



NOMBRE DE LA PRÁCTICA	Preparacion del	l datadet		NO.	1
ASIGNATURA:	SIMULACION	CARRERA:	INGENIERÍA EI SISTEMAS COMPUTACIONALES	DE LA	2

NOMBRE DEL ALUMNO:

GRUPO: 3502

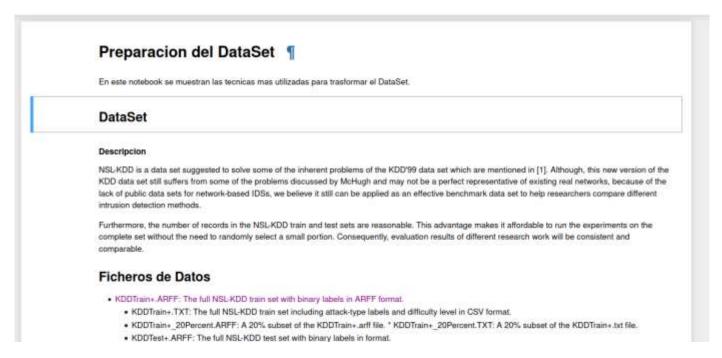
I. Competencia(s) específica(s):

•

- II. Lugar de realización de la práctica (laboratorio, taller, aula u otro):
 - Salon de clase
- III. Material empleado:

equipo de computo

PRACTICA



La "Preparación del DataSet" se refiere al proceso de acondicionamiento y transformación de un conjunto de datos antes de utilizarlo para entrenar un modelo de aprendizaje automático. Este proceso es esencial para garantizar que los datos sean adecuados y estén en la forma correcta





para que los algoritmos de machine learning puedan aprender patrones y realizar predicciones de manera efectiva.

Ficheros de Datos . KDDTrain+ ARFF: The full NSL-KDD train set with binary labels in ARFF format. . KDDTrain+.TXT: The full NSL-KDD train set including attack-type labels and difficulty level in CSV format. KDDTrain+_20Percent.ARFF: A 20% subset of the KDDTrain+.artf file. * KDDTrain+_20Percent.TXT: A 20% subset of the KDDTrain+.txt file. . KDDTest+.ARFF: The full NSL-KDD test set with binary labels in format. . KDDTest+.TXT: The full NSL-KDD test set including attack- labels and difficulty level in CSV format. . KDDTest-21.ARFF; A subset of the KDDTest+.arff file does not include records with difficulty level of 21 out of 21. . KDDTest-21.TXT; A subset of the KDDTest+.txt which does not include records with difficulty level of 21 out of 21. Descarga de los ficheros de datos. https://www.unb.ca/cic/datasets/index.html Referencias adicionales sobre el conjunto de datos M. Tavallase, E. Bagheri, W. Lu, and A. Ghorbani, "A Detailed Analysis of the KDD CUP 99 Data Set," Submitted to Second IEEE Symposium on Computational Intelligence for Security and Defense Applications (CISDA), 2009. Imports. In [1]: 1 import arff import pandas as pd import numpy as np 4 from sklearn.model_selection import train_test_split

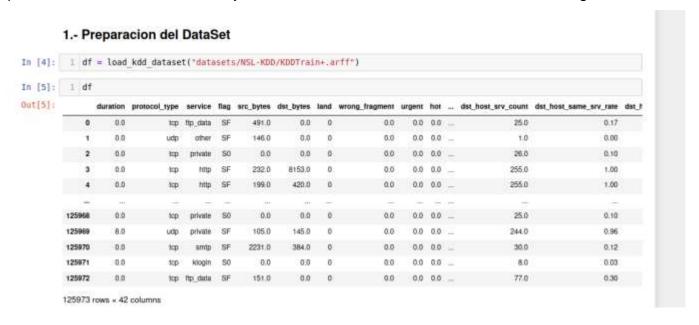
Los ficheros de datos están relacionados con el conjunto de datos NSL-KDD, que se utiliza comúnmente en el ámbito de la detección de intrusiones en redes

MANUAL DE PRACTICAS



```
Imports.
In [1]: 1 import arff
               import pandas as pd
               import numpy as np
            4 from sklearn, model selection import train test split
          Funciones Auxiliares.
           1 def load kdd dataset(data path):
                    ""Lectura del DataSet NSL-KDD.""
with open(data path, 'r') as train_set:
                          dataset = arff.load(train_set)
                          attributes = [attr[0] for attr in dataset["attributes"]]
                          return pd.DataFrame(dataset["data"], columns=attributes)
In [3]:
               # Construcion de una funcion que realize un particionado completo
def train val test split(df,rstate=42,shuffle=True,stratify=None):
                    sstrat = df[stratify] if stratify else None
train set, test set = train test split(
                    df. test size=0.4, random state=rstate,shuffle=shuffle, stratify =sstrat)
strat=test_set[stratify]if stratify else None
val_set, test_set=train_test_split(
                     test set, test size = 0.5, random state = rstate, shuffle=shuffle, stratify =strat)
                     return(train set, val set, test set)
              December of the del Detector
```

código importa varias bibliotecas en Python, cada una de las cuales desempeña un papel importante en el análisis de datos y la construcción de modelos de machine learning. Así mismo el código importan varias bibliotecas en Python, cada una de las cuales desempeña un papel importante en el análisis de datos y la construcción de modelos de machine learning.



 $funci\'on\ llamada\ \verb|load_kdd_dataset|\ para\ cargar\ un\ conjunto\ de\ datos\ desde\ un\ archivo\ ARFF.$





Realiza la división de un conjunto de datos (df) en conjuntos de entrenamiento, validación y prueba utilizando la función train val test split.

Asi mismo se separan las características de entrada (x_train) y las características de salida (y_train) del conjunto de entrenamiento ($train_tset$).

```
In [11]: 1 # comprobar si existe algun atributo con valores nulos
           x train.isna().any()
Out[11]: duration
                                         False
         protocol type
                                         False
         service
                                         False
         flag
                                         False
         src_bytes
dst_bytes
                                          True
                                          True
         land
                                         False
         wrong fragment
                                         False
                                         False
         hot.
                                         False
         num_failed_logins
                                         False
         logged in
                                         False
         num compromised
                                         False
         root shell
                                         False
         su_attempted
                                         False
         num_root
                                         False
         num file creations
                                         False
         num shells
                                         False
         num access files
                                         False
         num outbound cmds
                                         False
         is host login
                                         False
         is guest login
                                         False
         count
                                         False
         sry count
                                         False
         serror rate
                                         False
         srv serror rate
                                         False
                                         False
```

Se verifica si hay algún atributo (columna) en el conjunto de características de entrada x_train que contiene valores nulos.

MANUAL DE PRACTICAS



```
dtype: bool
In [12]: 1 # Selections to file que contienen valores nulos.
2 files_valores_nulos = [x_train.isnult().any(axis =1)]
            3 filas valores nulos
Out[12]: [113467
                         True
            31899
                       False
            108116
                        True
            89913
                       False
            106319
                       False
            64559
                        False
            67272
                         True
            32452
                        False
            112657
                        False
            99030
                        False
           Length: 75583, dtype: bool]
          Opcion 1: Eliminar las filas con valores nulos.
In [13]: 1 # Copiar el DataSet para no alterar el original
            2 x_train_copy = x_train.copy()
In [14]: 1 # Eliminar las filas con valores nulos
            z x train copy.dropna(subset=["src bytes", "dst bytes"], inplace = True)
            3 x train copy
Out[14]:
                   duration protocol_type service flag src_bytes dat_bytes land wrong_fragment urgent hot ... dat_host_count dat_host_srv_count dat_host_srv_count dat_host_srv_count
```

seleccionar las filas en el DataFrame x_{train} que contienen al menos un valor nulo asi mismo creamos una copia del DataFrame x_{train} y asignándolo a la variable x_{train}_{copy} . Esto se hace para asegurarse de que cualquier modificación realizada en x_{train}_{copy} no afecte al DataFrame original x_{train}

utilizamos el método dropna para eliminar las filas que contienen valores nulos en las columnas "src bytes" y "dst bytes" de x train copy





```
1 # Contar el numero de filas eliminadas
2 print(" El numero de filas eliminadas es ", len(x_train)-len(x_train_copy))
In [15]:
            El numero de filas eliminadas es 9886
           Opcion 2: eliminar los atributos con valores nulos.
In [16]:
             1 # Copiar el DataSet para no alterar el original
             2 x train copy = x train.copy()
            1 # Eliminar los atributos con valores nulos
In [17]:
                x_train_copy.drop(["src_bytes", "dst_bytes"], axis = 1, inplace = True)
             3 x train copy
Out[17]:
                   duration protocol_type service flag land wrong_fragment urgent hot num_failed_logins logged_in ... dst_host_count dst_host_srv_count dst_h
            113467
                        0.0
                                                                       0.0
                                                                              00 00
                                                                                                    0.0
                                                                                                                                9.0
                                                                                                                                                 255.0
                                      ton
                                             http
             31899
                                          private
                                                                              0.0 0.0
                                                                                                                              255.0
                                      tcp
            108116
                        0.0
                                             http
                                                                       0.0
                                                                              0.0 0.0
                                                                                                    0.0
                                                                                                                               39.0
                                                                                                                                                 255.0
             89913
                                                                                                                0 ...
                        0.0
                                                                       0.0
                                                                              0.0 0.0
                                                                                                    0.0
                                                                                                                              255.0
                                                                                                                                                  15.0
                                      tcp
                                          private
                                                  SO
            106319
                                                         0
                        0.0
                                                  SF
                                                                       0.0
                                                                              0.0 0.0
                                                                                                    0.0
                                                                                                                0 ...
                                                                                                                                2.0
                                            eco i
                                                                                                                                                   7.0
                                    icmp
                                           systat
             64559
                        0.0
                                      tcp
                                                  S0
                                                         0
                                                                       0.0
                                                                              0.0 0.0
                                                                                                    0.0
                                                                                                                0 ...
                                                                                                                              255.0
                                                                                                                                                  20.0
                                                                                                                              119.0
             67272
                        0.0
                                      tcp
                                             http
                                                  SF
                                                                       0.0
                                                                              0.0 0.0
                                                                                                    0.0
                                                                                                                                                 255.0
             32452
                                                                              0.0 0.0
                                                                                                                               111.0
                                                                                                                                                 155.0
            112657
                                                                       0.0
                                                                              00 00
                                                                                                    0.0
                                                                                                                              255 0
                                                                                                                                                 255.0
```

calculando el número de filas eliminadas al restar la longitud del DataFrame original x_train menos la longitud del DataFrame después de la eliminación x_train_copy

sto crea una copia del DataFrame x_train y asigna esa copia a la variable x_train_copy. Puedes realizar modificaciones y experimentos en x_train_copy sin afectar al conjunto de datos original x_train se eliminan las columnas "src_bytes" y "dst_bytes" del DataFrame x_train_copy utilizando el método drop

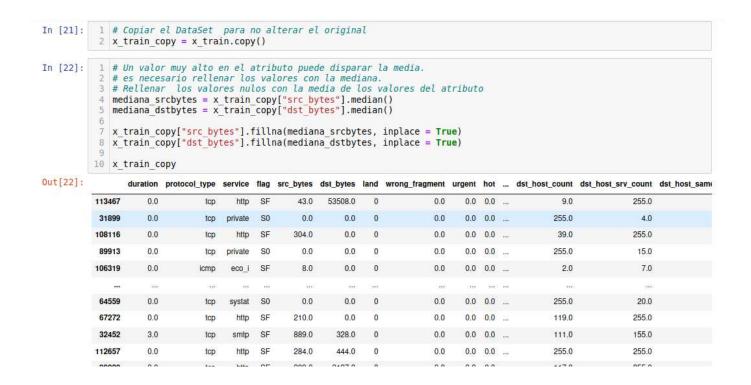
MANUAL DE PRACTICAS



Opcion 3: Rellenar los valores nulos con un valor deternimado.

```
In [19]:
                                     1 # Copiar el DataSet para no alterar el original
                                     2 x train copy = x train.copy()
In [20]:
                                             # Rellenar los valores nulos con la media de los valores del atributo
                                             media_srcbytes = x_train_copy["src_bytes"].mean()
                                             media dstbytes = x train copy["dst bytes"].mean()
                                             x train copy["src bytes"].fillna(media srcbytes, inplace = True)
                                             x_train_copy["dst_bytes"].fillna(media_dstbytes, inplace = True)
                                           x train copy
Out[20]:
                                                       duration protocol_type service flag
                                                                                                                                                                 src_bytes
                                                                                                                                                                                                     dst_bytes land wrong_fragment urgent hot ... dst_host_count dst_host_srv_count dst_host_s
                                  113467
                                                                    0.0
                                                                                                                              http
                                                                                                                                             SF
                                                                                                                                                        66914.530762 53508.000000
                                                                                                                                                                                                                                        0
                                                                                                                                                                                                                                                                               0.0
                                                                                                                                                                                                                                                                                                    0.0 0.0 ...
                                                                                                                                                                                                                                                                                                                                                                  9.0
                                                                                                                                                                                                                                                                                                                                                                                                                 255.0
                                    31899
                                                                    0.0
                                                                                                                                                                   0.000000
                                                                                                                                                                                                                                        0
                                                                                                                                                                                                                                                                               0.0
                                                                                                                                                                                                                                                                                                    0.0 0.0 ...
                                                                                                                                                                                                                                                                                                                                                             255.0
                                                                                                                                                                                                                                                                                                                                                                                                                      4.0
                                                                                                                                             SO
                                                                                                                                                                                                        0.000000
                                                                                                         top
                                                                                                                       private
                                  108116
                                                                    0.0
                                                                                                                              http
                                                                                                                                             SF
                                                                                                                                                              304.000000
                                                                                                                                                                                               9181.334754
                                                                                                                                                                                                                                                                               0.0
                                                                                                                                                                                                                                                                                                    0.0 0.0 ...
                                                                                                                                                                                                                                                                                                                                                                39.0
                                                                                                                                                                                                                                                                                                                                                                                                                 255.0
                                                                                                          tcp
                                     89913
                                                                    0.0
                                                                                                         tcp
                                                                                                                       private
                                                                                                                                             SO
                                                                                                                                                                   0.000000
                                                                                                                                                                                                        0.000000
                                                                                                                                                                                                                                        0
                                                                                                                                                                                                                                                                               0.0
                                                                                                                                                                                                                                                                                                    0.0 0.0 ...
                                                                                                                                                                                                                                                                                                                                                             255.0
                                                                                                                                                                                                                                                                                                                                                                                                                    15.0
                                  106319
                                                                    0.0
                                                                                                                                                                   8.000000
                                                                                                                                                                                                        0.000000
                                                                                                                                                                                                                                                                               0.0
                                                                                                                                                                                                                                                                                                    00 00
                                                                                                                                                                                                                                                                                                                                                                                                                      7.0
                                                                                                       icmp
                                                                                                                          eco i
                                                                                                                                             SE
                                                                                                                                                                                                                                        0
                                                                                                                                                                                                                                                                                                                                                                  20
                                     64559
                                                                    0.0
                                                                                                                          systat
                                                                                                                                             SO
                                                                                                                                                                   0.000000
                                                                                                                                                                                                        0.000000
                                                                                                                                                                                                                                                                               0.0
                                                                                                                                                                                                                                                                                                    0.0 0.0 ...
                                                                                                                                                                                                                                                                                                                                                             255.0
                                                                                                                                                                                                                                                                                                                                                                                                                    20.0
                                                                                                         tcp
                                                                                                                                                                                                                                        0
                                     67272
                                                                    nn
                                                                                                                              http
                                                                                                                                             SE
                                                                                                                                                              210 000000 0181 334754
                                                                                                                                                                                                                                                                               nn
                                                                                                                                                                                                                                                                                                    00 00
                                                                                                                                                                                                                                                                                                                                                              1190
```

utilizando un imputador (como el SimpleImputer de scikit-learn) para rellenar los valores nulos en tu conjunto de datos x train copy



MANUAL DE PRACTICAS



ellenando los valores nulos en las columnas "src_bytes" y "dst_bytes" con las medianas respectivas de esas columnas.

Existen otra alternativa para la opcion 3 que consiste en usar la clase imputer de Sklearn

```
1 # Copiar el DataSet para no alterar el original
          2 x train copy = x train.copy()
In [24]:
         1 from sklearn. impute import SimpleImputer
          3 imputer = SimpleImputer(strategy = "median")
In [25]: 1 # La clase imputer no admite valores categoricos eliminar losatributos categoricos
          2 x_train_copy_num = x_train_copy.select_dtypes(exclude = ['object'])
          3 x_train_copy_num.info()
         <class 'pandas.core.frame.DataFrame'>
         Int64Index: 75583 entries, 113467 to 99030
         Data columns (total 34 columns):
         # Column
                                           Non-Null Count Dtype
         0
             duration
                                           75583 non-null
                                                           float64
             src bytes
          1
                                           73696 non-null
                                                           float64
                                           67572 non-null
             dst_bytes
                                                           float64
          3
             wrong fragment
                                                           float64
                                           75583 non-null
             urgent
                                           75583 non-null
                                                           float64
             hot
                                           75583 non-null
                                                           float64
             num failed logins
                                           75583 non-null
                                                           float64
             num compromised
                                                           float64
                                           75583 non-null
              root shell
          8
                                           75583 non-null
                                                           float64
              su attemnted
                                           75583 non-null
                                                           float64
```

- from sklearn.impute import SimpleImputer: Importa la clase SimpleImputer desde el módulo impute de scikit-learn. Esta clase se utiliza para imputar valores nulos en un conjunto de datos.
- imputer = SimpleImputer(strategy="median"): Crea una instancia de SimpleImputer con la estrategia de imputación establecida como "median". La estrategia "median" significa que los valores nulos se rellenarán con las medianas de las columnas correspondientes.

Se seleccionan las columnas numéricas del conjunto de datos x_train_copy y eliminando aquellas que tienen tipos de datos categóricos ('object').





```
In [26]:
            1 # Se le tiene que proporcionar los atributos numericos para que calcule los valores
             2 imputer.fit(x train copy num)
Out[26]: SimpleImputer(strategy='median')
           In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.
           On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.
In [27]:
            1 #rellenar los valores nulos
              x_train_copy_num_nonan = imputer.transform(x train copy num)
In [28]:
            1 # Transformar el resultado a un DataFrame de pandas
            x train copy = pd.DataFrame(x train copy num nonan, columns = x train copy num.columns)
In [29]: 1 x train copy.head(10)
Out[29]:
              duration src_bytes dst_bytes wrong_fragment urgent hot num_failed_logins num_compromised root_shell su_attempted ... dst_host_count dst_host_:
                           43.0
                                                   0.0
                                                          0.0 0.0
                                                                                                                     0.0
                  0.0
                           0.0
                                     0.0
                                                   0.0
                                                                                                                     0.0 ...
                                                                                                                                    255.0
                                                          0.0 0.0
                                                                               0.0
                                     0.0
           2
                  0.0
                          304.0
                                                   0.0
                                                          00 00
                                                                              0.0
                                                                                               0.0
                                                                                                         0.0
                                                                                                                     0.0
                                                                                                                                     390
           3
                  0.0
                           0.0
                                     0.0
                                                   0.0
                                                          0.0 0.0
                                                                              0.0
                                                                                               0.0
                                                                                                         0.0
                                                                                                                     0.0 ...
                                                                                                                                    255.0
                  0.0
                           8.0
                                     0.0
                                                                               0.0
                                                                                                                                      2.0
                                                   0.0
                                                          0.0 0.0
                                                                                               0.0
                                                                                                         0.0
                                                                                                                     0.0 ...
                                                                                                                     0.0 ...
                  0.0
                           46.0
                                   130 0
                                                   0.0
                                                          00 00
                                                                              0.0
                                                                                                         0.0
                                                                                                                                    255.0
                                   363.0
                                                          00 00
                                                                                                                     0.0
                                                                                                                                    141 0
                                                                                                         0.0
```

método fit del imputador (SimpleImputer) para calcular los valores necesarios basándote en las columnas numéricas del conjunto de datos x_train_copy_num. stamos utilizando x train copy imputed, que es el resultado de la transformación del imputador

APIs de Sklearn

Antes de continuar es necesario hacer una pequeña reseña sobre como funcionan las APIs de sklearn:

- Estimators: Cualquier objeto que pueda estimar cualquier parametro
 - El propio estimator se forma mediante el metodo fit (), que siempre toma un data set como argumento
 - cualquier otro parametro de este metodo es un hipermetro
- Transformers: Son estimadores capaces de trasformar el DataSet (como imputer)
 - La trasformacion se realiza mediante el metodo trasform().
 - Reciben un DataSet como parametro de entrada
- Predictors Son estimadores capaces de realizar predicciones
 - La prediccion se realiza por el metodo predict ().
 - · Reciben un DataSet como entrada
 - Retornan un DataSet con las predicciones
 - Tienen un metodo score() para evaluar el resultado de la prediccion.

4.- Transformacion de atributos categoricos a numericos.

Antes de comenzar, a recuperar el DataSet el DataSet limpio es necesario separar las etiquetas del resto de los datos, no necesariamente se quiere aplicar las mismas trasformaciones en ambos conjuntos.

```
In [30]: 1 # Copiar el DataSet para no alterar el original
2 x train = train_set.drop("class", axis = 1)
3 y train = train_set!"slace" | copy()
```

MANUAL DE PRACTICAS



APIs de scikit-learn" se hace referencia a las interfaces de programación de aplicaciones (API, por sus siglas en inglés) proporcionadas por la biblioteca scikit-learn (sklearn). Scikit-learn es una biblioteca de machine learning en Python que ofrece herramientas simples y eficientes para análisis predictivo de datos.

```
In [32]:
             protocol_type = x_train['protocol_type']
              protocol_type_encoded, categorias = protocol_type.factorize()
          1 #Muestra por pantalla como sehan codificado
             for i in range(10):
                  print(protocol type.iloc[1], "=", protocol type encoded[i])
          tcp = 0
          tcp = 0
          tcp = 0
          tcp = 0
         tcp = 1
tcp = 2
tcp = 0
          tcp = 0
         tcp = 0
          tcp = 0
In [34]: 1 print(categorias)
         Index(['tcp', 'icmp', 'udp'], dtype='object')
```

Se esta utilizando la función factorize de pandas para convertir las etiquetas categóricas en valores numéricos asi mismo se esta mostrando cómo se han codificado las primeras 10 filas de la columna 'protocol type'.

HAVIOTHIAVIONOV ATAILEANAV HIVMIAINO VINVAIN

Ordinal Encoding.

Realiza la misma codificacion que el metodo factorize de pandas

```
In [35]: 1 from sklearn.preprocessing import OrdinalEncoder
             protocol_type = x_train[['protocol type']]
          5 ordinal_encoder = OrdinalEncoder()
          6 protocol_type_encoded = ordinal_encoder.fit_transform(protocol_type)
In [36]:
          1 # Mostrar por pantalla como se han modificado
           2 for i in range(10):
                 print( protocol type['protocol type'].iloc[i], "=", protocol type encoded[i])
         tcp = [1.]
         tcp = [1.]
         tcp = [1.]
         tcp = [1.]
         icmp = [0.]
         udp = [2.]
         tcp = [1.]
         tcp = [1.]
         tcp = [1.]
         tcp = [1.]
```

MANUAL DE PRÁCTICAS



```
In [37]: 1 print(ordinal_encoder.categories_)
    [array(['icmp', 'tcp', 'udp'], dtype=object)]
```

El problema de este tipo de codificacion radica en que ciertos algoritmos de ML que funcionan midiendo la similitud de dos puntos por distancia, van a conciderar que 1 esta mas cerca del 2, que del 3, y en este caso (para estos valores categoricos), no tiene sentido por ello se utiliza otros metodos de categorizacion como por ejemplo, OneHot Encoding

OneHot Encoding

Genera para cada categoria del atributo categorico una matriz binaria que reprecenta el valor

el OrdinalEncoder de scikit-learn para codificar las variables categóricas en tu conjunto de datos x train

```
1 # Convertir Esparse Matrix aun array de Numpy
           protocol type oh.toarray()
Out[39]: array([[0., 1., 0.],
                 [0., 1., 0.],
[0., 1., 0.],
                  [0., 1., 0.],
                 [0., 1., 0.],
[0., 1., 0.]])
In [40]: 1 # Mostrar en pantalla como se ha codificado
           2 for i in range(10):
                  print(protocol type["protocol type"].iloc[i], "=", protocol type oh.toarray()[i])
          tcp = [0. 1. 0.]
          icmp = [1. 0. 0.]
          udp = [0. 0. 1.]
          tcp = [0. 1. 0.]
          tcp = [0. 1. 0.]
          tcp = [0. 1. 0.]
          tcp = [0. 1. 0.]
In [41]: 1 print (ordinal encoder categories )
```

el método toarray () para convertir la matriz dispersa (sparse matrix) protocol_type_oh en un array de NumPy. Asi mismo muestra cómo se ha codificado la columna 'protocol_type' en las primeras 10 filas después de aplicar la codificación one-hot

MANUAL DE PRÁCTICAS



```
In [41]: 1 print (ordinal_encoder.categories_)
[array(['icmp', 'tcp', 'udp'], dtype=object)]
```

En muchas ocaciones al particionar el DataSet o al realizar una prediccion con nuevos ejemplos aparecen nuevos valores para determinadas categorias que produciran un error en la funcion **transfor()**. La claces OneHotEncoding proporciona el parametro **handle_uknown** ya sea para generar un error o ingnorar una caracteristica categorica desconocida que este precente durante la transformacion (el valor predeterminado es lanzar un error). Cuando este parametro se establece en "ignorar" y se encuentra una categoria desconocida durante la transformacion, las columnas codificadas resultantes para esta caracteristica seran todos ceros en la trasformacion inversa una categoria desconocida se denotara como None.

```
In [42]: 1 oh_encoder = OneHotEncoder(handle_unknown = 'ignore')
```

Get Dummes

Get Dummies es un metodo sencillo de utilizar que permite aplicar One-HOt Encoding a un data frame de pandas

OrdinalEncoder, puedes acceder a las categorías originales utilizando el atributo categories
OneHotEncoder con el parámetro handle_unknown='ignore', estás especificando
cómo manejar categorías desconocidas durante la transformación

5.- Escalado del DataSet.

Antes de comenzar, vamos a recuperar el DataSet limpio y separar las etiquetas del resto de los datos no necesariamente se quiere applicar las mismas trasnsformaciones en ambos conjuntos.

```
In [44]: 1 x_train = train_set.drop("class", axis = 1)
2 y_train = train_set["class"].copy()
```

Por norma general los algoritmos de ML no se comportan adecuadamente si los valores de las características que reciben como entrada se encuentran en rangos muy dispares por ello se utilizan distintas tecnicas de escalado. Importante:Tener en cuenta que estos mecanismos de escalado no deben aplicarse sobre etiquetas.

NORMALIZACION: Los valores del atributo se escalan para adquirir un valor entre 0 y 1. ESTANDARIZACION: Los valores del atributo se escalan y reciben un valor similar pero no se encuentra dentro de un rango

Es importante que para provar estos valores se realizen las trasformaciones solo sobre el DataSet de entrenamiento. Despues se aplicaran , sobre el DataSet para testear.

```
In [47]: 1
    from sklearn.preprocessing import RobustScaler
    scale_attrs = x_train[['src_bytes', 'dst_bytes']]
    robust_scaler = RobustScaler()
    x_train_scaled = robust_scaler.fit_transform(scale_attrs)

    x_train_scaled = pd.DataFrame(x_train_scaled, columns = ['src_bytes', 'dst_bytes'])
```

vamos a recuperar el DataSet limpio y separar las etiquetas del resto de los datos no necesariamente se quiere applicar las mismas trasnsformaciones en ambos conjuntos.

RobustScaler de scikit-learn para realizar una transformación robusta de escalado a las columnas 'src_bytes' y 'dst_bytes' de tu conjunto de datos x train





	erc hytee	dst_bytes
_		
	1.324818	
1	-0.160584	0.000000
2	0.948905	1.211429
3	-0.160584	0.000000
4	-0.131387	0.000000
5	0.007299	0.264762
6	6.372263	0.691429
7	2.500000	0.626667
8	0.591241	2.841905
9	1.058394	0.000000
10	37.755474	0.000000
11	-0.160584	0.000000
12	-0.160584	0.000000
13	1.065693	8.158095
14	-0.160584	0.000000
15	0.007299	0.142857
16	3.638686	0.628571

se muestran las primeras 20 filas del DataFrame x_{train_scaled} , utilizando el método head() de pandas

V. Conclusiones

la preparación del conjunto de datos es una fase crítica en el proceso de análisis y modelado de datos. El notebook proporciona una visión detallada de las técnicas más utilizadas para transformar el conjunto de datos, destacando la importancia de abordar problemas como la limpieza, la normalización de características. Al explorar diversas técnicas, desde el manejo de valores atípicos hasta la codificación de variables categóricas, se evidencia la versatilidad y la necesidad de adaptar las estrategias según las características específicas del conjunto de datos. La comprensión profunda de estas técnicas permite mejorar la calidad de los datos, lo que a su vez contribuye a la construcción de modelos más precisos y robustos