

INSTITUTO POLITÉCNICO NACIONAL
ESCUELA SUPERIOR DE CÓMPUTO
MATERIAL EDUCATIVO PARA LA UNIDAD DE APRENDIZAJE DE MINERÍA DE DATOS.

2022-2

Grupo 5CDM1

PRACTICA DE ÁRBOLES DE DECISIÓN

Nombres:

Angeles Lomeli Felipe Alberto

García Rodríguez Diana Itzel

De Luna Ocampo Yanina

Medina Barreras Daniel Ivan

1. Descripción del conjunto de datos.

Autores del conjunto de datos:

Donante: Ronny Kohavi and Barry Becker. Data Mining and Visualization. Silicon Graphics. e-mail: ronnyk '@' live.com for questions.

Enlace de acceso: <https://archive.ics.uci.edu/ml/datasets/adult>

2. Objetivo de la práctica.

Realizar un árbol de decisión para predecir si una persona gana más de 50K al año.

3. Diccionario de datos.

Construya el diccionario de datos considerando la siguiente estructura.

No	Nombre	Tipo	Dominio
1	age	Numeric	–
2	workclass	String	Private, Self-emp-not-inc, Self-emp-inc, Federal-gov, Local-gov, State-gov, Without-pay, Never-worked.
3	fnlwgt	Numeric	
4	education	String	Bachelors, Some-college, 11th, HS-grad, Prof-school, Assoc-acdm, Assoc-voc, 9th, 7th-8th, 12th, Masters, 1st-4th,

			10th, Doctorate, 5th-6th, Preschool.
5	education-num	Numeric	continuous
6	marital-status	String	Married-civ-spouse, Divorced, Never-married, Separated, Widowed, Married-spouse-absent, Married-AF-spouse.
7	occupation	String	Tech-support, Craft-repair, Other-service, Sales, Exec-managerial, Prof-specialty, Handlers-cleaners, Machine-op-inspct, Adm-clerical, Farming-fishing, Transport-moving, Priv-house-serv, Protective-serv, Armed-Forces.
8	relationship	String	Wife, Own-child, Husband, Not-in-family, Other-relative, Unmarried.
9	race	String	White, Asian-Pac-Islander, Amer-Indian-Eskimo, Other, Black.
10	sex	String	Female, Male.
11	capital-gain	Numeric	continuous
12	capital-loss	Numeric	continuous
13	hours-per-week	Numeric	continuous
14	native-country	String	United-States, Cambodia, England, Puerto-Rico, Canada, Germany, Outlying-US(Guam-USVI-etc), India, Japan, Greece, South, China, Cuba, Iran, Honduras, Philippines, Italy, Poland, Jamaica, Vietnam, Mexico, Portugal, Ireland, France, Dominican-Republic, Laos, Ecuador, Taiwan, Haiti, Columbia, Hungary, Guatemala, Nicaragua, Scotland, Thailand, Yugoslavia, El-Salvador, Trinidad&Tobago, Peru,

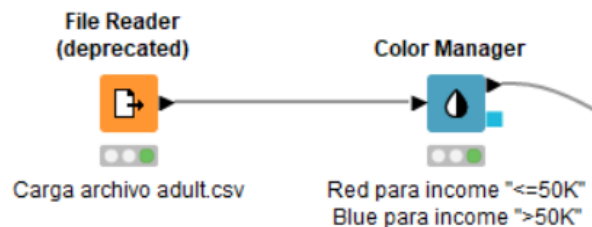
			Hong, Holand-Netherlands.
--	--	--	------------------------------

4. Resultados

Presente los resultados considerando lo siguiente:

1. Realice y describa los resultados de cinco consultas descriptivas en el conjunto de datos
2. Presente propiedades estadísticas del conjunto de datos
3. Describa las medidas generadas a partir de la matriz de confusión (archivo anexo)
4. Analice este comportamiento en función la cantidad de elementos de cada tipo que existen en el conjunto de datos
5. Anexe el modelo y las reglas generadas.

Use los siguientes nodos:






1. Realice y describa los resultados de cinco consultas descriptivas en el conjunto de datos

1.- Nivel Educativo por Sexo

<input type="checkbox"/>	RowID ↑↓	sex ↑↓	Unique concatenate with count(education) ↑↓
<input type="checkbox"/>	Row0	Female	Bachelors(1619), Masters(536), 9th(144), HS-grad(3390), Some-college(2806), Prof-school(92), 11th(432), Doctorate(86), Assoc-acdm(421), 10th(295), Assoc-voc(500), 1st-4th(46), Preschool(16), 5th-6th(84), 7th-8th(160), 12th(144)
<input type="checkbox"/>	Row1	Male	Bachelors(3736), HS-grad(7111), 11th(743), Some-college(4485), Assoc-acdm(646), Assoc-voc(882), 7th-8th(486), Doctorate(327), 9th(370), 5th-6th(249), 10th(638), Masters(1187), Prof-school(484), 1st-4th(122), 12th(289), Preschool(35)




(porcentaje por sexo, tratamiento de datos: separar y agrupar niveles educativos)

2.- Raza-nivel educativo

<input type="checkbox"/>	RowID 	race 	Unique concatenate with count(education) 
<input type="checkbox"/>	Row0	Amer-Indian-Eskimo	7th-8th(9), Some-college(79), 10th(16), HS-grad(119), Assoc-acdm(8), 11th(14), Bachelors(21), Assoc-voc(19), Prof-school(2), 9th(5), Masters(5), 5th-6th(2), 12th(5), Doctorate(3), 1st-4th(4)
<input type="checkbox"/>	Row1	Asian-Pac-Islander	Bachelors(289), Assoc-voc(38), Some-college(208), HS-grad(226), Masters(88), Doctorate(28), 11th(21), Assoc-acdm(29), Prof-school(41), 7th-8th(11), 9th(9), 12th(9), 5th-6th(18), 1st-4th(5), 10th(13), Preschool(6)
<input type="checkbox"/>	Row2	Black	11th(153), Bachelors(330), 9th(89), Some-college(746), Assoc-acdm(107), HS-grad(1174), Assoc-voc(112), 10th(133), 12th(70), 5th-6th(21), 1st-4th(16), 7th-8th(56), Masters(86), Prof-school(15), Doctorate(11), Preschool(5)
<input type="checkbox"/>	Row3	Other	Some-college(51), 11th(10), 7th-8th(17), Bachelors(33), HS-grad(78), 10th(9), Assoc-voc(6), 9th(8), Masters(7), 12th(14), 1st-4th(9), Assoc-acdm(8), 5th-6th(13), Prof-school(4), Doctorate(2), Preschool(2)
<input type="checkbox"/>	Row4	White	Bachelors(4682), HS-grad(8904), Masters(1537), 11th(977), Doctorate(369), Assoc-acdm(915), Some-college(6207), 9th(403), Assoc-voc(1207), Prof-school(514), 5th-6th(279), 7th-8th(553), 10th(762), 1st-4th(134), Preschool(38), 12th(335)

Porcentajes y tratamiento de datos

3.-Raza-Relaciones

<input type="checkbox"/>	RowID 	race 	Unique concatenate with count(relationship) 
<input type="checkbox"/>	Row0	Amer-Indian-Eskimo	Husband(92), Not-in-family(81), Own-child(48), Unmarried(58), Other-relative(13), Wife(19)
<input type="checkbox"/>	Row1	Asian-Pac-Islander	Husband(410), Wife(69), Other-relative(82), Unmarried(91), Not-in-family(214), Own-child(173)
<input type="checkbox"/>	Row2	Black	Husband(671), Wife(153), Not-in-family(812), Unmarried(769), Own-child(555), Other-relative(164)
<input type="checkbox"/>	Row3	Other	Wife(16), Other-relative(28), Unmarried(37), Husband(80), Not-in-family(73), Own-child(37)
<input type="checkbox"/>	Row4	White	Not-in-family(7125), Husband(11940), Wife(1311), Own-child(4255), Unmarried(2491), Other-relative(694)

Porcentajes

4.-Raza-Sexo-Pais

<input type="checkbox"/>	Row0	Amer-Indian-Eskimo	Female	United-States(114), Mexico(2), South(2), Columbia(1)
<input type="checkbox"/>	Row1	Amer-Indian-Eskimo	Male	Mexico(6), United-States(182), Germany(1), Puerto-Rico(1), Philippines(1), Hong(1)
<input type="checkbox"/>	Row2	Asian-Pac-Islander	Female	?(22), Philippines(73), United-States(115), England(1), Laos(8), South(28), India(10), China(20), Hong(5), Japan(10), Vietnam(22), Cambodia(2), Taiwan(14), Poland(1), Thailand(10), Germany(1), Canada(1), Portugal(1), Greece(1), Haiti(1)
<input type="checkbox"/>	Row3	Asian-Pac-Islander	Male	India(75), ?(61), South(49), United-States(177), Cambodia(16), Thailand(6), Taiwan(34), Philippines(115), China(53), Japan(28), Vietnam(43), Laos(10), Iran(6), Trinidad&Tobago(2), Germany(2), Hong(12), Puerto-Rico(1), Mexico(1), Dominican-Republic(1), Ireland(1)
<input type="checkbox"/>	Row4	Black	Female	Cuba(1), Jamaica(41), United-States(1429), Japan(2), Outlying-US(Guam-USVI-etc)(2), Haiti(19), Dominican-Republic(7), ?(32), Trinidad&Tobago(10), Germany(3), Honduras(1), England(1), Puerto-Rico(4), Cambodia(1), El-Salvador(1), France(1)
<input type="checkbox"/>	Row5	Black	Male	United-States(1403), Germany(5), Haiti(24), ?(64), Jamaica(34), Trinidad&Tobago(6), England(7), Outlying-US(Guam-USVI-etc)(4), Dominican-Republic(5), Mexico(4), Nicaragua(2), Puerto-Rico(5), India(2), Japan(1), Cuba(2), Philippines(1)
<input type="checkbox"/>	Row6	Other	Female	United-States(56), Puerto-Rico(9), Germany(1), Mexico(9), Jamaica(1), Columbia(4), ?(7), Dominican-Republic(8), Guatemala(3), Ecuador(3), El-Salvador(3), Taiwan(1), Trinidad&Tobago(1), Nicaragua(1), Cuba(1), Japan(1)
<input type="checkbox"/>	Row7	Other	Male	United-States(73), Puerto-Rico(12), ?(11), Dominican-Republic(10), Mexico(31), Guatemala(1), Ecuador(6), India(5), Nicaragua(3), El-Salvador(1), Cuba(1), Canada(1), Columbia(3), Peru(1), Iran(2), Japan(1)
<input type="checkbox"/>	Row8	White	Female	United-States(7968), ?(102), Honduras(6), England(30), Mexico(135), Columbia(19), Germany(55), France(11), Poland(16), Cuba(38), Italy(21), Guatemala(17), Dominican-Republic(20), El-Salvador(31), Canada(38), Peru(14), Puerto-Rico(39), Nicaragua(11), Portugal(11), Ireland(7), Iran(8), Ecuador(6), Outlying-US(Guam-USVI-etc)(5), Yugoslavia(3), Jamaica(1), Greece(4), Thailand(1), Scotland(5), China(1), Japan(7), Hong(1), Hungary(6), Vietnam(1), Holand-Netherlands(1), India(1)
<input type="checkbox"/>	Row9	White	Male	United-States(17653), Puerto-Rico(43), ?(284), Mexico(455), Cuba(52), Canada(81), Iran(27), Italy(52), Poland(41), Ecuador(13), Portugal(25), Dominican-Republic(19), El-Salvador(70), Guatemala(43), England(51), Philippines(8), Germany(69), Japan(12), Yugoslavia(13), Jamaica(4), Scotland(7), Greece(24), Nicaragua(17), Columbia(32), Ireland(16), Peru(16), France(17), Honduras(6), India(7), Hungary(7), Taiwan(2), Thailand(1), South(1), Outlying-US(Guam-USVI-etc)(3), China(1), Vietnam(1), Hong(1)

5.- Dinero-raza-sexo

<input type="checkbox"/>	RowID	race	sex	Mean(fnlwgt)
<input type="checkbox"/>	Row0	Amer-Indian-Eskimo	Female	112950.731092437
<input type="checkbox"/>	Row1	Amer-Indian-Eskimo	Male	125715.36458333336
<input type="checkbox"/>	Row2	Asian-Pac-Islander	Female	147452.07514450865
<input type="checkbox"/>	Row3	Asian-Pac-Islander	Male	166175.86580086604
<input type="checkbox"/>	Row4	Black	Female	212971.38778135023
<input type="checkbox"/>	Row5	Black	Male	242920.64499681254
<input type="checkbox"/>	Row6	Other	Female	172519.64220183485
<input type="checkbox"/>	Row7	Other	Male	213679.10493827166
<input type="checkbox"/>	Row8	White	Female	183549.9669058082
<input type="checkbox"/>	Row9	White	Male	188987.38614790735

2. Presente propiedades estadísticas del conjunto de datos

Row ID	D Min	D Max	D Mean	D Std. de...	D Variance	D Skewness	D Kurtosis	D Overall ...	I No. mis...	I No. NaNs	I No. +cos	I No. -cos	D Median	I Row co...	Histogram
age	17	90	38.582	13.64	186.061	0.559	-0.166	1,256,257	0	0	0	0	?	32561	
frlvgit	12,285	1,484,705	189,778.367	105,549.978	11,140,797,...	1.447	6.219	6,179,373,392	0	0	0	0	?	32561	
education-num	1	16	10.081	2.573	6.619	-0.312	0.623	328,237	0	0	0	0	?	32561	
capital-gain	0	99,999	1,077.649	7,385.292	54,542,539,...	11.954	154.799	35,089,324	0	0	0	0	?	32561	
capital-loss	0	4,356	87.304	402.96	162,376.938	4.595	20.377	2,842,700	0	0	0	0	?	32561	
hours-per-week	1	99	40.437	12.347	152.459	0.228	2.917	1,316,684	0	0	0	0	?	32561	

3. Describa las medidas generadas a partir de la matriz de confusión (archivo anexo)

Matriz de Confusión:

Row ID	I <=50K	I >50K
<=50K	6627	734
>50K	874	1440

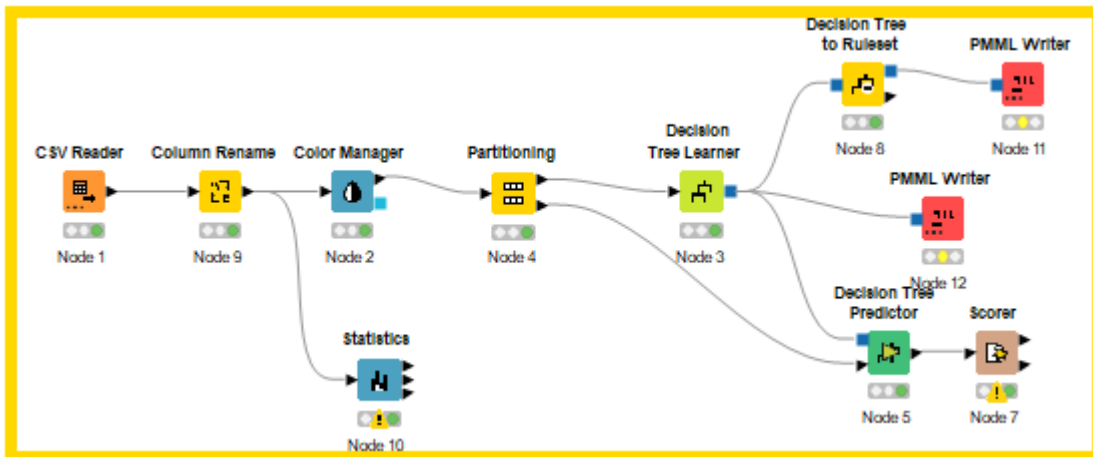
Row ID	I TruePositives	I FalsePositives	I TrueNegatives	I FalseNegatives	D Recall	D Precision	D Sensitivity	D Specificity	D F-measure	D Accuracy	D Cohen's kappa
<=50K	6627	874	1440	734	0.9	0.883	0.9	0.622	0.892	?	?
>50K	1440	734	6627	874	0.622	0.662	0.622	0.9	0.642	?	?
Overall	?	?	?	?	?	?	?	?	?	0.834	0.534

4. Analice este comportamiento en función la cantidad de elementos de cada tipo que existen en el conjunto de datos

Medida	Cálculo	Interpretación
Positivo verdadero	6627 =d	El valor real es negativo y la prueba predijo también que el resultado era negativo
Falso Positivo	734=b	El valor real es negativo, y la prueba predijo que el resultado es positivo.
Falso Negativo	874=c	El valor real es positivo, y la prueba predijo que el resultado es negativo
Verdaderos Negativos	1440=a	El valor real es positivo y la prueba predijo también que era positivo.
Tasa de exactitud	$(6627+1440)/(6627+734+874+1440) = 8067/9675 = 0.83$	La cantidad de predicciones positivas que fueron correctas fue del 83%
Tasa de error	$(874+734)/(6627+734+874+1440) = 1608/9675 = 0.16$	La cantidad de predicciones que fueron incorrectas es del 16%
Precisión	$1440/(734+1440) = 0.66$	El porcentaje de casos positivos detectados fue 66%
Sensibilidad (Recall)	$1440/(874+1440) = 0.62$	En este caso la sensibilidad apenas es capaz de detectar correctamente, con un porcentaje del 62%
Tasa de positivos falsos	$734/(6627+734) = 0.09$	La probabilidad de que se produzca una falsa alarma: que se dé un resultado positivo cuando el valor verdadero sea negativo es del 9%
Tasa de negativos falsos	$874/(1440+874) = 0.37$	La probabilidad de que la prueba pase por alto un verdadero positivo es del 37%
Especificidad	$1440/(734+1440) = 0.66$	La probabilidad de los casos negativos que el algoritmo ha clasificado correctamente es del 66%

5. Anexe el modelo y las reglas generadas.

Modelo



Consultas

