



Universidad de Oviedo

ESCUELA POLITÉCNICA DE INGENIERÍA DE GIJÓN

GRADO EN INGENIERÍA EN TECNOLOGÍA Y
SERVICIOS DE TELECOMUNICACIÓN

ÁREA DE INGENIERÍA TELEMÁTICA

Detección del trastorno del espectro autista en niños mediante
técnicas de aprendizaje automático

ANEXOS

D. Ordoñez Prendes Ángel
TUTOR: D. Corcoba Magaña Victor

FECHA: Julio de 2024

Índice de Contenido

A) Código del primer estudio.....	1
B) Código del segundo estudio.....	16



A) Código del primer estudio

Código correspondiente al estudio “Técnicas de 'Deep Learning' sobre métodos analíticos”

```
1      # -*- coding: utf-8 -*-
2      """TFG_et2_Comparativa
3
4      Automatically generated by Colab.
5
6      Original file is located at
7      https://colab.research.google.com/drive/14R2YZ-hd28j94bkOj4Xhi2aggmJulbJ0
8
9      # Instalación e Importación de Librerías
10     """
11
12     import tensorflow as tf
13     import tensorflow_decision_forests as tfdf
14
15     import pandas as pd
16     import numpy as np
17     import csv
18
19     from sklearn.model_selection import train_test_split
20     from sklearn.metrics import confusion_matrix
21     from sklearn import metrics
22     from sklearn.metrics import precision_recall_curve
23     from sklearn.metrics import PrecisionRecallDisplay
24
25     import matplotlib.pyplot as plt
26     import seaborn
27
28     from sklearn.model_selection import KFold
29     import time
30     import psutil
31     import os
32     import tensorflow_decision_forests as tfdf
33
34     print(tf.__version__)
35     print(tfdf.__version__)
```



```
27     """# Funciones de creación y entrenamiento de modelos
28
29     Modelos Secuenciales (En función del optimizador)
30     """
31
32     def sequential_model():
33         model = tf.keras.models.Sequential()
34         model.add(tf.keras.layers.Dense(units=15,
35             activation='sigmoid',input_shape=(15,)))
36         model.add(tf.keras.layers.Dense(units=10, activation='sigmoid'))
37         model.add(tf.keras.layers.Dropout(0.1))
38         model.add(tf.keras.layers.Dense(units=10, activation='relu'))
39         model.add(tf.keras.layers.Dense(units=2, activation='softmax'))
40
41         return model
42
43     def sequential_model_fit_by_optimizer(optimizer_type):
44
45         #Creación del Modelo
46
47         model = sequential_model()
48
49         #Compilación y Entrenamiento
50
51         model.compile(optimizer=optimizer_type,
52             loss='sparse_categorical_crossentropy',
53             metrics=['sparse_categorical_accuracy'])
54         start_time = time.time()
55
56         model_fit = model.fit(dataset_train_data, dataset_train_label,
57             epochs=30, validation_data=(dataset_validation_data ,
58             dataset_validation_label),verbose=0)
59
60         training_time = time.time() - start_time
61
62         #Evaluación
63         Test_loss, Test_accuracy = model.evaluate(dataset_test_data,
64             dataset_test_label)
65         dataset_test_predictions =
66             np.argmax(model.predict(dataset_test_data), axis=-1)
67
68         return model_fit , dataset_test_predictions , Test_loss,
69             Test_accuracy , training_time
```



```
50      """Funciones Secuenciales con Validación Cruzada"""
51      def sequential_model_fit_by_optimizer_kfold(optimizer_type):
52          kfolder = KFold(n_splits=10, shuffle=True)
53          folder_number = 1
54
55          acc_per_fold = []
56          loss_per_fold = []
57
58          data = np.concatenate((dataset_train_data, dataset_test_data),
59                                axis=0)
60          label = np.concatenate((dataset_train_label, dataset_test_label),
61                                 axis=0)
62
63          #Iteración por folder
64
65          for train, test in kfolder.split(data,label):
66
67              model = sequential_model()
68
69              model.compile(optimizer=optimizer_type,
70                            loss='sparse_categorical_crossentropy',
71                            metrics=['sparse_categorical_accuracy'])
72
73              print('-----')
74              print(f'-> Training for folder number {folder_number}')
75              model_fit = model.fit(data[train], label[train],
76                                   epochs=30,verbose=0)
77
78              values = model.evaluate(data[test], label[test])
79              print(f'-> Score for folder {folder_number}:
80                    {model.metrics_names[0]}
81                    of {values[0]}; {model.metrics_names[1]} of {values[1]*100}%')
82
83              acc_per_fold.append(values[1] * 100)
84              loss_per_fold.append(values[0])
85
86              folder_number = folder_number + 1
87
88          # Resumen por folder
89
90          print('-----')
91          print(f'-> Score per fold')
```



```
73     for i in range(0, len(acc_per_fold)):
74         print('-----')
75         print(f'--> Fold {i+1} - Loss: {loss_per_fold[i]} - Accuracy:
              {acc_per_fold[i]}%')

76     print('-----')
77     print('-> Average scores for all folds:')
78     print(f'--> Accuracy: {np.mean(acc_per_fold)} (+-
              {np.std(acc_per_fold)})')

79     print(f'--> Loss: {np.mean(loss_per_fold)}')
80     print('-----\n')
81     return np.mean(acc_per_fold) , np.mean(loss_per_fold)

82     """Modelos de Árboles de Decisión"""

83     def Decision_Tree_model_fit(model_type):

84         #Creación del Modelo

85         if model_type=="GB":
86             model = tfdf.keras.GradientBoostedTreesModel()
87         else:
88             model = tfdf.keras.RandomForestModel()

89         #Compilación y Entrenamiento
90         model.compile(metrics=['accuracy'])
91         start_time = time.time()

92         #import gc
93         model.fit(TF_train, validation_data=(TF_validation))

94         training_time = time.time() - start_time

95         #Evaluación
96         Test_loss, Test_accuracy = model.evaluate(TF_test)
97         dataset_test_predictions_RF = model.predict(TF_test)
98
99         dataset_test_predictions_RF =
              np.round(dataset_test_predictions_RF)

100    return model , dataset_test_predictions_RF , Test_loss,
              Test_accuracy , training_time
```



```
101      """Modelos de Árboles de Decisión con Validación Cruzada"""
102      def Decision_Tree_model_fit_kfold(model_type):
103          kfolder = KFold(n_splits=10, shuffle=True)
104          folder_number = 1
105
106          acc_per_fold = []
107          loss_per_fold = []
108
109          data = np.concatenate((dataset_train_data, dataset_test_data),
110                                axis=0)
111          label = np.concatenate((dataset_train_label, dataset_test_label),
112                                 axis=0)
113
114          #Iteración por folder
115
116          for train, test in kfolder.split(data,label):
117              if model_type=="GB":
118                  model = tfdf.keras.GradientBoostedTreesModel()
119              else:
120                  model = tfdf.keras.RandomForestModel()
121              model.compile(metrics=['accuracy'])
122
123              print('-----')
124              print(f'-> Training for folder number {folder_number}')
125              model_fit = model.fit(data[train], label[train],verbose=0)
126
127              values = model.evaluate(data[test], label[test])
128              print(f'-> Score for folder {folder_number}:
129                    {model.metrics_names[0]}
130                    of {values[0]}; {model.metrics_names[1]} of {values[1]*100}%')
131
132              acc_per_fold.append(values[1] * 100)
133              loss_per_fold.append(values[0])
134
135              folder_number = folder_number + 1
136
137          #Resumen por folder
138
139          print('-----')
140          print('-> Score per fold')
141
142          for i in range(0, len(acc_per_fold)):
143              print('-----')
144              print(f'-> Fold {i+1} - Loss: {loss_per_fold[i]} - Accuracy:
145                    {acc_per_fold[i]}%')
146          print('-----')
```



```
131         print('-> Average scores for all folds:')
132         print(f'--> Accuracy: {np.mean(acc_per_fold)} (+-
              {np.std(acc_per_fold)})')
133         print(f'--> Loss: {np.mean(loss_per_fold)}')
134         print('-----\n')

135         return np.mean(acc_per_fold) , np.mean(loss_per_fold)
136         """# Funciones de Creación de Gráficos"""

137     def generate_conf_matrix(test_label,test_predictions,title):
138
139         conf_matrx=metrics.confusion_matrix(test_label,test_predictions)
140         labels=np.unique(dataset_test_label)
141         plt.figure(figsize=(6,4))
142         seaborn.heatmap(conf_matrx,annot
143                         =True,fmt='d',cmap='Blues',xticklabels=
144                         labels,yticklabels=labels)
145         plt.xlabel('Predictions')
146         plt.ylabel('True Values')
147         plt.title(title)
148         plt.show()

149     def generate_accuracy_loss_graphs(model,title):
150         plt.figure(figsize=(12, 4))
151         plt.subplot(1, 2, 1)
152
153         plt.plot(model.history['sparse_categorical_accuracy'],
154                 label='Accuracy (Training)')
155         plt.plot(model.history['val_sparse_categorical_accuracy'],
156                 label='Accuracy (Validation)')
157         plt.xlabel('Epoch')
158         plt.ylabel('Precisión')
159         plt.title(title+' Accuracy Evolution')
160         plt.legend()
161         plt.grid(True)

162         plt.subplot(1, 2, 2)

163         plt.plot(model.history['loss'], label='Loss (Training)')
164         plt.plot(model.history['val_loss'], label='Loss (Validation)')
165         plt.xlabel('Epoch')
166         plt.ylabel('Loss')
167         plt.title(title+' Loss Evolution')
168         plt.legend()

169         plt.tight_layout()
170         plt.grid(True)
```




```
165         plt.show()
166     def generate_accuracy_graphs(model,title):
167         plt.figure(figsize=(12, 4))
168         plt.plot(model.history['sparse_categorical_accuracy'],
169                 label='Accuracy (Training)')
170
171         plt.plot(model.history['val_sparse_categorical_accuracy'],
172                 label='Accuracy (Validation)')
173         plt.xlabel('Epoch')
174         plt.ylabel('Precisión')
175         plt.title(title+' Accuracy Evolution')
176         plt.legend()
177
178         plt.tight_layout()
179         plt.grid(True)
180         plt.show()
181
182     def generate_loss_graphs(model,title):
183         plt.figure(figsize=(12, 4))
184
185         plt.plot(model.history['loss'], label='Loss (Training)')
186         plt.plot(model.history['val_loss'], label='Loss (Validation)')
187         plt.xlabel('Epoch')
188         plt.ylabel('Loss')
189         plt.title(title+'Loss Evolution')
190         plt.legend()
191
192         plt.tight_layout()
193         plt.grid(True)
194         plt.show()
195
196     def Decision_Tree_generate_accuracy_loss_graphs(model):
197         logs = model.make_inspector().training_logs()
198         plt.figure(figsize=(12, 6))
199         plt.subplot(1, 2, 1)
200         plt.plot([log.num_trees for log in logs], [log.evaluation.accuracy for
201             log in logs])
202         plt.xlabel("Number of trees")
203         plt.ylabel("Accuracy (out-of-bag)")
204         plt.grid(True)
205         plt.subplot(1, 2, 2)
206         plt.plot([log.num_trees for log in logs], [log.evaluation.loss
207             for log in logs])
208         plt.xlabel("Number of trees")
209         plt.ylabel("Logloss (out-of-bag)")
```



```
200     plt.grid(True)
201     plt.show()

202     """"---

203     # Carga y Preprocesado del Dataset

204     Leer el DataSet desde un archivo .csv

205     IMPORTANTE: Cambios en los valores del DataSet en las variables
formato string

Autism: Yes - 1 // No - 0
Gender: f - 1 // m - 0
Ethnicity:

206     'Middle eastern' - 1 'White European' - 2 'Hispanic' - 3 'Black' - 4
'Asian'
- 5 'South asian' - 6 'Native Indian' - 7 'Others' - 8 'Latino' - 9 'Mixed'
-
    10 'Pacifica' - 11

Family Member with ASD:

207     'Family Member' - 0 'Health Care Professional' - 1 'Self' - 2

208     Se elimina la columna de QChat-10-Score pq causa over-feating en
Random Forest

209     ---

210     **Pre Procesado de DataSet**
211     """"

212     #Montar en Google Drive (SOLO para google colab)

213     #from google.colab import drive
214     #drive.mount('/content/drive')
215     #dataset = pd.read_csv("drive/MyDrive/TFG/Toddler Autism dataset
July 2018_Text.csv").drop(columns=["Case_No", "Qchat-10-
Score", "Who completed the test"])

216     dataset = pd.read_csv("Toddler Autism dataset July
2018_Text.csv").drop(columns=["Case_No", "Qchat-10-
Score", "Who completed the test"]) #DS Original
217     dataset.head()
```



```
218     for counter, element in dataset.iterrows():
219         if element["Class/ASD_Traits"] == "Yes":
220             dataset.iloc[counter,15] = 1
221         elif element["Class/ASD_Traits"] == "No":
222             dataset.iloc[counter,15] = 0
223
224         if element["Sex"] == "f":
225             dataset.iloc[counter,11] = 1
226         elif element["Sex"] == "m":
227             dataset.iloc[counter,11] = 0
228
229         if element["Family_mem_with_ASD"] == "yes":
230             dataset.iloc[counter,14] = 1
231         elif element["Family_mem_with_ASD"] == "no":
232             dataset.iloc[counter,14] = 0
233
234         if element["Jaundice"] == "yes":
235             dataset.iloc[counter,13] = 1
236         elif element["Jaundice"] == "no":
237             dataset.iloc[counter,13] = 0
238
239         if element["Ethnicity"] == "White European":
240             dataset.iloc[counter,12] = 1
241         elif element["Ethnicity"] == "Hispanic":
242             dataset.iloc[counter,12] = 2
243         elif element["Ethnicity"] == "black":
244             dataset.iloc[counter,12] = 3
245         elif element["Ethnicity"] == "asian":
246             dataset.iloc[counter,12] = 4
247         elif element["Ethnicity"] == "south asian":
248             dataset.iloc[counter,12] = 5
249         elif element["Ethnicity"] == "Native Indian":
250             dataset.iloc[counter,12] = 6
251         elif element["Ethnicity"] == "Others":
252             dataset.iloc[counter,12] = 7
253         elif element["Ethnicity"] == "Latino":
254             dataset.iloc[counter,12] = 8
255         elif element["Ethnicity"] == "mixed":
256             dataset.iloc[counter,12] = 9
257         elif element["Ethnicity"] == "Pacifica":
258             dataset.iloc[counter,12] = 10
259         elif element["Ethnicity"] == "middle eastern":
260             dataset.iloc[counter,12] = 11
```



```
257     print(dataset["Class/ASD_Traits"].unique())
258     print(dataset["Sex"].unique())
259     print(dataset["Family_mem_with_ASD"].unique())
260     print(dataset["Jaundice"].unique())
261     print(dataset["Ethnicity"].unique() )

262     """Partición de datos para Training (80%) Test & Validation (20%)
        [80% Test & 20% Validation]"""

263     dataset_train , dataset_test_validation =
        train_test_split(dataset.astype(np.float32) , test_size=0.2,
        random_state=42)
264     dataset_test , dataset_validation =
        train_test_split(dataset_test_validation , test_size=0.2,
        random_state=42)

265     """Separación de Datos y Labels"""

266     dataset_train_data , dataset_train_label = dataset_train.iloc[:, :-1] ,
        dataset_train.iloc[:, -1]
267     dataset_test_data , dataset_test_label = dataset_test.iloc[:, :-1] ,
        dataset_test.iloc[:, -1]
268     dataset_validation_data , dataset_validation_label =
        dataset_validation.iloc[:, :-1] , dataset_validation.iloc[:, -1]

269     """DS Original"""

270     dataset_text_train , dataset_text_test_validation =
        train_test_split(dataset , test_size=0.2, random_state=42)
271     dataset_text_test , dataset_text_validation =
        train_test_split(dataset_text_test_validation , test_size=0.2,
        random_state=42)

272     """Conversión DS en formato TF"""

273     TF_train =
        tfdf.keras.pd_dataframe_to_tf_dataset(dataset_text_train,
        label="Class/ASD_Traits")
274     TF_test =
        tfdf.keras.pd_dataframe_to_tf_dataset(dataset_text_test,
        label="Class/ASD_Traits")
275     TF_validation =
        tfdf.keras.pd_dataframe_to_tf_dataset(dataset_text_validation,
        label="Class/ASD_Traits")
```



```
276     """---
277     # Implementación de Modelos sin Validación Cruzada

278     Modelos basados en Decision Tree
279     """

280     model_fit_RF, dataset_test_predictions_RF , test_loss_RF ,
        test_accuracy_RF , training_time_RF
        = Decision_Tree_model_fit("RF") #Random Forest

281     model_fit_GB, dataset_test_predictions_GB , test_loss_GB ,
        test_accuracy_GB , training_time_GB
        = Decision_Tree_model_fit("GB") #Gradient Boosted

282     """Modelos Secuenciales"""

283     model_fit_Nadam , dataset_test_predictions_Nadam ,
        test_loss_Nadam ,test_accuracy_Nadam , training_time_Nadam
        = sequential_model_fit_by_optimizer('Nadam')

284     model_fit_Adam , dataset_test_predictions_Adam , test_loss_Adam ,
        test_accuracy_Adam , training_time_Adam
        = sequential_model_fit_by_optimizer('Adam')

285     model_fit_RMSprop , dataset_test_predictions_RMSprop ,
        test_loss_RMSprop,test_accuracy_RMSprop ,
        training_time_RMSprop
        = sequential_model_fit_by_optimizer('RMSprop')

286     model_fit_Adamax , dataset_test_predictions_Adamax ,
        test_loss_Adamax, test_accuracy_Adamax , training_time_Adamax
        = sequential_model_fit_by_optimizer('Adamax')

287     model_fit_Adagrad , dataset_test_predictions_Adagrad ,
        test_loss_Adagrad , test_accuracy_Adagrad , training_time_Adagrad
        = sequential_model_fit_by_optimizer('Adagrad')

288     model_fit_Adadelta , dataset_test_predictions_Adadelta ,
        test_loss_Adadelta ,test_accuracy_Adadelta ,training_time_Adadelta
        = sequential_model_fit_by_optimizer('Adadelta')
```



```
289      """ ---
290      # Comparación de Modelos sin Validación Cruzada
291      """
292      #Decision Tree models
293      print("Test accuracy with Random Forest: ",test_accuracy_RF)
294      print("Test accuracy with Grandient Boosted: ",test_accuracy_GB)
295      #Sequential Models
296      print("Test accuracy with Adam: ",test_accuracy_Adam)
297      print("Test accuracy with Nadam: ",test_accuracy_Nadam)
298      print("Test accuracy with RMSprop: ",test_accuracy_RMSprop)
299      print("Test accuracy with Adamax: ",test_accuracy_Adamax)
300      print("Test accuracy with Adagrad: ",test_accuracy_Adagrad)
301      print("Test accuracy with Adadelata: ",test_accuracy_Adadelata)
302      #Sequential Models
303      print("Loss with Adam: ",test_loss_Adam)
304      print("Loss with Nadam: ",test_loss_Nadam)
305      print("Loss with RMSprop: ",test_loss_RMSprop)
306      print("Loss with Adamax: ",test_loss_Adamax)
307      print("Loss with Adagrad: ",test_loss_Adagrad)
308      print("Loss with Adadelata: ",test_loss_Adadelata)
309      #Decision Tree models
310      print("Test accuracy with Random Forest: ",training_time_RF)
311      print("Test accuracy with Grandient Boosted: ",training_time_GB)
312      #Sequential Models
313      print("Test accuracy with Adam: ",training_time_Adam)
314      print("Test accuracy with Nadam: ",training_time_Nadam)
315      print("Test accuracy with RMSprop: ",training_time_RMSprop)
316      print("Test accuracy with Adamax: ",training_time_Adamax)
317      print("Test accuracy with Adagrad: ",training_time_Adagrad)
318      print("Test accuracy with Adadelata: ",training_time_Adadelata)
```



```
319     """# Implementación de Modelos con Validación Cruzada
320     Decision Trees con Validación Cruzada (K=10)
321     """
322     print('\n-----')
323     print('RANDOM FOREST')
324     print('-----')
325     acc_RF , loss_RF = Decision_Tree_model_fit_kfold("RF")
326     print('\n-----')
327     print('GRADIENT BOOSTED')
328     print('-----')
329     acc_GB , loss_GB = Decision_Tree_model_fit_kfold("GB")
330     """Modelos Secuenciales con Validación Cruzada (k=10)"""
331     print('\n-----')
332     print('ADAM')
333     print('-----')
334     acc_Adam , loss_Adam =
335     sequential_model_fit_by_optimizer_kfold('Adam')
336     print('\n-----')
337     print('NADAM')
338     print('-----')
339     acc_Nadam , loss_Nadam =
340     sequential_model_fit_by_optimizer_kfold('Nadam')
341     print('\n-----')
342     print('RMSPROP')
343     print('-----')
344     acc_RMSProp , loss_RMSProp =
345     sequential_model_fit_by_optimizer_kfold('RMSProp')
346     print('\n-----')
347     print('ADAMAX')
348     print('-----')
349     acc_Adamax , loss_Adamax =
350     sequential_model_fit_by_optimizer_kfold('Adamax')
351     print('\n-----')
352     print('ADAGRAD')
353     print('-----')
354     acc_Adagrad , loss_Adagrad =
355     sequential_model_fit_by_optimizer_kfold('Adagrad')
356     print('\n-----')
357     print('ADADELTA')
358     print('-----')
359     acc_Adadelta , loss_Adadelta =
360     sequential_model_fit_by_optimizer_kfold('Adadelta')
```



```
355     """# Comparación de Modelos con Validación Cruzada"""
356     #Decision Tree
357     print("Test accuracy with Random Forest: ",acc_RF)
358     print("Test accuracy with Gradient Bosted: ",acc_GB)
359     #Sequential Models
360     print("Test accuracy with Adam: ",acc_Adam)
361     print("Test accuracy with Nadam: ",acc_Nadam)
362     print("Test accuracy with RMSProp: ",acc_RMSProp)
363     print("Test accuracy with Adamax: ",acc_Adamax)
364     print("Test accuracy with Adagrad: ",acc_Adagrad)
365     print("Test accuracy with Adadelata: ",acc_Adadelata)
366     #Decision Tree
367     print("Test loss with Random Forest: ",loss_RF)
368     print("Test loss with Gradient Bosted: ",loss_GB)
369     #Sequential Models
370     print("Test loss with Adam: ",loss_Adam)
371     print("Test loss with Nadam: ",loss_Nadam)
372     print("Test loss with RMSProp: ",loss_RMSProp)
373     print("Test loss with Adamax: ",loss_Adamax)
374     print("Test loss with Adagrad: ",loss_Adagrad)
375     print("Test loss with Adadelata: ",loss_Adadelata)
376     """---
377     **Matrices de Confusión**
378     # Matrices de Confusión
379     Matrices de confusión de modelos de Decision Tree
380     """
381     generate_conf_matrix(dataset_test_label,dataset_test_predictions_
RF,'Random Forest')
382     generate_conf_matrix(dataset_test_label,dataset_test_predictions_
GB,'Gradient Boosted')
```



```
383     """
384     ---
385     Matrices de confusión de modelos secuenciales"""

386     generate_conf_matrix(dataset_test_label,dataset_test_predictions_A
387                          dam,'Adam')
387     generate_conf_matrix(dataset_test_label,dataset_test_predictions_N
388                          adam,'Nadam')
388     generate_conf_matrix(dataset_test_label,dataset_test_predictions_R
389                          MSprop,'RMSprop')
389     generate_conf_matrix(dataset_test_label,dataset_test_predictions_A
390                          damax,'Adamax')
390     generate_conf_matrix(dataset_test_label,dataset_test_predictions_A
391                          dagrad,'Adagrad')
391     generate_conf_matrix(dataset_test_label,dataset_test_predictions_A
392                          dadelta,'Adadelta')

392     """---

393     # Gráficas de Precisión y Pérdidas

394     Modelos Decision Tree
395     """

396     Decision_Tree_generate_accuracy_loss_graphs(model_fit_RF)
397     Decision_Tree_generate_accuracy_loss_graphs(model_fit_GB)

398     """
399     ---
400     Modelos Secuenciales"""

401     generate_accuracy_loss_graphs(model_fit_RMSprop,'RMSprop')
402     generate_accuracy_loss_graphs(model_fit_Nadam,'Nadam')
403     generate_accuracy_loss_graphs(model_fit_Adam,'Adam')
404     generate_accuracy_loss_graphs(model_fit_Adamax,'Adamax')
405     generate_accuracy_loss_graphs(model_fit_Adagrad,'Adagrad')
406     generate_accuracy_loss_graphs(model_fit_Adadelta,'Adadelta')
```



B) Código del segundo estudio

Código correspondiente al estudio “Técnicas de 'Deep Learning' sobre datos 'Eye Tracking”

```
1      # -*- coding: utf-8 -*-
2      """TFG_Et3_v05_29_2.ipynb

3      # Importación de Librerías
4      """

5      import pandas as pd
6      import numpy as np
7      import os
8      import glob
9      import csv
10     import math
11     import time
12     import psutil

13     import tensorflow as tf
14     from tensorflow.keras.layers import LSTM, Dense
15     import matplotlib.pyplot as plt

16     #conda env config vars set
17     SKLEARN_ALLOW_DEPRECATED_SKLEARN_PACKAGE_INST
18     ALL=True
19     from sklearn.model_selection import train_test_split
20     from sklearn.model_selection import KFold
21     from sklearn.model_selection import train_test_split

22     """# Asignación de Datos y Metadatos
23     Este conjunto de código se ejecuta una sola vez para generar un
24     único dataset con los valores de entrenamiento y test reales

25     ---
26     Unión del conjunto de datasets en uno único
27     """

28     dataset_union = [i for i in glob.glob('DS/*.{}'.format('csv'))]
29     dataset = pd.concat([pd.read_csv(f) for f in dataset_union ])
30     dataset.head()
```



```
28     """Lectura de los archivos correspondientes al dataset y al
29     metadata"""
30     dataset = pd.read_csv("DataSet.csv")
31     metadata = pd.read_csv("Metadata_Participants.csv")
32
33     """Preprocesado Inicial
34     ---
35     Detección de valores del DataSet sin identificar (Valores
36     pertenecientes a sujetos no identificados)
37
38     Detecta el conjunto de datos perteneciente a un mismo participante
39     (Utilizando la variable ID) y asigna el valor de etiqueta
40     correspondiente
41
42     Si no detecta el ID del participante asigna una etiqueta "NODATA"
43     para su posterior eliminación
44     """
45
46     dataset_with_labels = dataset
47     counter=1
48     previous_row=0
49     number_values = len(dataset_with_labels)
50
51     for row in dataset_with_labels["Participant"]:
52         try:
53             if row != previous_row:
54                 dataset_with_labels.loc[dataset_with_labels["Participant"] ==
55                 row , "LABEL"] =
56                 metadata["Class"][metadata["ParticipantID"] ==
57                 int(row)].values[0]
58
59                 previous_row = row
60
61                 print("--> Row ", counter, " of ", number_values, " --- ", row,
62                       ", ",
63
64                 metadata["ParticipantID"][metadata["ParticipantID"] ==
65                 int(row)].values[0], " , ",
66
67                 metadata["Class"][metadata["ParticipantID"] ==
68                 int(row)].values[0])
69
70         except ValueError as error:
```



```
50         if row != previous_row:
51             dataset_with_labels.loc[dataset_with_labels["Participant"]
== row , "LABEL"] = "NODATA"
52             previous_row = row
53             print("--> Row ", counter, " of ", number_values, " --- ", row,
, " NOT VALID")
54         finally:

55             counter+=1

56         #Guardar el Archivo .csv
57
dataset_with_labels.to_csv("DataSet_With_Labels_NoData.csv"
, index=False)

58         """Elimina las variables etiquetadas como NODATA"""

59         dataset_with_labels_and_nodata =
pd.read_csv("DataSet_With_Labels_NoData.csv")

60         with open("DataSet_With_Labels.csv", 'w', newline=") as
DataSet_With_Labels_csv:
61             dataset_with_labels = csv.writer(DataSet_With_Labels_csv)
62             dataset_with_labels.writerow
(dataset_with_labels_and_nodata[1:2])

63         counter = 0

64         for row in dataset_with_labels_and_nodata["LABEL"]:
65             if not row == "NODATA":
66                 dataset_with_labels.writerow
(dataset_with_labels_and_nodata.iloc[counter].values)
67                 print("--> Adding Row number ", counter, " --- ", row)

68             else:
69                 print("--> NODATA Detected")
70                 counter+=1

71         """# Pre Procesado de Datos

72         Cargar el DataSet con los metadatos (Etiquetas)
73         """

74         dataset_with_labels =
pd.read_csv("DataSet_With_Labels.csv").replace('-',pd.NA)

75         """Borrado de las columnas despreciables para el estudio"""
```



```
76     columns_to_drop = [  
77         "Unnamed: 0", "RecordingTime [ms]",  
            "Time of Day [h:m:s:ms]", "Trial",  
            "Stimulus", "Export Start Trial Time [ms]",  
            "Export End Trial Time [ms]",  
            "Color", "Tracking Ratio [%]", "Category Group",  
            "Category Right", "Category Left", "Index Right", "Index Left",  
            "AOI Name Right", "AOI Name Left", "Annotation Name",  
            "Annotation Description", "Annotation Tags",  
            "Mouse Position X [px]", "Mouse Position Y [px]",  
            "Scroll Direction X", "Scroll Direction Y", "Content", "Port Status",  
            "AOI Group Right", "AOI Scope Right", "AOI Order Right",  
            "AOI Group Left", "AOI Scope Left", "AOI Order Binocular",  
            "groupe d'enfants"  
78     ]  
  
79     dataset_with_labels =  
        dataset_with_labels.drop(columns=columns_to_drop)  
  
80     """Detección y aproximación de los datos nulos"""  
  
81     def value_is_valid(value_to_inspect):  
82         try:  
83             return float(value_to_inspect)  
84         except ValueError:  
85             return np.nan  
  
86     index_data_pending = []  
87     index_temp_data = []  
88     fixed_columns = 2  
89     length_dataset = len(dataset_with_labels)  
  
90     for index in range(length_dataset):  
91         print("-> Índice de Fila ", index,  
            " Valores Temporales Almacenados:  
            ", len(index_temp_data), " Valores Pendientes:  
            ", len(index_data_pending))  
92         nan_counter = 0  
93         nan_column = []  
94         last_data_of_participant = False  
  
95         if dataset_with_labels["Participant"].iloc[index] !=  
            dataset_with_labels["Participant"].iloc[index+1]:  
96             last_data_of_participant = True  
97         if dataset_with_labels["Participant"].iloc[index] !=  
            dataset_with_labels["Participant"].iloc[index-1]:
```



```
98         if (len(index_data_pending) > 0) or  
           (len(index_temp_data) > 0):  
  
99             print("--> Queda información sin recuperar en  
             subconjunto de datos")  
100            last_data_of_participant = False  
101            index_data_pending = []  
102            index_temp_data = []  
  
103            for column in range(len(dataset_with_labels.columns)):  
104                if (pd.isnull(dataset_with_labels.iloc[index,column])):  
105                    nan_counter += 1  
106                    nan_column.append(column)  
  
107            if (nan_counter == len(dataset_with_labels.columns) -  
            fixed_columns):  
108                dataset_with_labels =  
                dataset_with_labels.drop(index)  
109                length_dataset -= 1  
110                print("--> Fila Eliminada")  
111            else:  
112                for column in nan_column:  
  
113                    value = dataset_with_labels.iloc[index,column]  
114                    last_value = dataset_with_labels.iloc[index-1,column]  
115                    next_value =  
                    dataset_with_labels.iloc[index+1,column]  
  
116                    if (pd.isnull(value)):  
117                        value = np.nan  
  
118                    if (index != 0 and not pd.isna(last_value)):  
119                        last_value = value_is_valid(last_value)  
120                    else:  
121                        last_value = np.nan  
  
122                    if (index != length_dataset and  
                    not pd.isna(next_value)):  
123                        next_value =  
                        value_is_valid(next_value)  
124                    else:  
125                        next_value = np.nan
```



```
126         if (not math.isnan(last_value) and
127             not math.isnan(next_value)):
128             if (last_value == next_value):
129                 dataset_with_labels.iloc[index,column] =
130                     last_value
131                 # Valor NaN entre dos valores actualizado
132             else:
133                 dataset_with_labels.iloc[index,column] =
134                     (next_value + last_value)/2
135                 # Valor NaN entre dos valores aproximado
136
137         elif (not math.isnan(last_value) and
138             math.isnan(next_value)):
139             index_temp_data.append([index-1,column])
140             # Valor Temporal Almacenado
141
142         elif (last_data_of_participant == False) and
143             (math.isnan(last_value) and
144             math.isnan(next_value)):
145             index_data_pending.append([index,column])
146             # Valor de Columna Pendiente de actualización
147
148         elif (last_data_of_participant == True) or
149             (math.isnan(last_value) and not
150             math.isnan(next_value)):
151             index_row_last_data = -1
152
153         for element_temp_data in
154             reversed(index_temp_data):
155             if element_temp_data[1] == column:
156                 index_row_last_data =
157                     element_temp_data[0]
158                 break
159
160         if index_row_last_data == -1:
161             temp_data = np.nan
162         else:
163             temp_data = dataset_with_labels.iloc
164             [index_row_last_data,column]
```



```
151         if (math.isnan(temp_data) or temp_data ==
152             next_value):
153             dataset_with_labels.iloc[index,column] =
154             next_value
155
156             index_row_data_pending = []
157             index_data_pending_updated = []
158
159             for index, data_pending_column in
160                 index_data_pending:
161
162                 if data_pending_column == column:
163                     index_row_data_pending.append(index)
164                 else:
165                     index_data_pending_updated.append(
166                         (index, data_pending_column))
167
168             index_data_pending =
169             index_data_pending_updated
170
171             dataset_with_labels.iloc
172             [index_row_data_pending,column] =
173             next_value
174
175             #Valores Anteriores actualizados a valor existente
176
177             else:
178                 dataset_with_labels.iloc[index,column] =
179                 next_value
180                 index_row_data_pending = []
181
182                 for index, data_pending_column in
183                     index_data_pending:
184
185                     if data_pending_column[1] == column:
186                         index_row_data_pending.append(index)
187
188             len_index_row_data_pending =
189             len(index_row_data_pending)
190
191             delta = (next_value -
192                 temp_data)/len_index_row_data_pending
193             increment_counter = 1
```



```
172         for row in index_row_data_pending:
173             if (temp_data < next_value):
174                 dataset_with_labels.iloc[row,column] =
temp_data + increment_counter*delta
175             else:
176                 dataset_with_labels.iloc[row,column] =
temp_data - increment_counter*delta

177                 increment_counter += 1

178                 index_data_pending = [element_index_data
for element_index_data in
index_data_pending
if element_index_data !=
[index_row_data_pending, column]]

179     # Valores Anteriores aproximados a valor existente
180     index += 1 #Re-actualización de index

181     dataset_with_labels.to_csv("DataSet_With_Labels_Less_NaN.csv",
index=False)

182     """Detección de los valores atípicos por medio del rango
intercuartil"""

183     dataset_with_labels =
pd.read_csv("DataSet_With_Labels_Less_NaN.csv")

184     columns =
dataset_with_labels.drop(columns=["Participant","LABEL"]).columns

185     for column in columns:
186         mean_column = dataset_with_labels[column].mean()
187         dataset_with_labels[column] =
dataset_with_labels[column].fillna(mean_column)

188         Q1 = dataset_with_labels[column].quantile(0.25)
189         Q3 = dataset_with_labels[column].quantile(0.75)
190         IQR = Q3 - Q1

191         upper_limit = Q3 + 1.5 * IQR
192         lower_limit = Q1 - 1.5 * IQR

193     dataset_with_labels.loc[dataset_with_labels[column] > upper_limit,
column] = upper_limit
194     dataset_with_labels.loc[dataset_with_labels[column] < lower_limit,
column] = lower_limit
```



```
195 dataset_with_labels.to_csv("DataSet_With_Labels_Less_NaN
   _without_outliers.csv", index=False)

196 """# Selección de Características a Estudiar
197 Correlación de las características respecto al valor label
198 """

199 #from google.colab import drive
200 #drive.mount('/content/drive')

201 from scipy.stats import pearsonr

202 dataset_with_labels = pd.read_csv
   ("DataSet_With_Labels_Less_NaN_without_outliers.csv")

203 column_corr_pvalue = []
204 columns =
   dataset_with_labels.drop(columns=
   ["Participant","LABEL"]).columns

205 for column in columns:
206     corr, p_value = pearsonr(dataset_with_labels[column],
   dataset_with_labels["LABEL"].replace({'TD':0,'ASD':1}))

207     if not math.isnan(corr):
208         column_corr_pvalue.append([column,corr,p_value])
209     else:
210         column_corr_pvalue.append([column,0,0])

211 column_corr_pvalue.sort(key=lambda x:
   abs(x[1]),reverse=True)

212 for element in column_corr_pvalue:
213     print(f"-> {element[0]}: Correlation = {element[1]:.3f} ,
214           P-Value: {element[2]:.3f}")
```



```
215     """Selección de las características a estudiar
216     (según correlación)"""
217     features = [
218         "Gaze Vector Right Y",
219         "Gaze Vector Right Z",
220         "Gaze Vector Left Y",
221         "Gaze Vector Left Z",
222         "Point of Regard Right X [px]",
223         "Point of Regard Left X [px]",
224         "Point of Regard Left Y [px]",
225         "Point of Regard Right Y [px]",
226
227         "Pupil Diameter Right [mm]",
228         "Pupil Diameter Left [mm]",
229         "Eye Position Right X [mm]",
230         "Pupil Position Left X [px]",
231         "Eye Position Left Y [mm]",
232         "Eye Position Right Y [mm]",
233         "Eye Position Right Z [mm]",
234         "Eye Position Left Z [mm]",
235
236         "Pupil Position Right Y [px]",
237         "Pupil Position Left Y [px]",
238         "Pupil Position Right X [px]",
239         "Pupil Size Left Y [px]",
240         "Pupil Size Left X [px]",
241         "Pupil Size Right Y [px]",
242         "Pupil Size Right X [px]",
243         "Eye Position Left X [mm]"
244
245         # "Gaze Vector Right X", #Correlación 0
246         # "Gaze Vector Left X" #Correlación 0
247
248     ]
249
250     """**Pre Procesamiento de datos antes del entrenamiento**
251     ---
252     Separación e intercalación de ambos casos de estudio
253     para crear los dos subconjuntos de datos (train y test)
254     """
255
256     dataset_with_labels_ASD =
257     dataset_with_labels[dataset_with_labels["LABEL"]=="ASD"]
258     dataset_with_labels_TD =
259     dataset_with_labels[dataset_with_labels["LABEL"]=="TD"]
```



```
250 id_ASD = dataset_with_labels_ASD["Participant"].unique()
251 id_TD = dataset_with_labels_TD["Participant"].unique()
      [:len(id_ASD)]
252 id_dataset = [value for par in zip(id_ASD, id_TD)
      for value in par]

253 dataset_with_labels_train_mix =
      pd.DataFrame(columns=dataset_with_labels.columns)
254 dataset_with_labels_test_mix = []

255 len_train=int(len(id_dataset)//1.25)

256 for value in id_dataset[:len_train]:
257     dataset_with_labels_train_mix =
            pd.concat([dataset_with_labels_train_mix,
            dataset_with_labels[dataset_with_labels
            ["Participant"] ==int(value)])

258 for value in id_dataset[len_train:]:
259     dataset_with_labels_test_mix.append
            ([value,dataset_with_labels[dataset_with_labels
            ["Participant"]==int(value)])

260 """Separación entre datos y etiquetas para cada uno de los
      subconjuntos"""

261 dataset_train_label =
      dataset_with_labels_train_mix["LABEL"]
      .replace({'TD':0,'ASD':1})

262 dataset_train_data =
      dataset_with_labels_train_mix[features].astype(float)

263 dataset_test_data = []
264 dataset_test_label = []

265 for element in dataset_with_labels_test_mix:
266     dataset_test_label.append(element[1]["LABEL"]
      .replace({'TD':0,'ASD':1}))

267 dataset_test_data.append(element[1][features].astype(float))
```



```
268         """# Funciones para entrenamiento y test de los Modelos"""
269     def model_fit_kfold(model,epochs = 50,folder_init = 1):
270
271         kfolder = KFold(n_splits=10, shuffle=False)
272         folder_number = 1
273         acc_per_fold = []
274         loss_per_fold = []
275
276         #Iteración por folder
277
278         data = dataset_train_data
279         label = dataset_train_label
280
281         for test_data , test_label in
282             zip(dataset_test_data,dataset_test_label):
283
284             data = np.concatenate((data, test_data), axis=0)
285             label = np.concatenate((label, test_label), axis=0)
286
287             print("Use of ",len(features)," Features\n")
288             total_start_time = time.time()
289
290             for train , test in kfolder.split(data,label):
291                 participant_iter = participant_start
292                 if folder_number == folder_init:
293                     model.compile(loss=
294                         "sparse_categorical_crossentropy",
295                         optimizer="Adamax",
296                         metrics=["sparse_categorical_accuracy"] )
297
298                 print('-----')
299                 print(f'-> Training for folder number
300                     {folder_number}')
301
302                 folder_start_time = time.time()
303
304                 model_fit = model.fit(data[train],label[train],epochs =
305                     epochs,verbose=1)
306
307                 folder_training_time = time.time() - folder_start_time
```



```
290     print('--> Training Time in folder: ',folder_training_time)
291     print('-----')
292     print(' R E S U L T S Folder
        ',folder_number)
293     print('-----')

294     values = model.evaluate(data[test],label[test])
295     print(f'--> Score for folder {folder_number}:
        {model.metrics_names[0]} of {values[0]};
        {model.metrics_names[1]} of {values[1]*100}%\n')

296     acc_per_fold.append(values[1] * 100)
297     loss_per_fold.append(values[0])

298     folder_init += 1

299     folder_number = folder_number + 1

300     total_training_time = time.time() - total_start_time

301     #Resumen por folder

302     print('-----')
303     print('-> Score per fold')

304     for i in range(0, len(acc_per_fold)):
305         print('-----')
306         print(f'--> Fold {i+1} - Loss: {loss_per_fold[i]} - Accuracy:
        {acc_per_fold[i]}%')

307         print('-----')
308         print('-> Average scores for all folds:')
309         print(f'--> Accuracy: {np.mean(acc_per_fold)} (+-
        {np.std(acc_per_fold)})')
310         print(f'--> Loss: {np.mean(loss_per_fold)}')
311         print(f'--> Training Time: {total_training_time}')
312         print('-----\n')

313     def sequential_model():

314         model = tf.keras.models.Sequential()
315         model.add(tf.keras.layers.Dense(units=64, activation='relu'))
316         model.add(tf.keras.layers.Dropout(0.2))
317         model.add(tf.keras.layers.Dense(units=len(features),
        activation='relu'))
318         model.add(tf.keras.layers.Dense(units=2,
        activation='softmax'))
```



```
319         return model

320     def sequential_lstm_model():

321         model = tf.keras.models.Sequential()
322         model.add(LSTM(64, input_shape=(len(features),1)))
323         model.add(tf.keras.layers.Dense(units=64, activation='relu'))
324         model.add(tf.keras.layers.Dropout(0.2))
325         model.add(tf.keras.layers.Dense(units=len(features),
326         activation='relu'))
327         model.add(tf.keras.layers.Dense(units=2,
328         activation='softmax'))
329         return model

330     def sequential_convolutional_model():

331         model = tf.keras.models.Sequential()
332         model.add(tf.keras.Input(shape=(len(features),1)))

333         model.add(tf.keras.layers.Conv1D(filters=16, kernel_size=3,
334         activation='relu', padding='same'))

335         model.add(tf.keras.layers.MaxPooling1D())
336         model.add(tf.keras.layers.Conv1D(filters=32, kernel_size=3,
337         activation='relu', padding='same'))

338         model.add(tf.keras.layers.MaxPooling1D())
339         model.add(tf.keras.layers.Flatten())

340         model.add(tf.keras.layers.Dropout(0.2))
341         model.add(tf.keras.layers.Dense(units=32, activation='relu'))
342         model.add(tf.keras.layers.Dense(units=2,
343         activation='softmax'))
344         return model

345     def sequential_convolutional_lstm_model():

346         model = tf.keras.models.Sequential()
347         model.add(tf.keras.Input(shape=(len(features),1)))

348         model.add(tf.keras.layers.Conv1D(filters=16, kernel_size=3,
349         activation='relu', padding='same'))
350         model.add(tf.keras.layers.MaxPooling1D())

351         model.add(tf.keras.layers.Conv1D(filters=32, kernel_size=3,
352         activation='relu', padding='same'))
```



```
346     model.add(tf.keras.layers.MaxPooling1D())
347     model.add(LSTM(64, input_shape=(128,1)))
348     model.add(tf.keras.layers.Flatten())

349     model.add(tf.keras.layers.Dropout(0.2))

350     model.add(tf.keras.layers.Dense(units=32, activation='relu'))
351     model.add(tf.keras.layers.Dense(units=2,
352         activation='softmax'))
352     return model

353     def gru_model():

354         model = tf.keras.models.Sequential()
355         model.add(tf.keras.Input(shape=(len(features),1)))

356         model.add(tf.keras.layers.GRU(32))
357         model.add(tf.keras.layers.Dense(len(features)))
358         model.add(tf.keras.layers.Dropout(0.2))

359         model.add(tf.keras.layers.Dense(units=len(features),
360             activation='relu'))
360         model.add(tf.keras.layers.Dense(units=2,
361             activation='softmax'))
361         return model

362     #Función para parar el entorno de Google Colab

363     #def parar_entorno():
364     # from google.colab import runtime
365     # runtime.unassign()

366     """# Implementación y entrenamiento del Modelo Secuencial

367     Modelo Secuencial
368     """

369     sequential = sequential_model()
370     model_fit_kfold(sequential,50,1)

371     """Modelo LSTM"""

372     lstm = sequential_lstm_model()
373     model_fit_kfold(lstm,50,1)
```



```
374      """Modelo Convolutacional"""  
  
375      convolucional = sequential_convolutacional_model()  
376      model_fit_kfold(convolucional,50,1)  
  
377      """Modelo Convolutacional & LSTM"""  
  
378      convolucional_lstm = sequential_convolutacional_lstm_model()  
379      model_fit_kfold(convolucional_lstm,50,1)  
  
380      """Modelo GRU"""  
  
381      gru = gru_model()  
382      model_fit_kfold(gru,50,1)
```
