto tenemos que  $Y \perp A$ .  $\square$

## 10.1.2 Paridad demográfica versus Probabilidades igualadas

Un resultado análogo de exclusión mutua es válido para la paridad demográfica y el criterio de probabilidades igualadas. El enunciado, en este caso, es un poco más rebuscado y requiere la suposición adicional de que la variable  $Y$  sea binaria. También necesitamos que la variable  $\hat{Y}$  dependa de  $Y$ . Esta suposición es una relajación bastante suave, ya que cualquier función de clasificación útil tiene correlación con la variable  $Y$ .

**Lema 2.** Supongamos que  $Y$  es una variable binaria,  $A$  e  $Y$  son dependientes y además,  $Y$  también depende de  $\hat{Y}$ . Entonces la paridad demográfica ( $\hat{Y} \perp A$ ) y el criterio de las probabilidades igualadas ( $\hat{Y} \perp A \mid Y$ ) no pueden verificarse simultáneamente.

*Demostración.* El enunciado del lema es equivalente a la siguiente expresión,

$$Y \not\perp A \wedge Y \not\perp \hat{Y} \implies \neg(\hat{Y} \perp A \wedge \hat{Y} \perp A \mid Y).$$

Por el contrarrecíproco, deberemos demostrar que,

$$\hat{Y} \perp A \wedge \hat{Y} \perp A \mid Y \implies Y \perp A \vee Y \perp \hat{Y}.$$

Si se da la independencia entre las variables  $A$  e  $\hat{Y}$ , entonces se cumple,

$$P(\hat{Y} = \hat{y}, A = a) = P(\hat{Y} = \hat{y})P(A = a). \quad (13)$$

Sabemos que la independencia condicional dada por  $\hat{Y} \perp A \mid Y$  satisface que,

$$P(\hat{Y} = \hat{y} \mid A = a, Y = y) = P(\hat{Y} = \hat{y} \mid Y = y). \quad (14)$$

Aplicando el teorema de la probabilidad total (T.P.T) y la hipótesis de independencia condicional sobre  $P(\hat{Y} = \hat{y} \mid A = a)$ , obtenemos la siguiente expresión:

$$\begin{aligned} P(\hat{Y} = \hat{y} \mid A = a) &\stackrel{\text{T.P.T}}{=} \sum_y P(Y = y \mid A = a)P(\hat{Y} = \hat{y} \mid A = a, Y = y) \\ &\stackrel{(14)}{=} \sum_y P(Y = y \mid A = a)P(\hat{Y} = \hat{y} \mid Y = y). \end{aligned} \quad (15)$$

Usando la hipótesis de independencia entre las variables  $A$  e  $\hat{Y}$ ,

$$\begin{aligned} P(\hat{Y} = \hat{y} \mid A = a) &= \frac{P(\hat{Y} = \hat{y}, A = a)}{P(A = a)} \\ &\stackrel{(13)}{=} \frac{P(\hat{Y} = \hat{y})P(A = a)}{P(A = a)} \\ &= P(\hat{Y} = \hat{y}). \end{aligned} \quad (16)$$

Combinando las Ecuaciones (15) y (16), llegamos a la siguiente expresión:

$$P(\hat{Y} = \hat{y}) = \sum_y P(Y = y | A = a)P(\hat{Y} = \hat{y} | Y = y). \quad (17)$$

Por otro lado, aplicando el teorema de la probabilidad total sobre  $P(\hat{Y} = \hat{y})$ , tenemos que,

$$P(\hat{Y} = \hat{y}) = \sum_y P(Y = y)P(\hat{Y} = \hat{y} | Y = y). \quad (18)$$

Combinando las Ecuaciones (17) y (18), conseguimos la expresión dada por:

$$\sum_y P(Y = y | A = a)P(\hat{Y} = \hat{y} | Y = y) = \sum_y P(Y = y)P(\hat{Y} = \hat{y} | Y = y). \quad (19)$$

A continuación, y para que sea más cómodo de manipular la expresión anterior definiremos la siguiente notación:

$$\begin{aligned} p &= P(Y = 0), \\ p_a &= P(Y = 0 | A = a), \\ \hat{y}_y &= P(\hat{Y} = \hat{y} | Y = y). \end{aligned}$$

Por hipótesis del Lema,  $Y$  es una variable binaria (supongamos que puede tomar los valores 0 o 1) y por tanto, podemos reescribir la Ecuación (19) como:

$$p\hat{y}_0 + (1 - p)\hat{y}_1 = p_a\hat{y}_0 + (1 - p_a)\hat{y}_1.$$

Simplificando en la ecuación anterior, tenemos que,

$$p(\hat{y}_0 - \hat{y}_1) = p_a(\hat{y}_0 - \hat{y}_1),$$

lo cual es equivalente a,

$$(p - p_a)(\hat{y}_0 - \hat{y}_1) = 0. \quad (20)$$

La igualdad de la Ecuación (20), se satisface si se da alguno de los siguientes casos:

- $p = p_a$ , que es equivalente a  $P(Y = 0) = P(Y = 0 | A = a)$ , y por tanto tenemos,

$$\begin{aligned} P(Y = 1) &= 1 - P(Y = 0) \\ &= 1 - P(Y = 0 | A = a) \\ &= P(Y = 1 | A = a), \end{aligned}$$

donde la expresión anterior equivale a que  $Y \perp A$ .

- $\hat{y}_0 = \hat{y}_1$ , que equivale a  $P(\hat{Y} = \hat{y} | Y = 0) = P(\hat{Y} = \hat{y} | Y = 1)$ , y por tanto tenemos que  $\hat{Y} \perp Y$ .

□



## 10.1.3 Probabilidades igualadas versus Tasa de paridad predictiva

Por último, pasamos a la relación entre la tasa de paridad predictiva y el criterio de probabilidades igualadas. Ambas exigen una relación de independencia condicional no trivial entre las tres variables  $A$ ,  $\hat{Y}$  e  $Y$ . Imponer ambas simultáneamente conduce a un espacio de soluciones degenerado, como confirma el corolario siguiente.

**Corolario 2.** Supongamos que todos los sucesos en la distribución conjunta  $(A, \hat{Y}, Y)$  tienen probabilidad positiva y además,  $A$  depende de  $Y$ . Entonces el criterio de probabilidades igualadas ( $\hat{Y} \perp A \mid Y$ ) y la tasa de paridad predictiva ( $Y \perp A \mid \hat{Y}$ ) no pueden verificarse simultáneamente.

*Demostración.* El enunciado del corolario es análogo a la siguiente expresión,

$$Y \not\perp A \implies \neg(\hat{Y} \perp A \mid Y \wedge Y \perp A \mid \hat{Y}).$$

Por el contrarrecíproco, tenemos que probar que,

$$\hat{Y} \perp A \mid Y \wedge Y \perp A \mid \hat{Y} \implies Y \perp A.$$

Por la propiedad simétrica de la independencia condicional, tenemos que,

$$\begin{aligned} \hat{Y} \perp A \mid Y &= A \perp \hat{Y} \mid Y, \\ Y \perp A \mid \hat{Y} &= A \perp Y \mid \hat{Y}. \end{aligned}$$

La hipótesis dada por  $A \perp \hat{Y} \mid Y$  satisface que,

$$P(A = a \mid \hat{Y} = \hat{y}, Y = y) = P(A = a \mid Y = y). \quad (21)$$

Por otro lado, sabemos que  $A \perp Y \mid \hat{Y}$  cumple que,

$$P(A = a \mid Y = y, \hat{Y} = \hat{y}) = P(A = a \mid \hat{Y} = \hat{y}). \quad (22)$$

Para que tengan sentido las Ecuaciones (21) y (22), es necesaria la hipótesis de que  $P(\hat{Y} = \hat{y}, Y = y) > 0$  que a su vez, es consecuencia directa de que la distribución conjunta  $(A, \hat{Y}, Y)$  tenga probabilidad positiva.

Por la propiedad de unión débil de la independencia condicional, tenemos que las Ecuaciones (21) y (22) son equivalentes, y por tanto:

$$P(A = a \mid Y = y) = P(A = a \mid \hat{Y} = \hat{y}). \quad (23)$$

Usando la definición de probabilidad condicionada sobre  $P(A = a \mid Y = y)$ , tenemos que,

$$P(A = a \mid Y = y) = \frac{P(A = a, Y = y)}{P(Y = y)}, \quad (24)$$

queremos demostrar que  $Y \perp A$ , es decir,  $P(A = a, Y = y) = P(A = a)P(Y = y)$ , que aplicado a la Ecuación (24), equivale a probar:

$$P(A = a \mid Y = y) = P(A = a).$$

Usando el teorema de la probabilidad total (T.P.T) sobre  $P(A)$ , conseguimos la siguiente expresión:

$$\begin{aligned} P(A = a) &\stackrel{\text{T.P.T}}{=} \sum_{\hat{y}} P(\hat{Y} = \hat{y})P(A = a \mid \hat{Y} = \hat{y}) \\ &\stackrel{(23)}{=} \sum_{\hat{y}} P(\hat{Y} = \hat{y})P(A = a \mid Y = y) \\ &= P(A = a \mid Y = y) \sum_{\hat{y}} P(\hat{Y} = \hat{y}) \\ &= P(A = a \mid Y = y). \end{aligned}$$

□

Para un objetivo binario, la hipótesis de no degeneración del Corolario 2 establece que en todos los grupos, para todos los valores de predicción, tenemos casos positivos y negativos. En caso de que el clasificador sea binario, podemos debilitar la suposición exigiendo que el clasificador haga al menos una predicción falsa positiva. Lo atractivo de la afirmación resultante es que su prueba se basa esencialmente en una relación bastante popular entre la tasa de verdaderos positivos (*Recall*) y el valor predictivo positivo (*Precision*).

**Lema 3.** Supongamos que  $\hat{Y}$  es una variable binaria que toma valores de un clasificador con una tasa de falsos positivos no nula y además,  $A$  dependiente de  $Y$ . Entonces el criterio de probabilidades igualadas ( $\hat{Y} \perp A \mid Y$ ) y la tasa de paridad predictiva ( $Y \perp A \mid \hat{Y}$ ) no pueden verificarse simultáneamente.

*Demostración.* Dado que  $Y \not\perp A$ , existirán dos grupos, a los que llamaremos  $p_0$  y  $p_1$  cumpliendo que  $p_0 \neq p_1$ , donde  $p_a = P(Y = 1 \mid A = a)$ .

Supondremos que se satisface el criterio de las probabilidades igualadas. Por hipótesis el clasificador tendrá la misma tasa para todos los grupos de falsos positivos  $FPR > 0$  y de verdaderos positivos  $TPR > 0$ . Procederemos a demostrar que la tasa de paridad predictiva no se satisface en estas condiciones.

En el caso binario, la tasa de paridad predictiva implica que todos los grupos tienen el mismo de PPV (ver Apartado 7.4.3). El valor predictivo positivo en el grupo  $a$ , denotado  $PPV_a$  satisface

$$PPV_a = \frac{TPR p_a}{TPR p_a + FPR(1 - p_a)}.$$

En la expresión anterior podemos ver que  $PPV_0 = PPV_1$  si, y solo si,  $TPR = 0$  o  $FPR = 0$ . Descartaremos esto último por hipótesis. Por tanto, se debe cumplir que  $TPR = 0$ . Sin embargo, podemos deducir que  $NPV_0 \neq NPV_1$  a partir de la siguiente expresión,

$$NPV_a = \frac{(1 - FPR)(1 - p_a)}{(1 - TPR)p_a + (1 - FPR)(1 - p_a)}.$$

Por tanto, la tasa de paridad predictiva no se satisface.  $\square$

#### 10.1.4 Enunciado y demostración

Una vez demostrados los resultados previos, estamos preparados para enunciar y demostrar la versión del teorema de imposibilidad para la equidad de grupo.

**Teorema 5** (Teorema de imposibilidad de la equidad). Consideremos un problema de clasificación binaria con una tasa de falsos positivos no nula donde se cumple la siguiente relación de dependencia entre las variables:  $Y \not\perp A$  e  $\hat{Y} \not\perp Y$ . Si existe una asignación de riesgo, entonces ésta no puede satisfacer los criterios de paridad demográfica, probabilidades igualadas y paridad predictiva simultáneamente dos a dos.

*Demostración.* Bajo las hipótesis del teorema, aplicamos los Lemas 1, 2 y 3.  $\square$

Podemos consultar una versión diferente de la demostración del Teorema 5 desde la perspectiva de la inferencia causal en el artículo propuesto por Saravanakumar [2021].

Añadir si hay tiempo los criterios de mejora que propone Saravakumar

---

## MEDIDAS CAUSALES

---

En este capítulo realizaremos un estudio del modelo contrafactual y cómo sirve de base para diferentes medidas causales, en particular para la equidad contrafactual que discutiremos más detalladamente.

### 11.1 CONTRAFACTUALES

Una vez definidos los modelos causales estructurales, podemos formular preguntas más delicadas que el mero efecto de una acción. En concreto, preguntas contrafactuales como: ¿Habría evitado el atasco si hubiese tomado otra ruta diferente? o ¿Me habrían concedido el préstamo si mi raza o edad fuesen distintas? Podemos dar respuesta a estas preguntas a partir de un modelo causal estructural. Sin embargo, el procedimiento para extraer la respuesta del modelo necesita del cálculo de *contrafactuales*.

#### 11.1.1 Ejemplo: Modelo de decisión contrafactual

Supongamos un problema de decisión entre dos modelos de caja negra para resolver un problema. Denotaremos por  $X$  a la variable indicadora de cada algoritmo. Si el problema es irresoluble ( $U = 1$ ) ninguno de los algoritmos podrá encontrar una solución. Si el problema es resoluble ( $U = 0$ ), un algoritmo obtendrá mejores resultados que el otro. El rendimiento es independientemente en cualquiera de los dos modelos con una probabilidad de  $\frac{1}{2}$ . Definiremos dos variables aleatorias  $U_0, U_1$  que nos informarán del rendimiento para los algoritmos  $X = 0$  y  $X = 1$ , respectivamente. El modelo seleccionado entre los dos existentes será elegido al azar por una variable  $U_X$  con probabilidad  $\frac{1}{2}$ . Supondremos también una variable  $Y \in \{0, 1\}$  que nos dirá si el algoritmo ha encontrado una solución óptima ( $Y = 0$ ) o no ( $Y = 1$ ). A continuación especificaremos el modelo discutido con el siguiente programa:

**Algoritmo 4:** Programa distribución contrafactual.

---

Muestras de variables aleatorias independientes de Bernoulli:

$$U, U_0, U_1, U_X \sim \text{Bernoulli}\left(\frac{1}{2}\right);$$

$$X := U_X;$$

$$Y := X \cdot \max\{U, U_1\} + (1 - X) \cdot \max\{U, U_0\};$$


---

El grafo asociado al modelo anterior viene dado por:

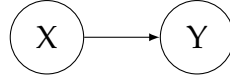


Figura 23: Grafo causal asociado al ejemplo.

Supongamos que elegimos el algoritmo  $X = 1$  y observamos que no consigue una solución  $Y = 1$ . A continuación, nos hacemos la siguiente pregunta: ¿Habría sido mejor elegir el otro algoritmo? Para responder a ello, calcularemos la probabilidad  $P_{M[X:=0]}(Y = 0)$ . Dada la sustitución  $X := 0$  en nuestro modelo, para que el algoritmo encontrase una solución óptima necesitamos que  $\max\{U, U_0\} = 0$ . Esto sólo ocurre cuando  $U = 0$  (con probabilidad  $\frac{1}{2}$ ) y  $U_0 = 0$  (también con probabilidad  $\frac{1}{2}$ ). Concluimos que  $P_{M[X:=0]}(Y = 0) = \frac{1}{2} \cdot \frac{1}{2} = \frac{1}{4}$ .

Aunque pudiera parecerlo, esta no es la respuesta correcta a nuestra pregunta. La razón es que hicimos los cálculos sin considerar la decisión previa  $\{X = 1, Y = 1\}$ . A partir de esta observación ciertas configuraciones de las variables de ruido ( $U, U_0, U_1$ ) ya no son factibles. En concreto si  $U$  y  $U_1$  hubieran sido ambos cero, habríamos encontrado una solución correcta con el algoritmo  $X = 1$ , pero esto es contrario a nuestra observación  $Y = 1$ . De hecho, la configuración  $\{X = 1, Y = 1\}$  sólo permite los siguientes valores para  $U$  y  $U_1$ :

$U$	$U_1$
0	1
1	0
1	1

Cuadro 4: Posibles valores de las variables de ruido dada la evidencia observada.

Cada uno de estos tres casos es igualmente probable, lo que en particular significa que el evento  $U = 1$  tiene una probabilidad de  $\frac{2}{3}$ . Sin considerar la observación, recordemos que  $U = 1$  tenía una probabilidad de  $\frac{1}{2}$ . Esto significa que la evidencia observada  $\{X = 1, Y = 1\}$  ha sesgado la distribución de la variable de ruido  $U$  hacia el valor 1. Utilizaremos la letra  $U'$  para referirnos a esta versión sesgada de  $U$ .

Podemos volver a considerar el efecto de la acción  $X := 0$  sobre el resultado  $Y$  trabajando ahora con la nueva variable  $U'$ . Para  $Y = 0$  necesitamos que  $\max\{U', U_0\} = 0$ .

Esto significa que  $U' = 0$ , un suceso que ahora tiene probabilidad  $\frac{1}{3}$  y  $U_0 = 0$  (con probabilidad  $\frac{1}{2}$ , igual que antes). Por lo tanto, obtenemos que una vez actualizado el modelo,  $P_{M'[X:=0]}(Y = 0) = \frac{1}{3} \cdot \frac{1}{2} = \frac{1}{6}$ .

Si comparamos el nuevo valor con el resultado anterior, tenemos que la incorporación de las observaciones disponibles en nuestro cálculo disminuyó la probabilidad de que encontrásemos la solución óptima con el otro algoritmo,

$$P_{M'[X:=0]}(Y = 0) = \frac{1}{6} < \frac{1}{4} = P_{M[X:=0]}(Y = 0).$$

La razón de ello es que el evento observado sesga la distribución de las variables de ruido haciendo que fuese incluso más probable que en general, el problema no tuviese solución. Llamaremos al resultado que acabamos de calcular: el valor contrafactual al elegir el algoritmo alternativo dado que el algoritmo seleccionado no encontró una solución óptima.

### 11.1.2 Formalización del cálculo contrafactual

Sea  $M = (U, V, F)$  un modelo causal estructural, la especificación de  $F$  es un supuesto fuerte que permite el cálculo de valores contrafactuales. Si queremos obtener el valor de  $Y$  si  $Z$  hubiera tomado valor  $z$ , para dos variables observadas  $Z$  e  $Y$ . El valor contrafactual se modela como la solución de  $Y$  para un  $U := u$  dado en el que las ecuaciones de  $Z$  se sustituyen por  $Z := z$ . Denotaremos la expresión anterior por  $Y_{Z \leftarrow z}(U)$  (Pearl [2000]).

La *inferencia contrafactual*, especificada por un modelo causal  $M$  dada la evidencia  $E$ , equivale al cálculo de  $P(Y_{Z \leftarrow z}(U) \mid E = e)$ , donde  $E, Z, Y \subseteq V$ . Existen tres pasos esenciales en el cálculo de contrafactuales: En primer lugar, incorporamos las observaciones sesgando las variables de ruido mediante una operación de condicionamiento. En segundo lugar, realizamos una operación *do* en el modelo causal después de sustituir las variables de ruido sesgadas. Finalmente, calculamos la distribución para una variable objetivo. Estos tres pasos suelen denominarse *abducción*, *acción* y *predicción* y se describen de la siguiente manera:

**Definición 82** (Cálculo del contrafactual). Dado un modelo causal estructural  $M$ , un evento observado  $E$ , una intervención  $Z := z$  y una variable objetivo  $Y$ , definimos el *contrafactual*  $Y_{Z \leftarrow z}(E)$  mediante los siguientes pasos:

- *Abducción*: Ajustar las variables de ruido al evento observado. Formalmente, condicionar la distribución conjunta de  $U = (U_1, \dots, U_n)$  al suceso  $E$ . Esto da lugar a una distribución sesgada  $U' = P(U \mid E = e)$ .
- *Acción*: Utilizar el operador *do* para realizar la intervención  $Z := z$  en el modelo causal estructural  $M$  obteniendo el modelo  $M' = M[Z := z]$ .
- *Predicción*: Calcular el objetivo contrafactual  $Y_{Z \leftarrow z}(E)$  usando  $U'$  como semilla aleatoria en  $M'$ .

## 11.2 EQUIDAD CONTRAFACTUAL

Todos los criterios de equidad definidos en el Capítulo 7 tienen limitaciones en su aplicación a problemas reales. Mientras que la equidad por desconocimiento es un criterio insuficiente debido a la gran cantidad de características correlacionadas con los atributos sensibles, la equidad individual tiene el problema de ser directamente dependiente de una distancia fiable entre individuos. Por otro lado, los criterios de equidad de grupo son observacionales y no pueden utilizarse para encontrar la causa de la disparidad entre grupos.

Como solución a estos problemas, aparecen las *medidas causales* donde, en este trabajo, profundizaremos en la *equidad contrafactual* (Russell et al. [2017]) que puede ser considerada como una subclase de las mismas. Este concepto considera que para un individuo, una decisión es justa si coincide en el mundo real y en un mundo "contrafactual" en el que el individuo perteneciese a un grupo demográfico diferente. Esta suposición construye un método para comprobar el tratamiento dispar que surge al sustituir únicamente el atributo sensible y además, proporciona una explicación del impacto del sesgo a través de un grafo causal.

**Definición 83** (Equidad contrafactual). Dado  $A \in \mathcal{A}$  un atributo sensible multivaluado,  $g: \mathcal{X} \rightarrow \mathcal{Y}$  un clasificador arbitrario,  $X \in \mathcal{X} \setminus \mathcal{A}$  una variable observada cualquiera y  $(U, V, F)$  un modelo causal donde  $V \equiv A \cup X$ . Se dice que  $g$  satisface la *equidad contrafactual* si, y solo si,

$$P(\hat{Y}_{A \leftarrow a}(U) = y \mid X = x, A = a) = P(\hat{Y}_{A \leftarrow a'}(U) = y \mid X = x, A = a),$$

para todo  $y$  y  $a' \neq a$ .

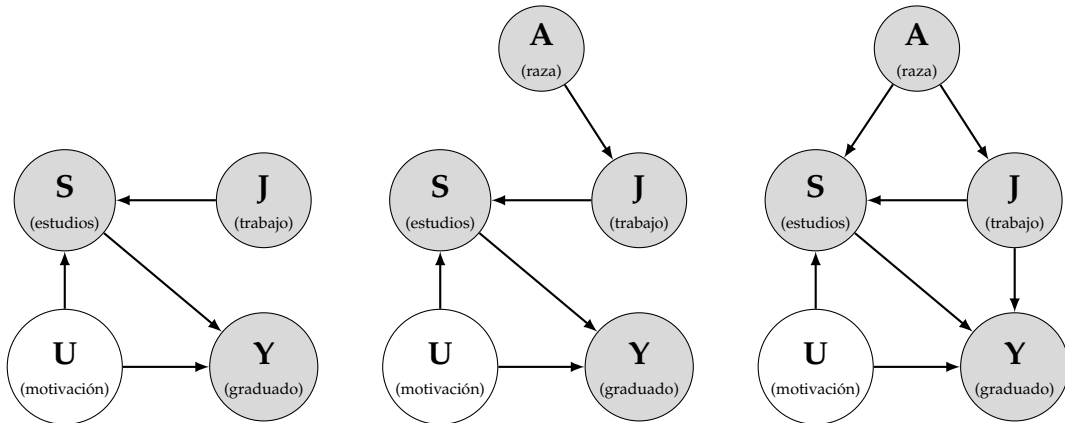


Figura 24: Ejemplo de grafo causal en un problema real.

La Figura 24 muestra varios ejemplos sobre un caso relacionado con la admisión a la universidad (Russell et al. [2017]). Si sustituimos el atributo sensible (raza) por su valor contrafactual, todas las características correlacionadas con él también se verían influidas (estudios y trabajo), propagándose hacia abajo en el grafo causal a través de las ecuaciones estructurales. Cualquier atributo que no descienda del atributo sensible permanecerá igual.

### 11.2.1 Implicaciones de la definición de equidad

En este apartado, presentaremos algunas implicaciones de la definición de equidad contrafactual y algunos resultados que proporcionan una forma directa de satisfacer esta noción de justicia para un modelo dado (Kusner et al. [2018]).

**Lema 4.** Sea  $G$  el grafo causal del modelo dado por  $(U, V, F)$ ,  $A \in \mathcal{A}$  un atributo sensible multivaluado y  $g: \mathcal{X} \rightarrow \mathcal{Y}$  un clasificador arbitrario. Entonces  $g$  satisface la equidad contrafactual si es una función que no depende de los nodos descendientes de  $A$ .

*Demostración.* Sea  $W$  una variable no descendiente de  $A$  en  $G$ . Entonces  $W_{A \leftarrow a}(U)$  y  $W_{A \leftarrow a'}(U)$  tienen la misma distribución para los tres pasos de la Definición 82. Por lo tanto, la distribución de cualquier función  $g$  de los nodos no descendientes de  $A$  es invariante con respecto a los valores contrafactuales de  $A$ .  $\square$

### Solución a problemas de otras nociones de equidad

La equidad contrafactual también proporciona una respuesta a algunos problemas sobre la incompatibilidad de los criterios de equidad. Supongamos el caso en el que nos gustaría que nuestro clasificador cumpliera el criterio de igualdad de oportunidades y la tasa de paridad predictiva simultáneamente. Se demostró en el Capítulo 10 que esto es imposible. La equidad contrafactual nos aporta una solución en este escenario, sugiriendo que tanto la igualdad de oportunidades como la paridad predictiva pueden ser insuficientes si  $A$  e  $Y$  están asociados: suponiendo que las variables no son confusas, esto es el resultado de que  $A$  sea una causa de  $Y$ .

Como estamos en el ámbito de la equidad contrafactual, no deberíamos utilizar  $Y$  como base para nuestras decisiones, en lugar de ello, deberíamos buscar una función de variables que no sean causadas por  $A$  pero que puedan predecir  $Y$ , a este conjunto de funciones lo denotaremos por  $Y_{\perp A}$ . Definiremos entonces  $\hat{Y}$  como una estimación de la  $Y_{\perp A}$  más "cercana" a  $Y$  usando como guía alguna función de riesgo. Esto hace que la incompatibilidad entre el criterio de igualdad de oportunidades y paridad predictiva sea irrelevante, ya que  $A$  e  $Y_{\perp A}$  serán independientes por como las hemos construido dadas las suposiciones del modelo.



### *Tratamiento de los prejuicios históricos y la paradoja de la equidad existente*

La diferencia explícita entre  $\hat{Y}$  e  $Y$  nos permite abordar los sesgos históricos. Por ejemplo, supongamos que  $Y$  es un indicador de si un cliente no devuelve un préstamo, mientras que  $\hat{Y}$  es la decisión real de conceder el préstamo. Consideremos el grafo causal que se muestra en la Figura 25 con la inclusión explícita del conjunto  $U$  de variables de ruido independientes. En principio,  $Y$  es la medida objetivamente ideal para la toma de decisiones, ya que indica si el cliente dejó de pagar o no el préstamo. Si  $A$  es un atributo protegido, entonces el clasificador definido por  $\hat{Y} = Y = f_Y(A, U)$  no satisface la equidad contrafactual, siendo la flecha  $A \rightarrow Y$  el resultado de un mundo que perjudica a los individuos de una manera que está fuera de su control. Por tanto, el principio de equidad contrafactual nos prohíbe utilizar  $Y$ .

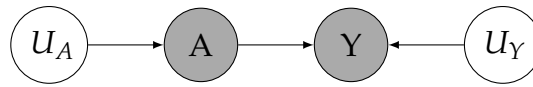


Figura 25: Grafo causal donde  $A$  es un atributo sensible e  $Y$  un resultado de interés.

Por el contrario, cualquier función de variables no descendientes de  $A$  puede utilizarse como base para la toma de decisiones justas. Esto significa que cualquier variable  $\hat{Y}$  definida por  $\hat{Y} = g(U)$  satisface la equidad contrafactual para cualquier función  $g(\cdot)$ . Por lo tanto, dado un modelo causal, el funcional definido por  $g(\cdot)$  que minimiza algún error de predicción para  $Y$  cumplirá la equidad contrafactual, como se propone en la Sección 12.1. En esencia, estamos aprendiendo una proyección de  $Y$  en el espacio de las decisiones justas, eliminando de paso los posibles sesgos históricos existentes.

### *Defectos y limitaciones*

El concepto de equidad contrafactual, teóricamente hablando, puede parecer una buena idea para eliminar todos los defectos que surgen del resto de criterios estudiados. En la práctica, la dependencia de conceptos como la inferencia causal o el estudio de contrafactuales lo hacen una de las nociones de equidad más complejas de implementar.

Otro problema surge cuando queremos acordar cómo debería ser el grafo causal o decidir qué características vamos a utilizar incluso disponiendo de dicho grafo. Esto se debe a que podríamos perder precisión si rechazamos variables de estudio que pudiesen ser problemáticas a la hora de realizar el estudio causal.

## Parte IV

### ANÁLISIS EXPERIMENTAL

Elaboración de un ejemplo práctico sobre la equidad contrafactual y su discusión frente a otras nociones de equidad estudiadas.

---

## DESCRIPCIÓN Y DISEÑO

---

En este capítulo se describirá un problema real de justicia en el ámbito de la educación y se propondrán varios diseños de modelos causales que puedan ser tratados por un algoritmo de equidad contrafactual con el objetivo de eliminar el tratamiento dispar en el problema.

### 12.1 ALGORITMO DE APRENDIZAJE JUSTO

El algoritmo propuesto por [Kusner et al. \[2018\]](#) parte de la necesidad de relacionar  $\hat{Y}$  con  $Y$ . Para ello restringiremos  $\hat{Y}$  para que funcione como una función parametrizada de los nodos no descendientes de  $A$  apoyándonos en el Lema 4.

Calculemos la predicción  $\hat{Y}$  a partir de un clasificador parametrizado por  $\theta$  al que denominaremos como  $g_\theta(U, X_{\neq A})$ , donde  $X_{\neq A} \subseteq X$  denota el conjunto de no descendientes de  $A$ . Dada una función de pérdida  $l(\cdot, \cdot)$  (*squared* o *logistic loss*) y un conjunto de datos  $\mathcal{D} = \{(A^{(i)}, X^{(i)}, Y^{(i)}) : i = 1, \dots, n\}$ . Definimos  $L(\theta) = \frac{1}{n} \sum_{i=1}^n \mathbb{E}[l(y^{(i')}, g_\theta(U^{(i')}, x_{\neq A}^{(i')})) | x^{(i)}, a^{(i)}]$  como la pérdida empírica al minimizar sobre  $\theta$ . Cada esperanza es respecto a la variable  $U^{(i)} \sim P_{\mathcal{M}}(U | x^{(i)}, a^{(i)})$  donde  $P_{\mathcal{M}}(U | x, a)$  es la distribución condicional de las variables de ruido dada por el modelo causal  $\mathcal{M}$ . En tiempo de predicción, crearemos una nueva variable  $\tilde{Y} = \mathbb{E}[\hat{Y}(U^*, x_{\neq A}^*) | x^*, a^*]$  para un nuevo elemento del conjunto de datos  $(a^*, x^*)$ .

Para calcular la esperanza, utilizaremos el *método de Monte Carlo basado en cadenas de Markov* (MCMC) para aproximarla ([Andrieu et al. \[2003\]](#)). Los métodos MCMC son una clase de algoritmos de simulación para el muestreo y estimación de distribuciones de probabilidad a posteriori. Al construir una cadena de Markov usando una distribución deseada como distribución de equilibrio de la cadena, se puede hacer un muestreo de la distribución registrando los diferentes estados del grafo. Cuantas más iteraciones realicemos con el método de Monte Carlo, la distribución de la cadena, se acercará más a la distribución deseada real.

**Algoritmo 5:** FairLearning( $\mathcal{D}, \mathcal{M}$ )**Entrada:**  $\mathcal{D}$ , conjunto de datos y  $\mathcal{M}$ , modelo causal.**Salida:** parámetros aprendidos  $\hat{\theta}$ .**para**  $i \in \mathcal{D}$  **hacer**    muestrear  $m$  ejemplos MCMC  $U_1^{(i)}, \dots, U_m^{(i)} \sim P_{\mathcal{M}}(U \mid x^{(i)}, a^{(i)})$ .**fin**Creamos  $\mathcal{D}'$  donde cada punto  $(a^{(i)}, x^{(i)}, y^{(i)})$  en  $\mathcal{D}$  es sustituido por los correspondientes  $m$  puntos  $\{(a^{(i)}, x^{(i)}, y^{(i)}, u_j^{(i)})\}$ ; $\hat{\theta} \leftarrow \operatorname{argmin}_{\theta} \sum_{i' \in \mathcal{D}'} l(y^{(i')}, g_{\theta}(U^{(i')}, x_{\neq A}^{(i')}))$ ;**devolver**  $g_{\hat{\theta}}$ ;

## 12.2 DISEÑO DEL MODELO CAUSAL DE ENTRADA

El modelo  $\mathcal{M}$  debe proporcionarse al Algoritmo 5. Recordemos que los modelos causales requieren de fuertes suposiciones cuando se hacen afirmaciones contrafactuales (Barocas et al. [2019]). Existen infinitas ecuaciones estructurales compatibles con la misma distribución observable por lo que, teóricamente estos modelos deberían ser propensos a modificaciones si, por ejemplo se incorporan nuevos datos anteriormente no observados que contradijesen el modelo actual.

En nuestro trabajo, no será necesario especificar un modelo totalmente determinista y relajaremos las ecuaciones estructurales a partir de la definición de distribución condicional. En particular, el concepto de equidad contrafactual es válido bajo los siguientes tres niveles que vienen especificados en el trabajo de Kusner et al. [2018] como:

- **Nivel 1:** predecir  $\hat{Y}$  usando sólo las variables observables no descendientes de  $A$ . Esta consideración no requiere de ninguna suposición causal, pero en la mayoría de los problemas la mayor parte de las variables observables serán descendientes de atributos protegidos y por tanto tendremos menos información manejable.
- **Nivel 2:** suponer variables de ruido que actúan como causas no deterministas de las variables observables, basadas en el conocimiento explícito del dominio y en algoritmos de aprendizaje. La información sobre  $X$  pasará a  $\hat{Y}$  a través de  $P(U \mid x, a)$ .
- **Nivel 3:** construir un modelo totalmente determinista con variables de ruido. Por ejemplo, la distribución  $P(V_i \mid pa_i)$  puede tratarse como un modelo de error aditivo,  $V_i = f_i(pa_i) + \epsilon_i$  (Peters et al. [2014]). El término de error  $\epsilon_i$  sirve como una entrada a  $\hat{Y}$  calculada a partir de las variables observadas. Esto maximizará la información extraída por el clasificador.

### 12.3 APLICACIÓN EN UN PROBLEMA REAL

Ilustraremos la aplicación de justicia contrafactual sobre un problema del mundo real que requiere equidad. El objetivo de este experimento es cuantificar el comportamiento del Algoritmo 5 con tamaños de muestra finitos sobre una suposición real compatible con un modelo sintético.

#### 12.3.1 Descripción del problema

El Consejo de Admisión de las Facultades de Derecho realizó una encuesta en 163 facultades de Derecho de Estados Unidos (Wightman [1998]). Contiene información sobre 21.790 estudiantes de Derecho, tales como las puntuaciones de su examen de acceso (LSAT), su media del expediente (GPA) antes de entrar en la facultad, y su nota media del primer año (FYA) en la carrera de Derecho.

A partir de estos datos, una escuela podría querer predecir si un solicitante tendrá un FYA alto. También podría ser interesante asegurarse de que estas predicciones no están sesgadas por la raza y el sexo del individuo. Sin embargo, es bastante probable que los resultados del LSAT, GPA y FYA estén sesgados por factores sociales. Nuestro trabajo consistirá en aplicar las herramientas aprendidas para equidad contrafactual en diversos escenarios y comparar su actuación con el rendimiento de otras nociones de equidad estudiadas.

#### 12.3.2 Escenarios de predicción

Utilizaremos el mismo escenario de experimentación que el propuesto en Kusner et al. [2018]. Dividiremos el conjunto de datos en un 80-20 (entrenamiento-test) para evaluar los modelos, preservando el equilibrio de las etiquetas. Para predecir los resultados utilizaremos un predictor basado en regresión lineal y mediremos la exactitud alcanzada por cada modelo a con la métrica RMSE.

Según hemos descrito en la Sección 12.2, existen tres aproximaciones a partir de las cuales podemos construir un predictor de FYA que satisfaga la equidad contrafactual:

- **Nivel 1:** usaremos cualquier característica que no sea descendiente de la raza y el sexo para la predicción. Como creemos que el LSAT, el GPA y el FYA están sesgados por la raza y el sexo, no podremos utilizar ninguna de las características observadas para construir un clasificador justo contrafactual.
- **Nivel 2 (Fair K):** supondremos que una variable de ruido: los conocimientos del estudiante ( $K$ ), afecta a las puntuaciones de GPA, LSAT y FYA. El gráfico

causal correspondiente a este modelo se muestra en la Figura 26. Emplearemos las siguientes distribuciones:

$$\begin{aligned} \text{GPA} &\sim \mathcal{N}(b_G + K\mathbf{w}_G^K + [A_R, A_S]\mathbf{w}_G^A, \sigma_G), \\ \text{LSAT} &\sim \text{Poisson}(\exp(b_L + K\mathbf{w}_L^K + [A_R, A_S]\mathbf{w}_L^A)), \\ \text{FYA} &\sim \mathcal{N}(K\mathbf{w}_F^K + [A_R, A_S]\mathbf{w}_F^A, 1), \\ K &\sim \mathcal{N}(0, 1). \end{aligned}$$

Realizamos la inferencia sobre este modelo utilizando un conjunto de entrenamiento observado para estimar la distribución posterior de  $K$ .

- **Nivel 3** (*Fair Add*): modelaremos las puntuaciones de GPA, LSAT y FYA como variables continuas con términos de error aditivos independientes de la raza y el sexo. Estimamos los términos de error  $\epsilon_G, \epsilon_L$  ajustando primero dos modelos que utilizan la raza y el sexo para predecir individualmente el GPA y LSAT. A continuación, calculamos los residuos de cada modelo (por ejemplo,  $\epsilon_G = \text{GPA} - \hat{Y}_{\text{GPA}}(A_R, A_S)$ ). Finalmente, usaremos las estimaciones residuales de  $\epsilon_G, \epsilon_L$  para predecir FYA. Este modelo se muestra en la Figura 26. En este caso las distribuciones vienen dadas por:

$$\begin{aligned} \text{GPA} &= b_G + [A_R, A_S]\mathbf{w}_G^A + \epsilon_G, \quad \epsilon_G \sim P(\epsilon_G) \\ \text{LSAT} &= b_L + [A_R, A_S]\mathbf{w}_L^A + \epsilon_L, \quad \epsilon_L \sim P(\epsilon_L), \\ \text{FYA} &= b_F + [A_R, A_S]\mathbf{w}_F^A + \epsilon_F, \quad \epsilon_F \sim P(\epsilon_F). \end{aligned}$$

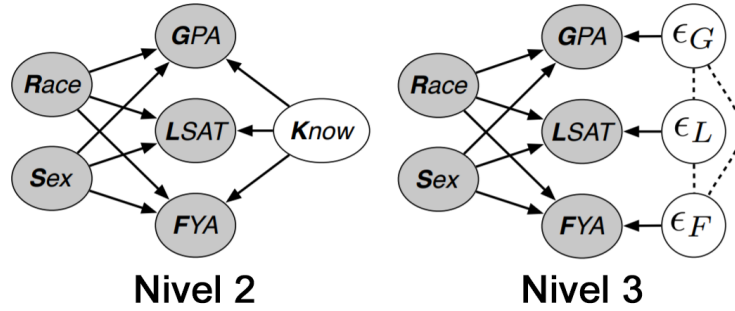


Figura 26: Grafos causales de los escenarios para los niveles 2 y 3, respectivamente.

Compararemos el escenario propuesto por la equidad contrafactual con dos líneas de base injustas:

- **Total**: utilizará todas las características disponibles para el individuo, incluidos los atributos sensibles.
- **Desconocimiento**: aplicaremos la noción de equidad por desconocimiento discutida en la Sección 7.2.

---

## IMPLEMENTACIÓN Y RESULTADOS

---

En este capítulo explicaremos los procedimientos empleados en la implementación de los diseños y algoritmos definidos en el capítulo anterior. Todo lo relativo a implementación se encuentra disponible en: [danibolanos/TFG-Guarantee\\_Fairness\\_in\\_ML](https://github.com/danibolanos/TFG-Guarantee_Fairness_in_ML)

### 13.1 OBTENCIÓN Y TRATAMIENTO DE LOS DATOS

El conjunto de datos ha sido extraído del repositorio Ethik, que contiene diferentes conjuntos de datos útiles para el estudio de equidad en aprendizaje automático y que está disponible en: <https://github.com/XAI-ANITI/ethik/tree/master/ethik/data>.

Es recomendable que antes de aplicar cualquier modelo o algoritmo sobre un conjunto de datos, se realice un preprocesamiento sobre estos. En nuestro caso, haciendo uso de la biblioteca Pandas (<https://pandas.pydata.org/docs/>) para Python, hemos realizado las siguientes modificaciones sobre el conjunto de datos original:

- Categorizar cada valor del atributo raza de manera que obtengamos una columna para cada tipo de raza que tome el valor 1 o 0 para cada individuo, indicando la pertenencia o no del mismo a la raza indicada por la columna en concreto (*get\_dummies()*).
- Sustituir el atributo sexo por dos columnas ('Male'-'Female') que indiquen con 1 o 0 la característica del individuo concreto.
- Discretizar el valor de 'LSAT' (convertir cada valor a tipo entero).

### 13.2 IMPLEMENTACIÓN DEL CÓDIGO

El lenguaje de programación seleccionado para la implementación del proyecto ha sido Python (<https://www.python.org/>) en su versión 3.8.5. La elección de Python se debe a que es un lenguaje muy popular en el ámbito de la ciencia de datos, con una gran cantidad de bibliotecas útiles para la visualización de datos y altamente compatible con otros lenguajes de programación que nos pueden ser útiles tanto como para el tratamiento de datos como para la construcción de modelos causales.

### 13.3 CONTRASTE DE LOS RESULTADOS

#### 13.4 CONDICIONES DE LA EXPERIMENTACIÓN

##### 13.4.1 *Entorno de ejecución*

Los experimentos realizados en este capítulo se han ejecutado en un equipo con las siguientes características:

- Arquitectura x86\_64.
- AMD Ryzen 7 4800H with Radeon Graphics 2.90 GHz.
- 8 núcleos con 2 hilos de procesamiento por núcleo.
- 16 GB RAM DDR4.
- Sistema Operativo: Ubuntu 20.04.2 LTS.

##### 13.4.2 *Entorno de programación*

Spyder.

##### 13.4.3 *Bibliotecas y herramientas auxiliares*

Algunas de las bibliotecas más importantes utilizadas junto al número de sus versiones han sido:

- Pandas 1.2.4
- Numpy 1.19.2
- Scikit-learn 0.24.2
- **pickle**,

Además se ha utilizado el módulo PyStan 2.19.1.1 (<https://pystan.readthedocs.io/en/latest>) que funciona como una interfaz de Python para el lenguaje de programación probabilística Stan (<https://mc-stan.org>). El cual permite generar muestras de datos a partir de métodos que trabajan con cadenas de Markov (como MCMC) y crear modelos causales sobre los que operar gracias a la inferencia estadística Bayesiana.

### 13.5 MANUAL DE EJECUCIÓN DEL EXPERIMENTO



## Parte V

### CONCLUSIONES Y VÍAS FUTURAS

Conclusiones extraídas a lo largo del desarrollo del proyecto y debate de posibles rutas de trabajo futuras para el mismo.

---

## CONCLUSIÓN

---

conclusiones obtenidas del trabajo en general

---

## TRABAJOS FUTUROS

---

Equidad individual se ha quedado sin tratar, trabajo futuro para TFM de matemáticas estudio de clase de funciones que puedan ser interesantes para implementar la métrica de distancia entre individuos

Pensar si incluir el esquema y prototipo en este apartado o en uno nuevo. Preguntar!

---

```
import numpy as np
import cv2
from matplotlib import pyplot as plt
import math

import library

def leeimagen(filename, flagColor):
    im = cv2.imread(filename)
    return im
```

---

## APÉNDICES



---

## HERRAMIENTAS PARA GARANTIZAR JUSTICIA EN AA

---

En este capítulo, hablaremos de algunas herramientas y técnicas que existen actualmente para garantizar la equidad en aprendizaje automático. Hablaremos de algunos de los más influyentes, entre los que destacaremos Aequitas por ser en el que nos vamos a basar en este trabajo.

### A.1 AEQUITAS

**NO REVISAR:** Se completará con los experimentos usando la biblioteca de Aequitas

Aequitas (Saleiro et al. [2019]) es una herramienta de auditoría desarrollada por el *Center for Data Science and Public Policy* de la Universidad de Chicago. Es una herramienta de código abierto que consta de diversas utilidades de soporte para la auditoría de sesgos creado para ser utilizado por analistas de todo tipo relacionados con el ámbito del aprendizaje automático y cuyo principal objetivo es auditar los modelos de *machine learning* con el fin de encontrar posibles discriminaciones en ellos y evitarlas en un futuro.

Aequitas nos permite detectar dos tipos de sesgos:

- Acciones sesgadas que no ocurren de forma representativa en la población.
- Resultados sesgados a causa de errores de clasificación de nuestro sistema con respecto a ciertos grupos de la población.

Para utilizar la herramienta, se necesitan aportar los siguientes datos:

- Datos sobre los atributos específicos (raza, sexo, etc.) que queramos auditar.
- El conjunto de personas de la población mencionada que el sistema de evaluación de riesgos seleccionó para una intervención.

#### A.1.1 Estructura de los datos de entrada y resultados

Podemos dividir en tres apartados (conformados por columnas en Aequitas) los datos que debemos aportar para el correcto funcionamiento de la herramienta.

- **score**: representa la conclusión a la que llega un modelo, puede ser binaria (0 o 1) o continua (decimal entre 0 y 1). Esta decisión representa si el sujeto es apto o no, por ejemplo, si se le concede un crédito bancario.
- **label\_value**: representa los datos reales, es decir, si la predicción realizada por el modelo fue correcta. Por ejemplo, el sujeto fue capaz de devolver el crédito en su totalidad. Es por esto, por lo que el modelo solo puede ser auditado después de su aplicación y no antes. Se representa como un valor binario, 1 significa que la predicción fue correcta, 0 que no lo fue.
- **attributes**: categorías de los atributos definidos por el usuario y utilizados para decidir la equidad del modelo. Algunos ejemplos de atributos son la raza, sexo, edad o ingresos.

score	label_value	race	sex	age_cat
0	1	Hispanic	Male	Less than 25
1	0	African-American	Female	25-45
0	0	Caucasian	Male	25-45

Cuadro 5: Ejemplo del *dataset* COMPAS aportado a Aequitas.

Para entender como funciona Aequitas, necesitamos presentar los siguientes conceptos preliminares definidos en: [https://dssg.github.io/aequitas/30\\_seconds\\_aequitas.html](https://dssg.github.io/aequitas/30_seconds_aequitas.html)

Nombre	Notación	Definición
Score	$S \in [0, 1]$	puntuación de valor real asignada a cada entidad por el clasificador.
Decision	$\hat{Y} \in \{0, 1\}$	predicción binaria asignada a una entidad dada ( <i>data point</i> ).
True Outcome	$Y \in \{0, 1\}$	etiqueta binaria de una entidad dada.
Attribute	$A = \{a_1, a_2, \dots, a_n\}$	atributo multivalor, por ejemplo., género={femenino,masculino,otro}
Group	$g(a_i)$	todas las entidades que comparten el mismo valor de atributo, p. ej., género=femenino
Reference group	$g(a_r)$	uno de los grupos de $A$ que es usado como referencia para calcular las métricas de sesgo.
Labeled Positive	$LP_g$	número de entidades etiquetadas como positivas dentro de un grupo.
Labeled Negative	$LN_g$	número de entidades etiquetadas como negativas dentro de un grupo.
Predicted Positive	$PP_g$	número de entidades dentro de un grupo cuya predicción es positiva, es decir, $\hat{Y} = 1$ .
Total Pred. Positive	$K = \sum_{A=a_1}^{A=a_n} PP_g(a_i)$	número total de entidades $PP_g$ a lo largo de los grupos definidos por $A$ .
Predicted Negative	$PN_g$	número de entidades dentro de un grupo cuya predicción es negativa, es decir $\hat{Y} = 0$ .
False Positive	$FP_g$	número de entidades de un grupo con $\hat{Y} = 1 \wedge Y = 0$ .
False Negative	$FN_g$	número de entidades de un grupo con $\hat{Y} = 0 \wedge Y = 1$ .
True Positive	$TP_g$	número de entidades de un grupo con $\hat{Y} = 1 \wedge Y = 1$ .
True Negative	$TN_g$	número de entidades de un grupo con $\hat{Y} = 0 \wedge Y = 0$ .

Cuadro 6: Conceptos preliminares de Aequitas.

La herramienta produce como resultado un informe en formato pdf que devuelve una interpretación descriptiva de los resultados junto con tres conjuntos de tablas con información relevante acerca de tres tipos de métricas:

- Métricas de grupo.
- Métricas de sesgo.

■ Medidas de equidad.

### A.1.2 Métricas usadas por Aequitas

A continuación, veremos como se definen las diferentes métricas a partir de los conceptos preliminares definidos en el Cuadro 6. Y mostraremos un ejemplo sobre el conjunto de datos COMPAS.

#### A.1.2.1 Métricas de grupo

Nombre	Notación	Definición
Prevalence	$Prev_g = LP_g /  g  = Pr(Y = 1 \mid A = a_i)$	fracción de entidades dentro de un grupo cuyo valor real de la etiqueta fue positivo.
Predicted Prevalence	$PPrev_g = PP_g /  g  = Pr(\hat{Y} = 1 \mid A = a_i)$	fracción de entidades dentro de un grupo que fue predicho como positivo.
Predicted Positive Rate	$PPR_g = PP_g / K = Pr(a = a_i \mid \hat{Y} = 1)$	fracción de entidades predichas como positivas que pertenecen a un determinado grupo.
False Discovery Rate	$FDR_g = FP_g / PP_g = Pr(Y = 0 \mid \hat{Y} = 1, A = a_i)$	fracción de falsos positivos de un grupo entre el $PP_g$ del grupo.
False Omission Rate	$FOR_g = FN_g / PN_g = Pr(Y = 1 \mid \hat{Y} = 0, A = a_i)$	fracción de falsos negativos de un grupo entre el $PN_g$ del grupo.
False Positive Rate	$FPR_g = FP_g / LN_g = Pr(\hat{Y} = 1 \mid Y = 0, A = a_i)$	fracción de falsos positivos de un grupo entre los etiquetados negativos del grupo.
False Negative Rate	$FNR_g = FN_g / LP_g = Pr(\hat{Y} = 0 \mid Y = 1, A = a_i)$	fracción de los falsos negativos de un grupo entre los etiquetados positivos del grupo.

Cuadro 7: Métricas de grupo de Aequitas.

Primero calcularemos el valor de las métricas de grupo para un atributo (p.ej. raza) teniendo en cuenta los conceptos del Cuadro 7.

**race**

Attribute Value	PPR	PPREV	FDR	FPR	FOR	FNR
0 Amer-Indian-Eskimo	0.01	0.93	0.94	0.97	0.57	0.4
1 Asian-Pac-Islander	0.03	0.74	0.88	0.89	0.68	0.66
2 Black	0.11	0.88	0.94	0.95	0.66	0.58
3 Other	0.01	0.92	0.96	0.98	0.83	0.62
4 White	0.83	0.75	0.89	0.9	0.69	0.67

Figura 27: Tabla con las principales métricas de grupo para el atributo *race*.

#### A.1.2.2 Métricas de sesgo

Mide la disparidad entre un grupo y el grupo de referencia. La disparidad se calcula a partir de la siguiente fórmula:

$$DisparityMeasure_{ProtectedGroup} = \frac{GroupMetric_{ProtectedGroup}}{GroupMetric_{ReferenceGroup}}$$

Donde *GroupMetric* hace referencia a una métrica de grupo del Cuadro 7. Es evidente que la disparidad de cualquier medida sobre el grupo de referencia siempre será 1. Si queremos calcular por ejemplo la disparidad del ratio de falsos negativos (*FNR*) sobre el grupo de raza negra, se calculará de la siguiente forma:

$$FNR_{Black} = \frac{FNR_{Black}}{FNR_{White}} = \frac{0,58}{0,67} = 0,86$$

Completando la tabla para todas las métricas de grupo obtenemos el siguiente resultado:

**race**

	Attribute Value	PPR Disparity	PPREV Disparity	FDR Disparity	FPR Disparity	FOR Disparity	FNR Disparity
0	Amer-Indian-Eskimo	0.01	1.24	1.05	1.08	0.83	0.59
1	Asian-Pac-Islander	0.04	0.98	0.98	0.99	0.99	0.98
2	Black	0.14	1.17	1.05	1.06	0.95	0.86
3	Other	0.01	1.22	1.07	1.1	1.21	0.93
4	White	1.0	1.0	1.0	1.0	1.0	1.0

Figura 28: Tabla con las métricas de sesgo para el atributo *race*.

*Comentario 5.* De forma predeterminada, Aequitas usa el grupo mayoritario dentro de cada atributo como grupo de referencia.

### A.1.2.3 Medidas de equidad

La equidad siempre se define en relación con un grupo de referencia. Podemos ver que el cálculo de la equidad, depende de la métrica de sesgo. En la evaluación del criterio de equidad, un grupo cumple con la paridad si

$$(1 - \varepsilon) \leq \text{DisparityMeasure}_{\text{group}_i} \leq \frac{1}{(1 - \varepsilon)}$$

donde  $\varepsilon$  es el umbral de equidad definido.

Tomando  $\varepsilon = 0,2$  cualquier métrica de sesgo se considerará justa en el intervalo

$$\left[1 - 0,2, \frac{1}{1 - 0,2}\right] = [0,8, 1,25]$$

En el informe de resultados final que devuelve Aequitas si todas las métricas de equidad contienen el flag *fair*, se evaluará el modelo actual como justo. De lo contrario, lo considerará injusto y enumerará los grupos afectados injustamente según los criterios de equidad dados.



race

Attribute Value	Statistical Parity	Impact Parity	FDR Parity	FPR Parity	FOR Parity	FNR Parity
0 Amer-Indian-Eskimo	Unfair	Fair	Fair	Fair	Fair	Unfair
1 Asian-Pac-Islander	Unfair	Fair	Fair	Fair	Fair	Fair
2 Black	Unfair	Fair	Fair	Fair	Fair	Fair
3 Other	Unfair	Fair	Fair	Fair	Fair	Fair
4 White	Ref	Ref	Ref	Ref	Ref	Ref

Figura 29: Tabla de medidas de equidad aplicado el umbral.

## The Bias Report

The Bias Report evaluates the current model as **unfair** using the following fairness criteria:

Fairness Criteria	Desired Outcome	Unfairly Affected Groups
0 Equal Parity	Each group is represented equally.	race:Amer-Indian-Eskimo, race:Asian-Pac-Islander, race:Black, race:Other
1 Proportional Parity	Each group is represented proportional to their representation in the overall population.	No Unfair Groups Found
2 False Positive Parity	Each group has proportionally equal false positive errors made by the model.	No Unfair Groups Found
3 False Negative Parity	Each group has proportionally equal false negative errors made by the model.	race:Amer-Indian-Eskimo

Figura 30: Resultado del modelo injusto.

En la Figura 30 podemos ver que el modelo no cumple con los criterios de Paridad Falsa Negativa (*False Negative Parity*) y de Igual Paridad (*Equal Parity*).

Podemos ver que los conceptos de equidad que utiliza AeQUITAS son los siguientes:

- Equal Parity.
- Proportional Parity.
- False Positive Parity.
- False Negative Parity.

La versión actual es un esquema de lo que se intentará mostrar con imágenes sacadas de la web que deben cambiarse. Idea probar en AeQUITAS con las medidas existentes el mismo dataset que usaré con contrafactual, si no se puede, usar COMPAS.

---

## ESTIMACIÓN DEL COSTE Y PLANIFICACIÓN

---

En este apéndice se realizará una estimación tanto del coste como de la planificación del trabajo durante el período de desarrollo con el objetivo de simular la valoración y presupuesto de un proyecto real en el ámbito laboral.

### B.1 ESTIMACIÓN DEL PRESUPUESTO DEL PROYECTO

Haremos un presupuesto del proyecto, en el que incluiremos las horas dedicadas a cada tema estudiado y realizaremos una estimación a precio 7 euros/hora. El análisis del presupuesto se puede observar en el Cuadro 8.

Para facilitar la estimación de los costes, hemos dividido el trabajo en tres partes:

- **Parte teórica:** recoge las tareas relacionadas con el análisis y el estudio de los conceptos de carácter teórico contenidos en el trabajo. Incluiremos la formalización de las nociones básicas para el proyecto y el desarrollo de las demostraciones matemáticas.
- **Parte práctica:** reúne las prácticas relacionadas con la programación, análisis y validación de los experimentos. Además tendremos en cuenta el equipo utilizado y los tiempos dedicados a instalación de bibliotecas y software empleados.
- **Parte general:** agrupa las labores de elaboración de la memoria y reuniones con los tutores (presencial u online).

El cómputo total es de 12.198 euros. Teniendo en cuenta que el período de trabajo útil ha sido de aproximadamente 7 meses, el sueldo medio mensual equivale a 1.742 euros brutos. Para un informático junior, este valor es bastante fiel a la realidad.

### B.2 PLANIFICACIÓN DEL TRABAJO

El diseño de la planificación se ha realizado con el software GanttProject (<https://www.ganttproject.biz>) que es un programa de código abierto utilizado para administrar proyectos pudiendo usar entre otras muchas herramientas, un diagrama de Gantt. Para el desarrollo de la planificación, hemos usado la misma división en partes que la presentada en la sección anterior.

En la planificación original, nos habría gustado presentar el trabajo para la convocatoria de septiembre, pero debido a algunos problemas derivados de la carga docente a lo largo del curso 2020-2021 y a la concreción de los experimentos asociados al proyecto, se ha acabado retrasando hasta el mes de noviembre.

Además es destacable, si comparamos las Figuras 31 y 32, que el tiempo dedicado tanto al estudio como la formalización de las diferentes medidas de equidad, ha sido mayor al previsto inicialmente.

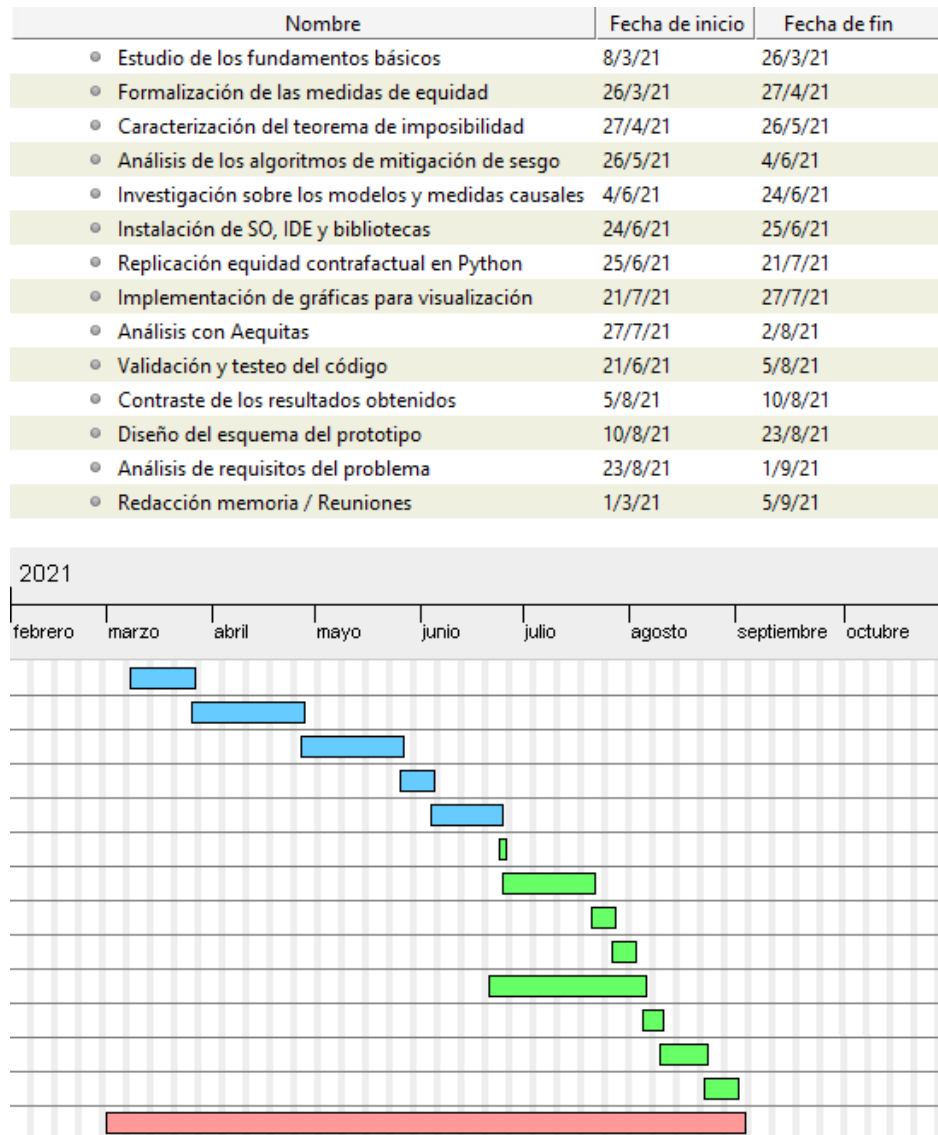


Figura 31: Planificación original del desarrollo del trabajo.

En ambos casos, la parte general es una tarea que se desarrolla a lo largo del período de trabajo, ya que la redacción de la memoria y las tutorías se realizan de forma simultánea al análisis teórico y práctico del proyecto.

En la planificación final, hemos establecido el mes de abril como inicio del proyecto. Además el periodo de exámenes de junio también ralentizó en gran medida el avance del proyecto propuesto. Todos estos factores se pueden intuir observando la Figura 32, donde indicamos en color azul la parte teórica, en verde la parte práctica y en rojo la parte general.

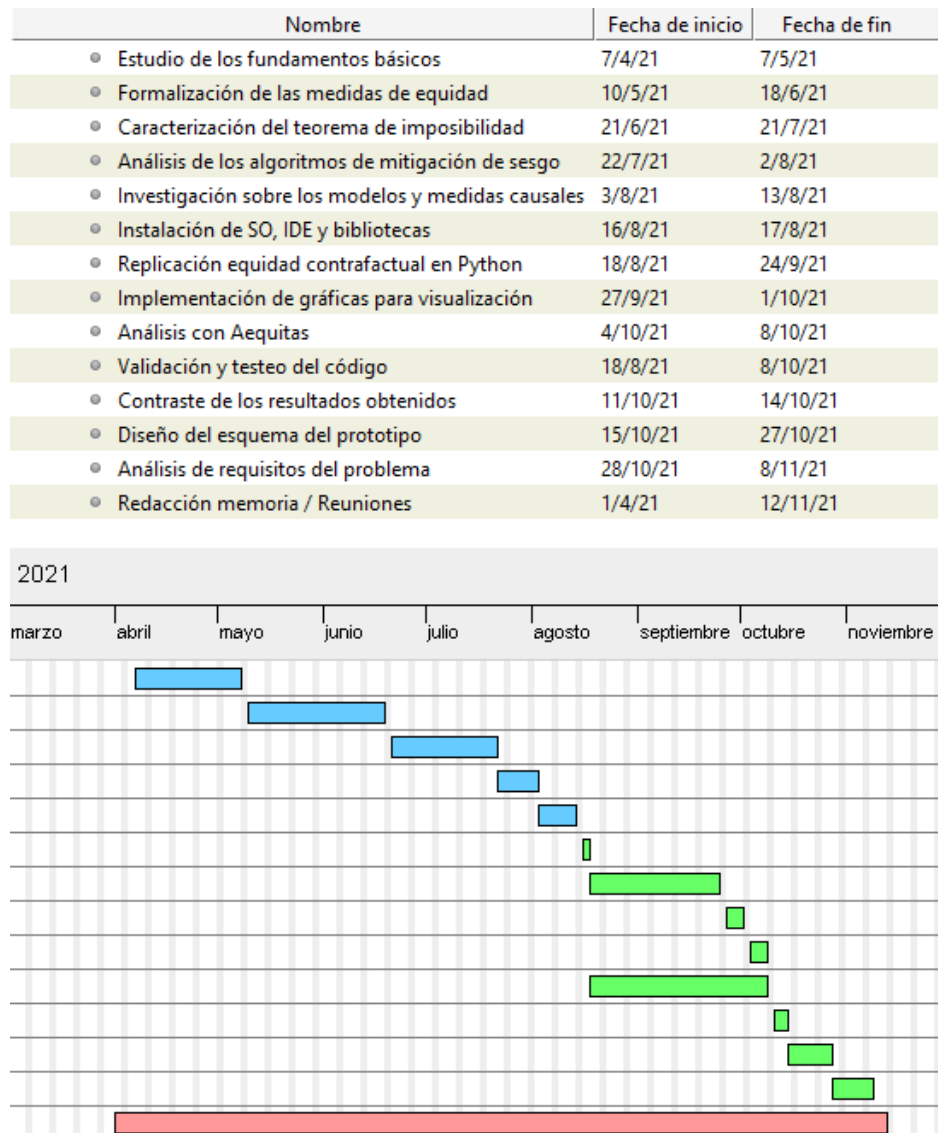


Figura 32: Planificación final del desarrollo del trabajo.

Concepto	Tiempo (horas)	Coste (euros/hora)	Coste total (euros)
<b>Parte general</b>			
Redacción de la memoria	890	7	6.230
Reuniones con los tutores	15	-	-
<b>Parte teórica</b>			
Estudio de los fundamentos básicos	80	7	560
Formalización de las medidas de equidad	150	7	1.050
Caracterización del teorema de imposibilidad	120	7	840
Análisis de los algoritmos de mitigación de sesgo	40	7	280
Investigación sobre los modelos y medidas causales	35	7	245
<b>Parte práctica</b>			
Equipo de trabajo: ASUS TUF	-	-	949
Instalación de SO, IDE y bibliotecas	2	7	14
Replicación de un modelo de equidad contrafactual en Python	70	7	490
Implementación de gráficas para visualización	10	7	70
Análisis con Aequitas	10	7	70
Validación y testeo del código	90	7	630
Contraste de los resultados obtenidos	15	7	105
Diseño del esquema del prototipo	50	7	350
Análisis de requisitos del problema	45	7	315
<b>Cómputo total del proyecto</b>	<b>1.622</b>	<b>-</b>	<b>12.198</b>

Cuadro 8: Estimación del coste del proyecto.

---

## NOTACIÓN

---

$P(A, B)$	Equivale a $P(A \cap B)$ .
$X$	Variable aleatoria observada.
<b><math>X</math></b>	En negrita, indica una variable aleatoria multivariante $(X_1, \dots, X_n)$ .
i.i.d.	Referido a variables aleatorias independientes e idénticamente distribuidas.

---

## BIBLIOGRAFÍA

---

- Yaser Abu-Mostafa, Malik Magdon-Ismael, and Hsuan-Tien Lin. *Learning from Data: A Short Course*. AMLBook, 2012. ISBN 1600490069.
- Adverse Impact Analysis / Four-Fifths Rule. Adverse impact analysis / four-fifths rule. <https://www.prevuehr.com/resources/insights/adverse-impact-analysis-four-fifths-rule/>, Center for Data Science and Public Policy, 2009.
- Christophe Andrieu, Nando De Freitas, Arnaud Doucet, and Michael I. Jordan. An introduction to mcmc for machine learning, 2003.
- Julia Angwin, Jeff Larson, Surya Mattu, and Lauren Kirchner. Machine bias: The software used across the country to predict future criminals. and it's biased against blacks. *ProPublica*, 2016. <https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing>.
- Solon Barocas and Andrew D. Selbst. Big data's disparate impact. *California Law Review*, 2016.
- Solon Barocas, Moritz Hardt, and Arvind Narayanan. *Fairness in Machine Learning: Limitations and Opportunities*. fairmlbook.org, 2019. <http://www.fairmlbook.org>.
- Joseph Berkson. Limitations of the application of fourfold table analysis to hospital data. *International Journal of Epidemiology* 43, no. 2, 2014.
- Reuben Binns. Fairness in machine learning: Lessons from political philosophy, 2021.
- Christopher M. Bishop. *Pattern Recognition and Machine Learning*. Springer, 2006. ISBN 0387310738.
- Flavio P. Calmon, Dennis Wei, Karthikeyan Natesan Ramamurthy, and Kush R. Varshney. Optimized data pre-processing for discrimination prevention, 2017.
- Michael Collins. Convergence proof for the perceptron algorithm. *Columbia University*, 2012. <http://www.cs.columbia.edu/~mcollins/courses/6998-2012/notes/perc.converge.pdf>.
- Sam Corbett-Davies, Emma Pierson, Avi Feller, Sharad Goel, and Aziz Huq. Algorithmic decision making and the cost of fairness. *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2017.

- Frederick H. Cramer. *Random Variables and Probability Distributions*. Cambridge University Press, 2004. ISBN 0521604869.
- Bart Custers, Toon Calders, Bart Schermer, and Tal Zarsky. *Discrimination and Privacy in the Information Society: Data Mining and Profiling in Large*. Springer, 2012. ISBN 3642304869.
- Amir Dembo. *Probability Theory: STAT310/MATH230*. CreateSpace Independent Publishing Platform, 2014. ISBN 1502955652.
- Nicholas Diakopoulos. *Algorithmic-accountability: the investigation of black boxes. Tow Center for Digital Journalism*, 2014.
- Cynthia Dwork, Moritz Hardt, Toniann Pitassi, Omer Reingold, and Rich Zemel. *Fairness through awareness*, 2011.
- Roland Fryer, Glenn Loury, and Tolga Yuret. An economic analysis of color-blind affirmative action. *Journal of Law, Economics and Organization*, 2008.
- Andreas Fuster, Paul Goldsmith-Pinkham, Tarun Ramadorai, and Ansgar Walther. Predictably unequal? the effects of machine learning on credit markets. *Journal of Finance*, 2018.
- Pratik Gajane and Mykola Pechenizkiy. On formalizing fairness in prediction with machine learning, 2018.
- Chris Godsil and Gordon Royle. *Algebraic Graph Theory*. Springer, 2001. ISBN 0387952209.
- Moritz Hardt, Eric Price, and Nathan Srebro. Equality of opportunity in supervised learning, 2016.
- L.L. Humphrey, B.K.S. Chan, and H.C. Sox. Postmenopausal hormone replacement therapy and the primary prevention of cardiovascular disease. *Annals of Internal Medicine* 137, no. 4, 2002.
- Ricardo Vélez Ibarrola and Alfonso García Pérez. *Principios de inferencia estadística*. UNED, 2012. ISBN 8436265777.
- Jon Kleinberg, Sendhil Mullainathan, and Manish Raghavan. Inherent trade-offs in the fair determination of risk scores, 2016.
- Matt J. Kusner, Joshua R. Loftus, Chris Russell, and Ricardo Silva. Counterfactual fairness, 2018.
- Claire C. Miller. Can an algorithm hire better than a human? *New York Times*, 2015.



- Shira Mitchell, Eric Potash, Solon Barocas, Alexander D’Amour, and Kristian Lum. Algorithmic fairness: Choices, assumptions, and definitions. *Annual Review of Statistics and Its Application*, 8(1):141–163, Mar 2021. ISSN 2326-831X. doi: 10.1146/annurev-statistics-042720-125902. URL <http://dx.doi.org/10.1146/annurev-statistics-042720-125902>.
- Judea Pearl. *Causality: Models, Reasoning and Inference*. Cambridge University Press, Cambridge, UK, 2000. ISBN 0521773628.
- Jonas Peters, Joris Mooij, Dominik Janzing, and Bernhard Schölkopf. Causal discovery with continuous additive noise models, 2014.
- Benjamin Recht, Rebecca Roelofs, Ludwig Schmidt, and Vaishaal Shankar. Do image-net classifiers generalize to imagenet?, 2019.
- Frank Rosenblatt. The perceptron: A perceiving and recognizing automaton (project para). *Cornell Aeronautical Laboratory*, 1957.
- Chris Russell, Matt J. Kusner, Joshua R. Loftus, and Ricardo Silva. When worlds collide: Integrating different counterfactual assumptions in fairness, 2017.
- Pedro Saleiro, Benedict Kuester, Loren Hinkson, Jesse London, Abby Stevens, Ari Anisfeld, Kit T. Rodolfa, and Rayid Ghani. Aequitas: A bias and fairness audit toolkit, 2019.
- Kailash Karthik Saravanakumar. The impossibility theorem of machine fairness – a causal perspective, 2021.
- Xinyue Shen, Steven Diamond, Yuantao Gu, and Stephen Boyd. Disciplined convex-concave programming, 2016.
- Masashi Sugiyama, Anton Schwaighofer, Neil D. Lawrence, and Joaquin Quiñero-Candela. Dataset shift in machine learning. *The MIT Press*, 2017.
- Title VII of the Civil Rights Act: Equal Employment Opportunities. Title vii of the civil rights act of 1964: Equal employment opportunities. <https://www.eeoc.gov/statutes/title-vii-civil-rights-act-1964>, United States, 2 de julio de 1964.
- Sahil Verma and Julia Rubin. Fairness definitions explained. *Indian Institute of Technology Kanpur and University of British Columbia*, 2018.
- Hao Wang, Berk Ustun, and Flavio P. Calmon. Repairing without retraining: Avoiding disparate impact with counterfactual distributions, 2019.
- Linda F. Wightman. Lsac national longitudinal bar passage study, 1998.
- Blake Woodworth, Suriya Gunasekar, Mesrob I. Ohannessian, and Nathan Srebro. Learning non-discriminatory predictors, 2017.

- Muhammad Bilal Zafar, Isabel Valera, Manuel Gomez Rodriguez, and Krishna P. Gummadi. Fairness beyond disparate treatment & disparate impact. *Proceedings of the 26th International Conference on World Wide Web*, Apr 2017a. doi: 10.1145/3038912.3052660. URL <http://dx.doi.org/10.1145/3038912.3052660>.
- Muhammad Bilal Zafar, Isabel Valera, Manuel Gomez Rodriguez, and Krishna P. Gummadi. Fairness constraints: Mechanisms for fair classification, 2017b.
- Muhammad Bilal Zafar, Isabel Valera, Manuel Gomez Rodriguez, Krishna P. Gummadi, and Adrian Weller. From parity to preference-based notions of fairness in classification, 2017c.
- Richard Zemel, Yu L. Wu, Kevin Swersky, Toniann Pitassi, and Cynthia Dwork. Learning fair representations, 2013.
- Ciyu Zhu, Richard H. Byrd, Jorge Nocedal, and Peihuang Lu. Algorithm 778: L-bfgs-b: Fortran subroutines for large-scale bound-constrained optimization, 1997.