Ciudad de México a 25 de octubre de 2021

Grupo: _	3CV19		Equipo No	6	
Integrantes	:				
		Medina Granados Alan Alejan	dro		
Guzman Gutierrez Manuel					
		Castro Cruces Jorge Eduard	lo		

Bhumika Gupta, Aditya Rawat, Akshay Jain, Arpit Arora, Naresh Dhami. (2017). Analysis of Various Decision Tree Algorithms for Classification in Data Mining. *International Journal of Computer Applications* (0975–8887). Volume 163 – No 8, April 2017.

Leer el artículo mencionado y responder las siguientes preguntas:

Ejercicio de clase:

Ejercicio No. 1

a) Describa la minería de datos

Es el proceso de clasificación de grandes conjuntos de datos para hayar patrones, información relevante y útil; involucran métodos en la intersección de inteligencia, aprendizaje automático, estadísticas y base de datos sistemas.

- b) Explique las razones por las cuales se utilizan los árboles de decisión
 - Son fáciles de comprender e interpretar.
 - Rápida presentación.
 - Tienen un costo no tan alto
 - Se pueden usar tanto para datos numéricos como categóricos.
 - Maneja problemas de múltiple salida.
 - Se explican fácilmente mediante lógica de boole.
- c) ¿Cuál es la diferencia entre un árbol de clasificación y un árbol de regresión?

En el árbol de clasificación las variable dependientes son categóricas y no ordenadas

En el árbol de regresión las variables dependientes son continuas u ordenadas (valores completos).

Ejercicio No. 2

Considerando los algoritmos: ID3, C4.5, CART y Random Forest. Realice un cuadro comparativo que considere los siguientes aspectos: descripción del algoritmo, criterio de partición que utiliza, si utiliza poda o no, tipo de datos que utiliza, ventajas y desventajas.

Criterio	ID3	C4 5	CART	Random Forest
Criterio Descripción	Utiliza una búsqueda codiciosa de arriba hacia abajo a través del conjuntos dados, donde cada atributo en cada nodo del árbol se prueba para seleccionar el atributo que sea mejor para la clasificación de un	Es mucho más rápido, más eficiente en memoria y se utiliza para Construye árboles de decisión más pequeños y produce interpretaciones más intuitivas.	CART Significa árboles de clasificación y regresión	Random Forest Un bosque aleatorio es una colección de árboles simples predictores, de modo que cada árbol produce una respuesta cuando un conjunto de los valores predictores se dan como entrada. También funciona tanto para
	conjunto dado			clasificación y problemas de regresión.
Criterio partición	Las instancias de resultado que son posibles se examinan si pertenecen a la misma clase o no. Para los casos de la misma clase, se utiliza una clase de un solo nombre para denotar, de lo contrario, el las instancias se clasifican sobre la base del atributo de división.	Para la división de atributos categóricos, C4.5 sigue el enfoque similar a los algoritmos ID3. Continuo los atributos siempre generan divisiones binarias. Seleccionar el atributo con la relación de ganancia más alta. Estos pasos se aplican repetidamente a nuevas ramas de árboles.	El árbol de clasificación es construido por CART mediante la división binaria del atributo. La función de regresión de CART se puede utilizar al pronosticar una variable dependiente dado un conjunto de predictores variable durante un período de tiempo determinado	Resuelve la clasificación problemas, la respuesta o el resultado aparece en forma de un membresía de clase, que asocia o clasifica, un conjunto de valores de predictores independientes con la categoría coincidente presente en la variable dependiente.
Poda	El atributo con la mayor ganancia de información se puede seleccionar como atributo de prueba del nodo actual. ID3 es basado en la navaja de Occam.	La ganancia de información sesga el atributo con más número de valores. Por lo tanto, C4.5 usa Gain Ratio que es un criterio de selección menos sesgado.	El índice de Gini se utiliza para seleccionar el atributo de división.	Cada árbol produce una respuesta cuando un conjunto de los valores predictores se dan como entrada.
Tipos de datos que utiliza	Categóricos	Continuos y discritos	Admite tanto continuo como	Admite tanto continuo como

	da ⁻	atos de	datos de	atributos
	atr	ributos	nominales	
	no	ominales		

Criterio	ID3	C4.5	CART	Random Forest
Ventajas	Ārbol rápido y conto Busca en todo el conjunto de datos Encuentra los nodos hoja permitiend o así que los datos de prueba sean podado y reduciendo el número de pruebas.	• C4.5 es fácil de impleme ntar. • C4.5 crea modelos que se pueden interpret ar fácilment e. • Puede manejar valores categóric os y continuo s. • Puede lidiar con el ruido y lidiar con el valor faltante atributos.	• CART puede manejar los valores perdidos automáti camente usando divisione s sustituta s. • Utiliza cualquier combina ción de variables continua s / discretas . • CART realiza automáti camente la selección de variables . • CART puede establece r interacci ones entre variables	Random Forest Reconoce valores atípicos y anomalías en conocimien tos datos. Es uno de los algoritmos de aprendizaje más precisos. disponible. Para muchos conjuntos de datos, produce una gran clasificació n de clasificador es precisos. Da una estimación de las variables importante s en clasificació n.
Desventajas	 Para una muestra pequeña, los datos pueden estar sobre ajustados o sobre clasificado. Para tomar una decisión, solo se prueba un 	Una pequeña variación en los datos puede llevar a diferente s decisione s de árboles cuando	 CART puede tener árboles de decisión inestable s. CART se divide solo por una variable. 	 A veces la clasificació n hecha por Random Forest son dificiles de interpretar por los humanos. Random Forest a veces se sobreajuste con

atributo en	se usa	• No	conjuntos
un	C4.5.	paramétr	de datos
instantáne	• Para un	ico.	con Tareas
0	conjunto		ruidosas de
consumien	de		clasificació
do así	entrena		n /
mucho	miento		regresión.
tiempo.	pequeño,		
	C4.5 no		
	funciona		
	muy		
	bien.		

Ejercicio 3

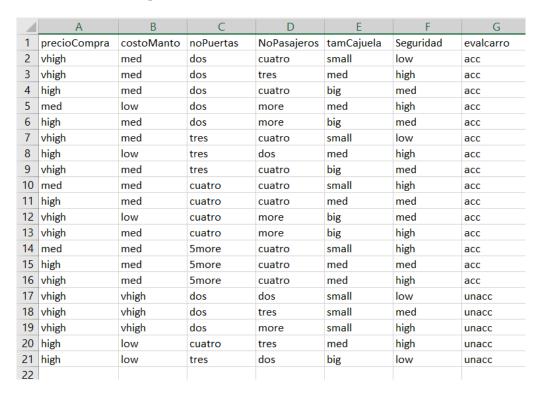
Considerando los siguientes criterios de selección de atributo para particionamiento: *Entropy (Information Gain), Gain Ratio* and *Gini Index.* Realice una descripción con sus propias palabras de cada uno de ellos.

Criterio	Descripción
Entropy (Information Gain)	Valor entre 0-1 que índica la incertidumbre en una variable aleatoria.
Gain Ratio	Medida de ganancia de información, es mejor seleccionar atributos con una gran cantidad de valores
Gini Index	Medida de impureza de un conjunto de datos, son los atributos de división.

Ejercicio 4

Considerando el ejercicio de evaluación del artículo (laptop), aplique el proceso de cálculo de medidas de evaluación al conjunto de datos de carros (accesible, no accesible).

Realizar el proceso de los cinco (visto en clase) pasos de los cálculos para la elección del atributo de particionamiento.



Aplicación del proceso del cálculo de medidas de evaluación al conjunto de datos del carro (accesible, no accesible)

Paso 1: Cálculo de la entropía total

Evaluación Carro	
Accesible	No accesible
15	5

E(Total) =
$$-\frac{15}{20}log_2\left(\frac{15}{20}\right) - \frac{5}{20}log_2\left(\frac{5}{20}\right) = 0.8112$$

Paso 2: Dividir el conjunto de datos en los diversos atributos

Atributo objetivo	Evaluación Carro
Atributo	Dominio
Precio Compra	Muy alto
	Alto
	Medio
Costo Mantenimiento	Muy alto
	Medio
	Bajo
Número Puertas	Dos
	Tres
	Cuatro
	Cinco o más
Número Pasajeros	Dos
	Tres
	Cuatro
	Más
Tamaño Cajuela	Grande
	Mediana
	Pequeña
Seguridad	Alta
	Mediana
	Baja

Paso 3: Se calcula la entropía de cada rama y se suman proporcionalmente para calcular la entropía total

Para Precio Compra

Precio Compra	Muy alto	Accesible	7
		No accesible	3
	Alto	Accesible	5
		No accesible	2
	Mediano	Accesible	3
		No accesible	0

E(Muy alto) =
$$-\frac{7}{20}log_2(\frac{7}{20}) - \frac{3}{20}log_2(\frac{3}{20}) = 0.9406$$

E(Alto) =
$$-\frac{5}{20}log_2(\frac{5}{20}) - \frac{2}{20}log_2(\frac{2}{20}) = 0.8321$$

E(Precio Compra, Evaluación) = $\frac{10}{20}(0.9406) + \frac{7}{20}(0.8321) + \frac{3}{20}(0.4105) = 0.7615$

Para Costo Mantenimiento

Costo	Muy alto	Accesible	0
Mantenimiento		No accesible	3
	Medio	Accesible	12
		No accesible	0
	Bajo	Accesible	3
		No accesible	2

E(Muy alto) = 0

E(Alto) = 0

E(Mediano) =
$$-\frac{3}{20}log_2(\frac{3}{20}) - \frac{2}{20}log_2(\frac{2}{20}) = 0.7427$$

E(Precio Compra, Evaluación) = $\frac{5}{20}$ (0.7427) = 0.1856

Para Número Puertas

Precio Compra	Dos	Accesible	5
		No accesible	3
	Tres	Accesible	3
		No accesible	1
	Cuatro	Accesible	4
		No accesible	1
	Cinco o más	Accesible	3
		No accesible	0

E(Dos) =
$$-\frac{5}{20}log_2(\frac{5}{20}) - \frac{3}{20}log_2(\frac{3}{20}) = 0.9105$$

E(Tres) = $-\frac{3}{20}log_2(\frac{3}{20}) - \frac{1}{20}log_2(\frac{1}{20}) = 0.6266$
E(Cuatro) = $-\frac{4}{20}log_2(\frac{4}{20}) - \frac{1}{20}log_2(\frac{1}{20}) = 0.6804$

E(Tres) =
$$-\frac{3}{20}log_2(\frac{3}{20}) - \frac{1}{20}log_2(\frac{1}{20}) = 0.6266$$

E(Cuatro) =
$$-\frac{4}{20}log_2(\frac{4}{20}) - \frac{1}{20}log_2(\frac{1}{20}) = 0.6804$$

E(Cinco o más) = 0

E(Número Puertas, Evaluación) = $\frac{8}{20}(0.9105) + \frac{4}{20}(0.6266) + \frac{5}{20}(0.6804) + \frac{3}{20}(0.4105) =$ 0.6595

Para Número Pasajeros

	1		
Numero Pasajeros	Dos	Accesible	1
		No accesible	2
	Tres	Accesible	1
		No accesible	2
	Cuatro	Accesible	9
		No accesible	0
	Más	Accesible	4
		No accesible	1

E(Dos) =
$$-\frac{1}{20}log_2(\frac{1}{20}) - \frac{2}{20}log_2(\frac{2}{20}) = 0.5482$$

$$\begin{split} & \text{E(Tres)} = -\frac{1}{20}log_2(\frac{1}{20}) - \frac{2}{20}log_2(\frac{2}{20}) = 0.5482 \\ & \text{E(Cuatro)} = 0 \\ & \text{E(Cinco o más)} = -\frac{4}{20}log_2(\frac{4}{20}) - \frac{1}{20}log_2(\frac{1}{20}) = 0.6804 \\ & \text{E(Número Pasajeros, Evaluación)} = \frac{3}{20}(0.5482) + \frac{3}{20}(0.5482) + \frac{5}{20}(0.6804) = 0.3267 \end{split}$$

Para Tamaño de Cajuela

Costo	Grande	Accesible	5
Mantenimiento		No accesible	1
	Mediana	Accesible	6
		No accesible	1
	Pequeña	Accesible	4
	_	No accesible	3

$$\begin{split} & \text{E(Grande)} = -\frac{5}{20}log_2(\frac{5}{20}) - \frac{1}{20}log_2\left(\frac{1}{20}\right) = 0.715 \\ & \text{E(Mediana)} = -\frac{6}{20}log_2(\frac{6}{20}) - \frac{1}{20}log_2\left(\frac{1}{20}\right) = 0.7371 \\ & \text{E(Pequeña)} = -\frac{4}{20}log_2(\frac{4}{20}) - \frac{3}{20}log_2(\frac{3}{20}) = 0.8749 \\ & \text{E(Tamaño Cajuela, Evaluación)} = \frac{6}{20}(0.715) + \frac{7}{20}(0.7371) + \frac{7}{20}(0.8749) = 0.779 \end{split}$$

Para Seguridad

Seguridad	Alta	Accesible	7
		No accesible	2
	Mediana	Accesible	6
		No accesible	1
	Baja	Accesible	2
		No accesible	2

$$\begin{split} & \text{E(Grande)} = -\frac{7}{20}log_2(\frac{7}{20}) - \frac{2}{20}log_2\left(\frac{2}{20}\right) = 0.8622 \\ & \text{E(Mediana)} = -\frac{6}{20}log_2(\frac{6}{20}) - \frac{1}{20}log_2\left(\frac{1}{20}\right) = 0.7371 \\ & \text{E(Pequeña)} = -\frac{2}{20}log_2(\frac{2}{20}) - \frac{2}{20}log_2(\frac{2}{20}) = 0.6643 \\ & \text{E(Tamaño Cajuela, Evaluación)} = \frac{9}{20}(0.8622) + \frac{7}{20}(0.7371) + \frac{4}{20}(0.6643) = 0.7788 \end{split}$$

Paso 4: Se calcula la ganancia de información

Para Precio Compra: Gain(PrecioCompra, Evaluación) = 0.8112 – 0.7615 = 0.0497 Para Costo Mantenimiento: Gain(Costo Mantenimiento, Evaluación) = 0.8112 – 0.1856 = 0.6256

Para Número Puertas: Gain(Número Puertas, Evaluación) = 0.8112 - 0.6595 = 0.1517

Para Número Pasajeros: Gain(Número Pasajeros, Evaluación) = 0.8112 - 0.3267 = 0.4845

Para Tamaño de Cajuela: Gain(Tamaño Cajuela, Evaluación) = 0.8112 - 0.779 = 0.0322

Para Seguridad: Gain(Seguridad, Evaluación) = 0.8112 - 0.7788 = 0.0324

Paso 5: Elección del nodo de Decisión.

Variable	Ganancia	(Accesible, No Accesible)
Precio de Compra	0.0497	Muy Alto = (7, 3) Alto = (5, 2) Medio = (3, 0)
Costo	0.6256	Muy Alto = (0, 3) Medio = (12, 0) Bajo = (3, 2)
Mantenimiento		
Número Puertas	0.1517	Dos = (5, 3) Tres = (3, 1) Cuatro = (4, 1) Cinco
		o Más = (3, 0)
Número Pasajeros	0.4845	Dos = (1, 2) Tres = (1, 2) Cuatro = (9, 0) Más =
		(4, 1)
Tamaño Cajuela	0.0322	Grande = (5, 1) Mediana = (6, 1) Pequeña = (4,
		3)
Seguridad	0.0324	Alta = (7, 2) Mediana = (6, 1) Baja = (2, 2)

Costo Mantenimiento es la variable que da una mayor ganancia

Se calcula SplitInfo de Mantenimiento:

SplitInfoMantenimimento =
$$-\frac{3}{20}log_2(\frac{3}{20}) - \frac{12}{20}log_2(\frac{12}{20}) - \frac{5}{20}log_2(\frac{5}{20}) = 1.3527$$

GainRatio(Costo Mantenimiento) = $\frac{0.6256}{1.3527} = 0.4624$

Se calcula Gini Index con respecto a la variable Evaluación

Gini(Evaluación) =
$$1 - \left(\frac{15}{20}\right)^2 - \left(\frac{5}{20}\right)^2 = 0.375$$

Se Calcula Gini_{CostoMantenimiento∈(Bajo, Medio)} = $\frac{17}{20}$ Gini(Evaluación1) + $\frac{3}{20}$ Gini(Evaluación2) = $\frac{17}{20}(1 - \left(\frac{15}{17}\right)^2 - \left(\frac{2}{17}\right)^2) + \frac{3}{20}(1 - \left(\frac{3}{3}\right)^2 - \left(\frac{0}{3}\right)^2) = 0.1764$
Se Calcula Gini_{CostoMantenimiento∈(Bajo, Alto)} = $\frac{8}{20}$ Gini(Evaluación1) + $\frac{12}{20}$ Gini(Evaluación2) = $\frac{8}{20}(1 - \left(\frac{3}{8}\right)^2 - \left(\frac{5}{8}\right)^2) + \frac{12}{20}(1 - \left(\frac{12}{12}\right)^2 - \left(\frac{0}{12}\right)^2) = 0.1875$
Se Calcula Gini_{CostoMantenimiento∈(Medio, Muy alto)} = $\frac{15}{20}$ Gini(Evaluación1) + $\frac{5}{20}$ Gini(Evaluación2) = $\frac{15}{20}$ Gini(Evaluación2)

Proponemos realizar este proceso para los demás atributos, para ver si obtenemos valores menores de Gini Index

SplitInfoPrecioCompra =
$$-\frac{10}{20}log_2(\frac{10}{20}) - \frac{7}{20}log_2(\frac{7}{20}) - \frac{3}{20}log_2(\frac{3}{20}) = 1.4406$$

GainRatio(PrecioCompra) = $\frac{0.0497}{1.4406}$ = 0.0344 Se Calcula $Gini_{PrecioComptra \in (Medio, Alto)} = \frac{10}{20} Gini(Evaluación1) + \frac{10}{20} Gini(Evaluación2)$ $= \frac{10}{20} \left(1 - \left(\frac{8}{10}\right)^2 - \left(\frac{2}{10}\right)^2\right) + \frac{10}{20} \left(1 - \left(\frac{7}{10}\right)^2 - \left(\frac{3}{10}\right)^2\right) = 0.37$ $Gini_{PrecioCompra \in (Medio, Muy alto)} = \frac{13}{20} Gini(Evaluación1) +$ $\frac{7}{30}$ Gini(Evaluación2) $= \frac{13}{20} \left(1 - \left(\frac{10}{13}\right)^2 - \left(\frac{3}{13}\right)^2\right) + \frac{7}{20} \left(1 - \left(\frac{5}{7}\right)^2 - \left(\frac{2}{7}\right)^2\right) = 0.3736$ Se Calcula Gini_{PrecioCompra∈(Alto, Muy alto)} = $\frac{17}{20}$ Gini(Evaluación1) + $\frac{3}{20}$ Gini(Evaluación2) = $\frac{17}{20}$ $\left(1 - \left(\frac{12}{17}\right)^2 - \left(\frac{5}{17}\right)^2\right) + \frac{3}{20}\left(1 - \left(\frac{3}{3}\right)^2 - \left(\frac{0}{3}\right)^2\right) = 0.3529$ SplitInfoNumeroPuertas = $-\frac{8}{20}log_2(\frac{8}{20}) - \frac{4}{20}log_2(\frac{4}{20}) - \frac{5}{20}log_2(\frac{5}{20}) - \frac{3}{20}log_2(\frac{3}{20}) =$ 1.9037 GainRatio(NumeroPuestas) = $\frac{0.1517}{19037}$ = 0.0796 Se Calcula $Gini_{N\'umeroPuertas \in (Dos, Tres)} = \frac{12}{20} Gini(Evaluaci\'on1) + \frac{8}{20} Gini(Evaluaci\'on2)$ $= \frac{12}{20} \left(1 - \left(\frac{8}{12} \right)^2 - \left(\frac{4}{12} \right)^2 \right) + \frac{8}{20} \left(1 - \left(\frac{7}{8} \right)^2 - \left(\frac{1}{8} \right)^2 \right) = 0.3541$ Se Calcula $Gini_{N\'umeroPuertas \in (Tres, Cuatro)} = \frac{9}{20} Gini(Evaluaci\'on1) + \frac{11}{20} Gini(Evaluaci\'on2)$ $= \frac{9}{20} \left(1 - \left(\frac{7}{9}\right)^2 - \left(\frac{2}{9}\right)^2\right) + \frac{11}{20} \left(1 - \left(\frac{8}{11}\right)^2 - \left(\frac{3}{11}\right)^2\right) = 0.3737$ $Gini_{N\'umeroPuertas \in (Cuatro, Cinco o m\'as)} = \frac{8}{20} Gini(Evaluaci\'on1) +$ $\frac{12}{20}$ Gini(Evaluación2) $=\frac{8}{20}\left(1-\left(\frac{7}{8}\right)^2-\left(\frac{1}{8}\right)^2\right)+\frac{12}{20}\left(1-\left(\frac{8}{12}\right)^2-\left(\frac{4}{12}\right)^2\right)=0.3541$ $Gini_{N\'umeroPuertas \in (Dos, Cinco o m\'as)} = \frac{11}{20} Gini(Evaluaci\'on1) +$ $\frac{9}{20}$ Gini(Evaluación2) $= \frac{11}{20} \left(1 - \left(\frac{8}{11}\right)^2 - \left(\frac{3}{11}\right)^2\right) + \frac{9}{20} \left(1 - \left(\frac{7}{9}\right)^2 - \left(\frac{2}{9}\right)^2\right) = 0.3737$ Se Calcula $Gini_{N\'umeroPuertas \in (Dos, Cuatro)} = \frac{13}{20} Gini(Evaluaci\'on1) + \frac{7}{20} Gini(Evaluaci\'on2)$ $= \frac{13}{20} (1 - \left(\frac{9}{13}\right)^2 - \left(\frac{4}{14}\right)^2) + \frac{7}{20} (1 - \left(\frac{6}{7}\right)^2 - \left(\frac{1}{7}\right)^2) = 0.3626$ Se Calcula $Gini_{N\'umeroPuertas \in (Tres, Cinco)} = \frac{7}{20} Gini(Evaluaci\'on1) + \frac{13}{20} Gini(Evaluaci\'on2)$ $= \frac{7}{20} (1 - \left(\frac{6}{7}\right)^2 - \left(\frac{1}{7}\right)^2) + \frac{13}{20} (1 - \left(\frac{9}{13}\right)^2 - \left(\frac{4}{13}\right)^2) = 0.3626$ SplitInfoNúmeroPasajeros = $-\frac{3}{20}log_2(\frac{3}{20}) - \frac{3}{20}log_2(\frac{3}{20}) - \frac{9}{20}log_2(\frac{9}{20}) - \frac{5}{20}log_2(\frac{5}{20}) =$ 1.8394 GainRatio(NúmeroPasajeros) = $\frac{0.4845}{1.8394}$ = 0.2634 Se Calcula $Gini_{N\'umeroPasajeros \in (Dos, Tres)} = \frac{6}{20} Gini(Evaluaci\'on1) + \frac{14}{20} Gini(Evaluaci\'on2)$ = $\frac{6}{20} (1 - \left(\frac{2}{6}\right)^2 - \left(\frac{4}{6}\right)^2) + \frac{14}{20} (1 - \left(\frac{13}{14}\right)^2 - \left(\frac{1}{14}\right)^2) = 0.2261$

 $Gini_{N\'umeroPasajeross \in (Tres, Cuatro)} = \frac{12}{20} Gini(Evaluaci\'on1) +$ Se Calcula $\frac{8}{20}$ Gini(Evaluación2) $= \frac{12}{20} \left(1 - \left(\frac{10}{12} \right)^2 - \left(\frac{2}{12} \right)^2 \right) + \frac{8}{20} \left(1 - \left(\frac{5}{8} \right)^2 - \left(\frac{3}{8} \right)^2 \right) = 0.3541$ Se Calcula $Gini_{N\'umeroPasajeros \in (Dos, Cuatro)} = \frac{12}{20} Gini(Evaluaci\'on1) + \frac{8}{20} Gini(Evaluaci\'on2)$ Como son el número de elementos el número de pasajeros cuando son dos y tres, el resultado es el mismo, o sea Se Calcula $Gini_{N\'umeroPasajeros \in (Tres, M\'as)} = \frac{8}{20} Gini(Evaluaci\'on1) + \frac{13}{20} Gini(Evaluaci\'on2)$ Como son el número de elementos el número de pasajeros cuando son dos y tres, el resultado es el mismo $SplitInfoTamañoCajuela = -\frac{6}{20}log_2(\frac{6}{20}) - \frac{7}{20}log_2(\frac{7}{20}) - \frac{7}{20}log_2(\frac{7}{20}) = 1.5812$ $GainRatio(TamañoCajuela) = \frac{0.0322}{1.5812} = 0.0203$ $Gini_{Tama\~noCajuela \in (Peque\~na, Mediana)} = \frac{14}{20} Gini(Evaluaci\'on1) +$ $\frac{6}{20}$ Gini(Evaluación2) $= \frac{14}{20} \left(1 - \left(\frac{10}{14} \right)^2 - \left(\frac{4}{14} \right)^2 \right) + \frac{6}{20} \left(1 - \left(\frac{5}{6} \right)^2 - \left(\frac{1}{6} \right)^2 \right) = 0.3736$ $Gini_{Tama\~noCajuela \in (Mediana, Grande)} = \frac{13}{20} Gini(Evaluaci\'on1) +$ $\frac{7}{20}$ Gini(Evaluación2) $= \frac{13}{20} \left(1 - \left(\frac{11}{13} \right)^2 - \left(\frac{2}{13} \right)^2 \right) + \frac{7}{20} \left(1 - \left(\frac{4}{7} \right)^2 - \left(\frac{3}{7} \right)^2 \right) = 0.3406$ $Gini_{TamañoCajuela \in (Pequeña, Grande)} = \frac{13}{20} Gini(Evaluación1) +$ $\frac{7}{30}$ Gini(Evaluación2) $=\frac{13}{20}\left(1-\left(\frac{9}{12}\right)^2-\left(\frac{4}{12}\right)^2\right)+\frac{7}{20}\left(1-\left(\frac{6}{7}\right)^2-\left(\frac{1}{7}\right)^2\right)=0.3626$ SplitInfoSeguridad = $-\frac{9}{20}log_2(\frac{9}{20}) - \frac{7}{20}log_2(\frac{7}{20}) - \frac{4}{20}log_2(\frac{4}{20}) = 1.5128$ GainRatio(Seguridad) = $\frac{0.0324}{1.5128} = 0.0214$ Se Calcula $Gini_{Seguridad \in (Baja, Mediana)} = \frac{11}{20} Gini(Evaluación1) + \frac{9}{20} Gini(Evaluación2)$

SplitInfoSeguridad = $-\frac{9}{20}log_2(\frac{9}{20}) - \frac{7}{20}log_2(\frac{7}{20}) - \frac{4}{20}log_2(\frac{4}{20}) = 1.5128$ GainRatio(Seguridad) = $\frac{0.0324}{1.5128}$ = 0.0214 Se Calcula Gini_{Seguridad∈(Baja, Mediana)} = $\frac{11}{20}Gini(Evaluación1) + \frac{9}{20}Gini(Evaluación2)$ = $\frac{11}{20}(1 - (\frac{8}{11})^2 - (\frac{3}{11})^2) + \frac{9}{20}(1 - (\frac{7}{9})^2 - (\frac{2}{9})^2) = 0.3737$ Se Calcula Gini_{Seguridad∈(Mediana, Alta)} = $\frac{16}{20}Gini(Evaluación1) + \frac{4}{20}Gini(Evaluación2)$ = $\frac{16}{20}(1 - (\frac{13}{16})^2 - (\frac{3}{16})^2) + \frac{4}{20}(1 - (\frac{2}{4})^2 - (\frac{2}{4})^2) = 0.3437$ Se Calcula Gini_{Seguridad∈(Baja, Alta)} = $\frac{13}{20}Gini(Evaluación1) + \frac{7}{20}Gini(Evaluación2)$ = $\frac{13}{20}(1 - (\frac{9}{13})^2 - (\frac{4}{13})^2) + \frac{7}{20}(1 - (\frac{6}{7})^2 - (\frac{1}{7})^2) = 0.3626$

Ejercicio 5

Plantee un conjunto de datos, con 15 registros y 5 atributos, cuyo atributo objetivo sea dicotómico y aplique las actividades que realizó en el ejercicio número 4 de esta guía.

El conjunto de datos que planteamos es el siguiente:

4	Α	В	С	D	E	F
1	pecioCompra	costoManto	noLuces	noPasajeros	tamMotor	evalMoto
2	vhigh	med	dos	dos	small	acc
3	med	vhigh	una	tres	big	unacc
4	high	low	dos	uno	med	acc
5	med	low	una	tres	small	acc
6	vhigh	high	tres	uno	big	acc
7	med	vhigh	una	tres	med	unacc
8	med	low	una	uno	big	acc
9	high	med	dos	uno	small	acc
10	med	high	tres	tres	med	unacc
11	med	vhigh	una	uno	big	unacc
12	vhigh	low	dos	dos	small	acc
13	high	high	tres	uno	med	unacc
14	med	low	una	tres	big	acc
15	vhigh	med	dos	dos	small	acc
16	high	low	tres	dos	big	unacc

Aplicación del proceso del cálculo de medidas de evaluación al conjunto de datos de motos (accesible, no accesible)

Paso 1: Cálculo de la entropía total

Evaluación Moto	
Accesible	No accesible
9	6

E(Total) =
$$-\frac{9}{15}log_2\left(\frac{9}{15}\right) - \frac{6}{15}log_2\left(\frac{6}{15}\right) = 0.9709$$

Paso 2: Dividir el conjunto de datos en los diversos atributos

Atributo objetivo	Evaluación Motos
Atributo	Dominio
Precio Compra	Muy alto
	Alto
	Medio
Costo Mantenimiento	Muy alto
	Alto
	Medio
	Bajo
Número Luces	Uno
	Dos
	Tres
Número Pasajeros	Uno
	Dos
	Tres
Tamaño Motor	Pequeña
	Mediana
	Grande

Paso 3: Se calcula la entropía de cada rama y se suman proporcionalmente para calcular la entropía total

Para Precio Compra

Precio Compra	Precio Compra Muy alto	Accesible	4
		No accesible	0
	Alto	Accesible	2
		No accesible	2
	Mediano	Accesible	3
		No accesible	4

$$\begin{split} & \text{E(Muy alto)} = -\frac{4}{15}log_2\left(\frac{4}{15}\right) - \frac{0}{15}log_2\left(\frac{0}{15}\right) = 0.5085 \\ & \text{E(Alto)} = -\frac{2}{15}log_2\left(\frac{2}{15}\right) - \frac{2}{15}log_2\left(\frac{2}{15}\right) = 0.7751 \\ & \text{E(Mediano)} = -\frac{3}{15}log_2\left(\frac{3}{15}\right) - \frac{4}{15}log_2\left(\frac{4}{15}\right) = 0.9728 \end{split}$$

E(Precio Compra, Evaluación) =
$$\frac{4}{15}(0.5085) + \frac{4}{15}(0.7751) + \frac{7}{15}(0.9728) = 0.7962$$

Para Costo Mantenimiento

i ara costo mantenimiento			
Costo	Muy alto	Accesible	0
Mantenimiento		No accesible	3
	Alto	Accesible	1
		No accesible	2
	Medio	Accesible	3
		No accesible	0
	Bajo	Accesible	5
		No accesible	1

$$\begin{split} & \text{E(Muy alto)} = -\frac{0}{15}log_2\left(\frac{0}{15}\right) - \frac{3}{15}log_2\left(\frac{3}{15}\right) = 0.4643 \\ & \text{E(Alto)} = -\frac{1}{15}log_2\left(\frac{1}{15}\right) - \frac{2}{15}log_2\left(\frac{2}{15}\right) = 0.6480 \\ & \text{E(Mediano)} = -\frac{3}{15}log_2\left(\frac{3}{15}\right) - \frac{4}{15}log_2\left(\frac{4}{15}\right) = 0.9728 \\ & \text{E(Bajo)} = -\frac{3}{15}log_2\left(\frac{3}{15}\right) - \frac{0}{15}log_2\left(\frac{0}{15}\right) = 0.4643 \\ & \text{E(Precio Mantenimiento, Evaluación)} = \frac{3}{15}(0.4643) + \frac{3}{15}(0.6480) + \frac{3}{15}(0.9728) + \frac{6}{15}(0.4643) = 0.6027 \end{split}$$

Para Número Luces

Número Luces	Una	Accesible	4
		No accesible	2
	Dos	Accesible	5
		No accesible	0
	Tres	Accesible	1
		No accesible	3

$$\begin{split} & \text{E(Uno)} = -\frac{4}{15}log_2\left(\frac{4}{15}\right) - \frac{2}{15}log_2\left(\frac{2}{15}\right) = 0.8960 \\ & \text{E(Dos)} = -\frac{5}{15}log_2\left(\frac{5}{15}\right) - \frac{0}{15}log_2\left(\frac{0}{15}\right) = 0.5283 \\ & \text{E(Tres)} = -\frac{1}{15}log_2\left(\frac{1}{15}\right) - \frac{3}{15}log_2\left(\frac{3}{15}\right) = 0.7248 \\ & \text{E(Número Luces, Evaluación)} = \frac{6}{15}(0.8960) + \frac{5}{15}(0.5283) + \frac{4}{15}(0.7248) = 0.7277 \end{split}$$

Para Número Pasaieros

Número Pasajeros	Uno	Accesible	4
		No accesible	2
	Dos	Accesible	3
		No accesible	1
	Tres	Accesible	2
		No accesible	3

E(Uno)=
$$-\frac{4}{15}log_2\left(\frac{4}{15}\right) - \frac{2}{15}log_2\left(\frac{2}{15}\right) = 0.8960$$

E(Dos)= $-\frac{3}{15}log_2\left(\frac{3}{15}\right) - \frac{1}{15}log_2\left(\frac{1}{15}\right) = 0.7248$

E(Tres)=
$$-\frac{2}{15}log_2\left(\frac{2}{15}\right) - \frac{3}{15}log_2\left(\frac{3}{15}\right) = 0.8519$$

E(Número Pasajeros, Evaluación) = $\frac{6}{15}(0.8960) + \frac{4}{15}(0.7248) + \frac{5}{15}(0.8519) = 0.8356$

- 00100 10011100110 111000	_		
Tamaño Motor	Grande	Accesible	3
		No accesible	3
	Mediana	Accesible	1
		No accesible	3
	Pequeña	Accesible	5
		No accesible	0

$$\begin{split} & \text{E(Grande)} = -\frac{3}{15}log_2\left(\frac{3}{15}\right) - \frac{3}{15}log_2\left(\frac{3}{15}\right) = 0.9287 \\ & \text{E(Mediano)} = -\frac{1}{15}log_2\left(\frac{1}{15}\right) - \frac{3}{15}log_2\left(\frac{3}{15}\right) = 0.7248 \\ & \text{E(Pequeño)} = -\frac{5}{15}log_2\left(\frac{5}{15}\right) - \frac{0}{15}log_2\left(\frac{0}{15}\right) = 0.5283 \\ & \text{E(Tamaño Motor, Evaluación)} = \frac{6}{15}(0.9287) + \frac{4}{15}(0.7248) + \frac{5}{15}(0.5283) = 0.7408 \end{split}$$

Paso 4: Se calcula la ganancia de información

Para Precio Compra:

Gain(PrecioCompra, Evaluación) = 0.9709 - 0.7962 = 0.1747

• Para Costo Mantenimiento:

Gain(Costo Mantenimiento, Evaluación) = 0.9709 - 0.6027 = 0.3682

• Para Número Luces:

Gain(Número Puertas, Evaluación) = 0.9709 - 0.7277 = 0.2432

• Para Número Pasajeros:

Gain(Número Pasajeros, Evaluación) = 0.9709 - 0.8356 = 0.1353

Para Tamaño Motor:

Gain(Tamaño Cajuela, Evaluación) = 0.9709 - 0.7408 = 0.2301

Paso 5: Elección del nodo de Decisión.

Variable	Ganancia	(Accesible, No Accesible)
Precio de Compra	0.1747	Muy Alto = (4, 0) Alto = (2, 2) Medio = (3, 4)
Costo	0.3682	Muy Alto = $(0, 3)$ Alto = $\{1, 2\}$ Medio = $(3, 0)$
Mantenimiento		Bajo = (5, 1)
Número Luces	0.2432	Uno = (4, 2) Dos = (5, 0) Tres = (1, 3)
Número Pasajeros	0.1353	Uno = (4, 2) Dos = (3, 1) Tres = (2, 3)

Tamaño Motor	0.2301	Grande = (3, 3) Mediana = (1, 3) Pequeña = (5,
		0)

Costo Mantenimiento es la variable que da una mayor ganancia

Se calcula SplitInfo de Mantenimiento:

SplitInfoMantenimimento =
$$-\frac{3}{15}log_2\left(\frac{3}{15}\right) - \frac{3}{15}log_2\left(\frac{3}{15}\right) - \frac{3}{15}log_2\left(\frac{3}{15}\right) - \frac{6}{15}log_2\left(\frac{6}{15}\right) = 19219$$

 $GainRatio(Costo\ Mantenimiento) = \frac{0.3682}{1.9219} = 0.3993$

Se calcula Gini Index con respecto a la variable Evaluación

Gini(Evaluación) =
$$1 - \left(\frac{9}{15}\right)^2 - \left(\frac{6}{15}\right)^2 = 0.48$$

Se Calcula
$$Gini_{CostoMantenimiento \in (Bajo, Medio)} = \frac{17}{15} Gini(Evaluación1) + \frac{3}{15} Gini(Evaluación2)$$

$$= \frac{17}{20} \left(1 - \left(\frac{15}{17}\right)^2 - \left(\frac{2}{17}\right)^2\right) + \frac{3}{20} \left(1 - \left(\frac{3}{3}\right)^2 - \left(\frac{0}{3}\right)^2\right) = 0.1764$$

Se Calcula
$$Gini_{CostoMantenimiento \in (Bajo, Alto)} = \frac{8}{20} Gini(Evaluación1) + \frac{12}{20} Gini(Evaluación2)$$

$$\frac{12}{20}$$
 Gini(Evaluación2)

$$= \frac{8}{20} \left(1 - \left(\frac{3}{8}\right)^2 - \left(\frac{5}{8}\right)^2\right) + \frac{12}{20} \left(1 - \left(\frac{12}{12}\right)^2 - \left(\frac{0}{12}\right)^2\right) = 0.1875$$

Se Calcula
$$Gini_{CostoMantenimiento \in (Medio, Muy alto)} = \frac{15}{20} Gini(Evaluación1) + \frac{5}{20} Gini(Evaluación2)$$

$$= \frac{15}{20} \left(1 - \left(\frac{12}{15}\right)^2 - \left(\frac{3}{15}\right)^2\right) + \frac{5}{20} \left(1 - \left(\frac{3}{5}\right)^2 - \left(\frac{2}{5}\right)^2\right) = 0.36$$

Proponemos realizar este proceso para los demás atributos, para ver si obtenemos valores menores de Gini Index

$$SplitInfoPrecioCompra = -\frac{10}{20}log_2(\frac{10}{20}) - \frac{7}{20}log_2(\frac{7}{20}) - \frac{3}{20}log_2\left(\frac{3}{20}\right) = 1.4406$$

$$GainRatio(PrecioCompra) = \frac{0.0497}{1.4406} = 0.0344$$

$$GainRatio(PrecioCompra) = \frac{0.0497}{1.4406} = 0.0344$$

Se Calcula
$$Gini_{PrecioComptra \in (Medio, Alto)} = \frac{10}{20}Gini(Evaluación1) + \frac{10}{20}Gini(Evaluación2)$$

$$= \frac{10}{20} \left(1 - \left(\frac{8}{10}\right)^2 - \left(\frac{2}{10}\right)^2\right) + \frac{10}{20} \left(1 - \left(\frac{7}{10}\right)^2 - \left(\frac{3}{10}\right)^2\right) = 0.37$$

Se
$$Calcula$$
 $Gini_{PrecioCompra \in (Medio, Muy alto)} = \frac{13}{20}Gini(Evaluación1) + \frac{7}{2}Gini(Evaluación2)$

$$\frac{7}{20}$$
Gini(Evaluación2)

$$= \frac{13}{20} \left(1 - \left(\frac{10}{13} \right)^2 - \left(\frac{3}{13} \right)^2 \right) + \frac{7}{20} \left(1 - \left(\frac{5}{7} \right)^2 - \left(\frac{2}{7} \right)^2 \right) = 0.3736$$

Se Calcula
$$Gini_{PrecioCompra \in (Alto, Muy alto)} = \frac{17}{20}Gini(Evaluación1) + \frac{3}{20}Gini(Evaluación2)$$

$$= \frac{17}{20} \left(1 - \left(\frac{12}{17}\right)^2 - \left(\frac{5}{17}\right)^2\right) + \frac{3}{20} \left(1 - \left(\frac{3}{3}\right)^2 - \left(\frac{0}{3}\right)^2\right) = 0.3529$$

$$SplitInfoNumeroLuces = -\frac{8}{20}log_2(\frac{8}{20}) - \frac{4}{20}log_2(\frac{4}{20}) - \frac{5}{20}log_2\left(\frac{5}{20}\right) - \frac{3}{20}log_2\left(\frac{3}{20}\right) = 1.9037$$

$$GainRatio(NumeroLuces) = \frac{0.1517}{1.9037} = 0.0796$$

Se Calcula
$$Gini_{N\'umeroPuertas \in (Dos, Tres)} = \frac{12}{20} Gini(Evaluaci\'on1) + \frac{8}{20} Gini(Evaluaci\'on2)$$

$$= \frac{12}{20} (1 - \left(\frac{8}{12}\right)^2 - \left(\frac{4}{12}\right)^2) + \frac{8}{20} (1 - \left(\frac{7}{8}\right)^2 - \left(\frac{1}{8}\right)^2) = 0.3541$$

$$Se \ Calcula \ Gini_{N\'umeroPuertas \in (Tres, \ Cuatro)} = \frac{9}{20} \ Gini(Evaluaci\'on1) + \frac{11}{20} \ Gini(Evaluaci\'on2)$$

$$= \frac{9}{20} (1 - \left(\frac{7}{9}\right)^2 - \left(\frac{2}{9}\right)^2) + \frac{11}{20} (1 - \left(\frac{8}{11}\right)^2 - \left(\frac{3}{11}\right)^2) = 0.3737$$

$$Se \ Calcula \ Gini_{N\'umeroPuertas \in (Cuatro, \ cinco\ o\ m\'as)} = \frac{8}{20} \ Gini(Evaluaci\'on1) + \frac{12}{20} \ Gini(Evaluaci\'on2)$$

$$= \frac{8}{20} (1 - \left(\frac{7}{8}\right)^2 - \left(\frac{1}{8}\right)^2) + \frac{12}{20} (1 - \left(\frac{8}{12}\right)^2 - \left(\frac{4}{12}\right)^2) = 0.3541$$

$$Se \ Calcula \ Gini_{N\'umeroPuertas \in (Dos, \ cinco\ o\ m\'as)} = \frac{11}{20} \ Gini(Evaluaci\'on1) + \frac{9}{20} \ Gini(Evaluaci\'on2)$$

$$= \frac{11}{20} (1 - \left(\frac{8}{11}\right)^2 - \left(\frac{3}{11}\right)^2) + \frac{9}{20} (1 - \left(\frac{7}{9}\right)^2 - \left(\frac{2}{9}\right)^2) = 0.3737$$

$$Se \ Calcula \ Gini_{N\'umeroPuertas \in (Dos, \ Cuatro)} = \frac{13}{20} \ Gini(Evaluaci\'on1) + \frac{7}{20} \ Gini(Evaluaci\'on2)$$

$$= \frac{13}{20} (1 - \left(\frac{9}{11}\right)^2 - \left(\frac{4}{14}\right)^2) + \frac{7}{20} (1 - \left(\frac{6}{9}\right)^2 - \left(\frac{1}{1}\right)^2) = 0.3626$$

$$Se \ Calcula \ Gini_{N\'umeroPuertas \in (Tres, \ Cinco)} = \frac{7}{20} \ Gini(Evaluaci\'on1) + \frac{13}{20} \ Gini(Evaluaci\'on2)$$

$$= \frac{7}{20} (1 - \left(\frac{6}{7}\right)^2 - \left(\frac{1}{1}\right)^2) + \frac{13}{20} (1 - \left(\frac{9}{13}\right)^2 - \left(\frac{4}{13}\right)^2) = 0.3626$$

$$SplitInfoN\'umeroPasajeros = -\frac{3}{20} log_2(\frac{3}{20}) - \frac{3}{20} log_2(\frac{3}{20}) - \frac{9}{20} log_2(\frac{9}{20}) - \frac{5}{20} log_2(\frac{5}{20}) = 1.8394$$

$$GainRatio(N\'umeroPasajeros) = \frac{0.4845}{1.8394} = 0.2634$$

$$Se \ Calcula \ Gini_{N\'umeroPasajeros} = (0.2634)$$

Se Calcula $Gini_{N\'umeroPasajeros \in (Dos, Tres)} = \frac{6}{20} Gini(Evaluaci\'on1) + \frac{14}{20} Gini(Evaluaci\'on2)$ = $\frac{6}{20} (1 - (\frac{2}{6})^2 - (\frac{4}{6})^2) + \frac{14}{20} (1 - (\frac{13}{14})^2 - (\frac{1}{14})^2) = 0.2261$

$$= \frac{6}{20} \left(1 - \left(\frac{2}{6}\right)^2 - \left(\frac{4}{6}\right)^2\right) + \frac{14}{20} \left(1 - \left(\frac{13}{14}\right)^2 - \left(\frac{1}{14}\right)^2\right) = 0.2262$$

 $Gini_{N\'umeroPasajeross \in (Tres, Cuatro)} = \frac{12}{20} Gini(Evaluaci\'on1) +$

 $\frac{8}{20}$ Gini(Evaluación2)

$$= \frac{12}{20} \left(1 - \left(\frac{10}{12} \right)^2 - \left(\frac{2}{12} \right)^2 \right) + \frac{8}{20} \left(1 - \left(\frac{5}{8} \right)^2 - \left(\frac{3}{8} \right)^2 \right) = 0.3541$$

 $Gini_{N\'umeroPasajeros \in (Dos, Cuatro)} = \frac{12}{20} Gini(Evaluaci\'on1) +$ Se Calcula

 $\frac{8}{20}$ Gini(Evaluación2)

Como son el número de elementos el número de pasajeros cuando son dos y tres, el resultado es el mismo, o sea

Se Calcula $Gini_{N\'umeroPasajeros \in (Tres, M\'as)} = \frac{8}{20} Gini(Evaluaci\'on1) + \frac{13}{20} Gini(Evaluaci\'on2)$

Como son el número de elementos el número de pasajeros cuando son dos y tres, el resultado es el mismo

$$SplitInfoTama\~noMotor = -\frac{6}{20}log_2(\frac{6}{20}) - \frac{7}{20}log_2(\frac{7}{20}) - \frac{7}{20}log_2(\frac{7}{20}) = 1.5812$$

$$GainRatio(Tama\~noMotor) = \frac{0.0322}{1.5812} = 0.0203$$

 $Gini_{Tama\~noCajuela \in (Peque\~na, Mediana)} = \frac{14}{20} Gini(Evaluaci\'on1) + \frac{14}{20} Gini$ Se Calcula $\frac{6}{20}$ Gini(Evaluación2)

$$= \frac{14}{20} \left(1 - \left(\frac{10}{14}\right)^2 - \left(\frac{4}{14}\right)^2\right) + \frac{6}{20} \left(1 - \left(\frac{5}{6}\right)^2 - \left(\frac{1}{6}\right)^2\right) = 0.3736$$

$$\begin{array}{lll} Se & Calcula & Gini_{Tama\~noCajuela\in(Mediana,\ Grande)} = \frac{13}{20}Gini(Evaluaci\'on1) + \\ & \frac{7}{20}Gini(Evaluaci\'on2) \\ & = \frac{13}{20}(1-\left(\frac{11}{13}\right)^2-\left(\frac{2}{13}\right)^2) + \frac{7}{20}(1-\left(\frac{4}{7}\right)^2-\left(\frac{3}{7}\right)^2) = 0.3406 \\ Se & Calcula & Gini_{Tama\~noCajuela\in(Peque\~na,\ Grande)} = \frac{13}{20}Gini(Evaluaci\'on1) + \\ & \frac{7}{20}Gini(Evaluaci\'on2) \\ & = \frac{13}{20}(1-\left(\frac{9}{13}\right)^2-\left(\frac{4}{13}\right)^2) + \frac{7}{20}(1-\left(\frac{6}{7}\right)^2-\left(\frac{1}{7}\right)^2) = 0.3626 \\ \end{array}$$

Ejercicio 6

Suponga que tiene la siguiente matriz de confusión de la evaluación de carros.

Original/Predicción	Accesible	No accesible
Accesible	40	5
No accesible	2	3
	42	8

Realizamos una modificación al orden de los datos, para una mejor visualización.

Original/Predicción	No accesible	Accesible
No accesible	3	2
Accesible	5	40
	8	42

Calcule y explique las diversas medidas que puede generar con base en los datos de la matriz de confusión.

Medida	Cálculo	Explicación
Negativo verdadero	3	A son los ejemplares clasificados como negativos de forma adecuada (son negativos originalmente). Negativo verdadero
Positivo falso	2	B son los ejemplares clasificados como positivos de forma incorrecta, ya que son negativos originalmente. Falso positivo
Negativo falso	5	C son los ejemplares clasificados como negativos de forma incorrecta, ya que son positivos originalmente. Falso negativo
Positivo verdadero	40	D son los ejemplares clasificados como positivos de forma adecuada (son positivos originalmente). Positivo verdadero

Tasa de exactitud	$\frac{(\Box + \Box)}{(\Box + \Box + \Box + \Box)}$ $= \frac{(3+40)}{(3+2+5+40)}$ $= 86\%$	Tasa de exactitud Las clasificaciones que hizo de forma correcta
Tasa de error	$ \frac{(\Box + \Box)}{(\Box + \Box + \Box + \Box)} \\ = \frac{(2+5)}{(3+2+5+40)} \\ = 14\% $	La tasa de error Las clasificaciones que hizo de forma equivocada
Precisión	$\frac{d}{b+d} = \frac{40}{2+40} = 95\%$	Precisión La precisión (<i>Precision</i>) mide cuántos ejemplos clasificados como clase "positiva" son realmente "positivos".
Sensibilidad (<i>Recall</i>)	$\frac{d}{c+d} = \frac{40}{5+40} = 88\%$	Recall La sensibilidad (Sensitivity) (también conocida como recuerdo (recall)) evalúa qué tan bien el clasificador puede reconocer muestras positivas
Tasa de positivos falsos	$\frac{b}{a+b} = \frac{2}{3+2} = 40\%$	Nos revela el porcentaje de falsos positivos según la matriz de confusión.
Tasa de negativos falsos	$\frac{c}{c+d} = \frac{5}{5+40} = 11\%$	Nos revela el porcentaje de falsos negativos según la matriz de confusión.
Especificidad	$\frac{a}{a+b} = \frac{3}{3+2} = 60\%$	Especificidad La medida de especificidad (specificity) mide que tan bien el clasificador puede reconocer las muestras negativas