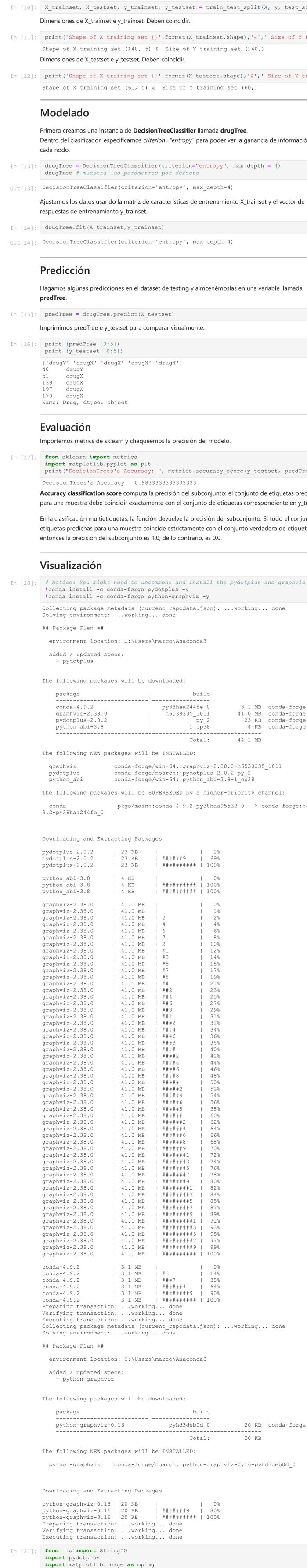
	• Desarrollar un modelo de clasificación usando el algoritmo de árboles de decisión Utilizaremos el algoritmo de clasificación de árboles de decisión para construir un modelo a partir de datos históricos de pacientes y su respuesta a diferentes medicaciones. Luego utilizará el árbol entrenado para predecir la clase de un paciente desconocido, o para encontrar la drogra apropiada para un nuevo paciente.					
	Tabla de contenido					
	 Acerca del dataset Descargando los datos Pre-procesamiento Estableciendo el árbol de decisión Modelando Predicción Evaluación Visualización 					
	 Importamos las siguientes librerías: numpy (as np) pandas DecisionTreeClassifier from sklearn.tree 					
In [1]:	<pre>import numpy as np import pandas as pd from sklearn.tree import DecisionTreeClassifier</pre>					
	Acerca del dataset Imagine que es un investigador médico que reúne datos para un estudio. Ha coleccionado datos acerca de un conjunto de pacientes; todos ellos sufrieron la misma enfermedad. Durante el curso de su tratamiento,					
	cada paciente respondió a 1 de 5 medicamentos: Drug A, Drug B, Drug c, Drug x e y. Parte de su trabajo es construir un modelo para encontrar qué droga es apropiada para un futuro paciente que tenga la misma enfermedad. El conjunto de características es Edad, Sexo, presión sanguínea y colesterol, y el objetivo es la droga a la que responde cada paciente. Es un ejemplo de clasificador multiclase. Se usa la parte de					
	de un paciente desconocido o p	entrenamiento del dataset para construir un árbol de decisión y luego éste se utiliza para predecir la clase de un paciente desconocido o para prescribirle la droga apropiada a un nuevo paciente. Descargando los datos				
In [2]:	<pre>import urllib.request url = 'https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBMDevelope filename = 'drug200.csv'</pre>					
Out[2]:	urllib.request.urlretrie ('drug200.csv', <http.cl: dataframe="" en="" leemos="" los="" par<="" td="" un=""><td>ient.HTTPM</td><td></td><td>db7109cd30>)</td><td></td></http.cl:>	ient.HTTPM		db7109cd30>)		
<pre>In [3]: Out[3]:</pre>	my_data[0:5] Age Sex BP Cholest	terol Na_to_	K Drug	=",")		
	1 47 M LOW H	HIGH 13.09	5 drugY 3 drugC 4 drugC 8 drugX			
In [5]:			3 drugY			
Out[5]:		0				
	 Utilizando my_data declararemos las siguientes variables: X como la matriz de características y como el vector respuesta (objetivo) 					
In [6]:	<pre>Removemos la columna que col X = my_data[['Age', 'Sex X[0:5]</pre>			•		
Out[6]:	array([[23, 'F', 'HIGH', [47, 'M', 'LOW', [47, 'M', 'LOW', [28, 'F', 'NORMAL [61, 'F', 'LOW',	'HIGH', 13 'HIGH', 10 ', 'HIGH', 'HIGH', 18	.093], .113999999999 7.7979999999 .043]], dtype	99999], =object)	ortunadamente los	
In [7]:	árboles de decisión de Sklearn r valores numéricos. pandas.get _ from sklearn import prep	no trabajan o _dummies() processing	con variables cat convierte variab	egóricas, así que debo	emos convertirlas a	
	<pre>le_sex = preprocessing.LabelEncoder() le_sex.fit(['F','M']) X[:,1] = le_sex.transform(X[:,1]) le_BP = preprocessing.LabelEncoder()</pre>					
	<pre>le_BP.fit(['LOW', 'NORMAL', 'HIGH']) X[:,2] = le_BP.transform(X[:,2]) le_Chol = preprocessing.LabelEncoder() le_Chol.fit(['NORMAL', 'HIGH'])</pre>					
Out[7]:	<pre>X[:,3] = le_Chol.transform(X[:,3]) X[0:5] array([[23, 0, 0, 0, 25.355],</pre>					
	[47, 1, 1, 0, 13.093], [47, 1, 1, 0, 10.1139999999999], [28, 0, 2, 0, 7.797999999999], [61, 0, 1, 0, 18.043]], dtype=object) Variable objetivo					
<pre>In [8]: Out[8]:</pre>	y[0:5] 0 drugY 1 drugC					
	<pre>drugC drugX drugY Name: Drug, dtype: object</pre>					
	Estableciendo el árbol de decisión Utilizaremostrain/test split					
In [9]:	from sklearn.model_selection import train_test_split Ahora, train_test_split devolverá 4 parámetros diferentes. Los llamaremos: X_trainset, X_testset, y_trainset, y_testset					
	train_test_split necesita los parámetros: X, y, test_size y random_state.					
In [10]:	random_state asegura que obte	Y son los arreglos antes de la división. test_size representa el cociente del dataset de testing. dom_state asegura que obtengamos las mismas divisionestrainset, X_testset, y_trainset, y_testset = train_test_split(X, y, test_size=0.3, z				
In [11]:	Dimensiones de X_trainset e y_trainset. Deben coincidir. print('Shape of X training set {}'.format(X_trainset.shape),'&',' Size of Y training Shape of X training set (140, 5) & Size of Y training set (140,)					
In [12]:	Dimensiones de X_testset e y_te print('Shape of X traini Shape of X training set	ing set {}	'.format(X_te		' Size of Y training s	
	Modelado Primero creamos una instancia de DecisionTreeClassifier llamada drugTree.					
In [13]:	Dentro del clasificador, especificamos criterion="entropy" para poder ver la ganancia de información en cada nodo. drugTree = DecisionTreeClassifier(criterion="entropy", max_depth = 4) drugTree # muestra los parámetros por defecto					
Out[13]:	DecisionTreeClassifier(criterion='entropy', max_depth=4) Ajustamos los datos usando la matriz de características de entrenamiento X_trainset y el vector de respuestas de entrenamiento y_trainset. drugTree.fit(X_trainset, y_trainset) DecisionTreeClassifier(criterion='entropy', max_depth=4) Predicción Hagamos algunas predicciones en el dataset de testing y almcenémoslas en una variable llamada predTree.					
<pre>In [14]: Out[14]:</pre>						
In [15]:	<pre>predTree = drugTree.predict(X_testset)</pre> Imprimimos predTree e y_testset para comparar visualmente.					
In [16]:	<pre>print (predTree [0:5]) print (y_testset [0:5]) ['drugY' 'drugX' 'drugX' 'drugX'] 40</pre>					
	197 drugX 170 drugX Name: Drug, dtype: object					
In [17]:	Evaluación Importemos metrics de sklearn g from sklearn import metr import matplotlib.pyplot	rics	os la precisión d	el modelo.		
	print ("DecisionTrees's Accuracy: ", metrics.accuracy_score (y_testset, predTree)) DecisionTrees's Accuracy: 0.9833333333333 Accuracy classification score computa la precisión del subconjunto: el conjunto de etiquetas predichas para una muestra debe coincidir exactamente con el conjunto de etiquetas correspondiente en y_true.					
	En la clasificación multietiquetas, la función devuelve la precisión del subconjunto. Si todo el conjunto de etiquetas predichas para una muestra coincide estrictamente con el conjunto verdadero de etiquetas, entonces la precisión del subconjunto es 1.0; de lo contrario, es 0.0.					
	Visualización # Notice: You might need to uncomment and install the pydotplus and graphviz librarie.					
In [20]:	<pre>!conda install -c conda-forge pydotplus -y !conda install -c conda-forge python-graphviz -y Collecting package metadata (current_repodata.json):working done Solving environment:working done</pre>					
	<pre>## Package Plan ## environment location: C:\Users\marco\Anaconda3 added / updated specs: - pydotplus</pre>					
	The following packages was package	 	bui			
	conda-4.9.2 graphviz-2.38.0 pydotplus-2.0.2 python_abi-3.8		py38haa244fe h6538335_10 py 1_cp	2_0 3.1 MI 11 41.0 MI 7_2 23 KI	3 conda-forge 3 conda-forge 3 conda-forge	
	The following NEW package graphviz conc pydotplus conc python_abi conc	da-forge/w	<pre>INSTALLED: in-64::graphv</pre>	riz-2.38.0-h653833		
	pydotplus conda-forge/noarch::pydotplus-2.0.2-py_2 conda-forge/win-64::python_abi-3.8-1_cp38 The following packages will be SUPERSEDED by a higher-priority channel: conda pkgs/main::conda-4.9.2-py38haa95532_0> conda-forge::conda-4.9.2-py38haa244fe_0					
	Downloading and Extracting pydotplus-2.0.2 23 pydotplus-2.0.2 23	3 KB	 #####9			
	python_abi-3.8 4 python_abi-3.8 4 python_abi-3.8 4 python_abi-3.8 4	3 KB KB KB KB	######### ######### ##########	100% 0% 100% 100%		
	graphviz-2.38.0 43 graphviz-2.38.0 43 graphviz-2.38.0 43 graphviz-2.38.0 43	1.0 MB 1.0 MB 1.0 MB 1.0 MB 1.0 MB	2 4 6 7 1	1% 2% 4% 6% 8%		
	graphviz-2.38.0 43 graphviz-2.38.0 43 graphviz-2.38.0 43 graphviz-2.38.0 43 graphviz-2.38.0 43 graphviz-2.38.0 43	1.0 MB 1.0 MB 1.0 MB 1.0 MB 1.0 MB 1.0 MB	9	10% 12% 14% 15% 17% 19%		
	graphviz-2.38.0 43 graphviz-2.38.0 43 graphviz-2.38.0 43 graphviz-2.38.0 43 graphviz-2.38.0 43	1.0 MB 1.0 MB 1.0 MB 1.0 MB 1.0 MB 1.0 MB	##	21% 23% 25% 27% 29% 31%		

Árboles de decisión

Objetivos



from sklearn import tree

graph.write_png(filename) img = mpimg.imread(filename) plt.figure(figsize=(100, 200))

entropy = 0.0

samples = 16

value = [16, 0, 0, 0, 0]

class = drugA

featureNames = my_data.columns[0:5]

plt.imshow(img,interpolation='nearest')

Out[22]: <matplotlib.image.AxesImage at 0x1db72fec220>

targetNames = my_data["Drug"].unique().tolist()

graph = pydotplus.graph_from_dot_data(dot_data.getvalue())

 $Age \leq 50.5$ entropy = 0.975

samples = 27

value = [16, 11, 0, 0, 0]

class = drugA

entropy = 0.0

samples = 11

value = [0, 11, 0, 0, 0]

class = drugB

out=tree.export_graphviz(drugTree, feature_names=featureNames, out_file=dot_data, class

True

 $BP \le 0.5$

entropy = 1.832

samples = 71

value = [16, 11, 11, 33, 0]

class = drugX

entropy = 0.0

samples = 11

value = [0, 0, 11, 0, 0]class = drugC

 $Na_to_K \le 14.615$ entropy = 1.929samples = 140value = [16, 11, 11, 33, 69] class = drugY

Cholesterol ≤ 0.5

entropy = 0.811

samples = 44

value = [0, 0, 11, 33, 0]class = drugX

 $BP \leq 1.5\,$

entropy = 0.998

samples = 21

value = [0, 0, 11, 10, 0]

class = drugC

False

entropy = 0.0

samples = 69

value = [0, 0, 0, 0, 69]

class = drugY

entropy = 0.0

samples = 10value = [0, 0, 0, 10, 0]

class = drugX

entropy = 0.0

samples = 23

value = [0, 0, 0, 23, 0]

class = drugX

%matplotlib inline

dot_data = StringIO() filename = "drugtree.png"

In [22]: