

## ▼ NLP - Ejercicio Final Integrador

Entrenar un clasificador binario para análisis de sentimiento de sentimiento utilizando diferentes técnicas de NLP.

1. Descargar el dataset del siguiente [link](#).
2. Pre-procesar el texto de manera básica:
  - a. Eliminar tags html: por ejemplo <br/>
  - b. Eliminar puntuaciones
  - c. Eliminar stopwords
3. Entrenar un clasificador utilizando BOW.
  - a. Definir el tamaño del vocabulario.
  - b. (opcional) Agregar pesos utilizando TFIDF.
  - c. Transformar los textos en vectores.
  - d. Crear los datasets de train, validation y test.
  - e. (opcional) Aplicar PCA para reducir la dimensión de los vectores.
  - f. Entrenar una red neuronal (seleccionar arquitectura, loss y optimizador).
  - g. Medir AUC.
4. Entrenar un clasificador utilizando words embeddings (probar con GloVe y con FastText).
  - a. Calcular el embedding de cada texto como el promedio de los embeddings de las palabras.
  - b. Crear los datasets de train, validation y test.
  - c. Entrenar una red neuronal (seleccionar arquitectura, loss y optimizador).
  - d. Medir AUC.
  - e. Tratar de mejorar los embeddings de los textos con:
    - Agregar TFIDF como pesos a las palabras.
    - En lugar de tomar el promedio probar tomando el max o el min de cada componente de los embeddings de las palabras.
    - Probar diferentes tamaños de embeddings.
5. Entrenar un clasificador utilizando celdas LSTM.
  - a. Convertir secuencias de palabras a secuencias de números (indexer).
  - b. Agregar padding para que cada elemento de entrenamiento tenga la misma longitud.
  - c. Armar una red LSTM con las siguientes capas:
    - Capa de embedding que transforma un número (index) en un embedding.
    - Agregar uno dos layers LSTM para obtener el embedding de la secuencias.
    - Agregar un layer denso para entrenar el clasificador.

## ▼ Desarrollo

```
1 # Para futura tabla comparativa (se agregó luego de entrenar y probar todos
2 # los modelos, por eso se incorporan manualmanete los resultados)
3 model_comparison_table = {}
```

### ▼ 1. Descarga y exploración de dataset.

1. Descargar el dataset del siguiente link.

```
1 !gdown --id 1zX2KM72Enhv8eWyJt1j4x6YZ75XwzHkm
2 !unzip "IMDB Dataset.csv.zip"
```

Downloading...

From: <https://drive.google.com/uc?id=1zX2KM72Enhv8eWyJt1j4x6YZ75XwzHkm>

To: /content/IMDB Dataset.csv.zip

27.0MB [00:00, 65.1MB/s]

Archive: IMDB Dataset.csv.zip

inflating: IMDB Dataset.csv

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
```

```
1 df = pd.read_csv("IMDB Dataset.csv")
2 df
```

```

1 df[df.isnull().any(axis=1)]

  review sentiment
2      I thought this was a wonderful way to spend ti... positive
1 df.sentiment.value_counts()

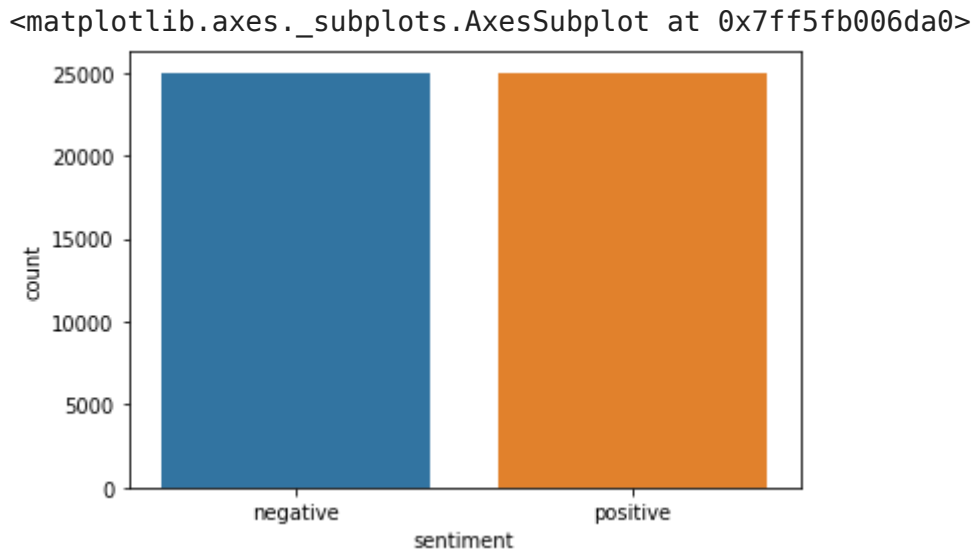
negative    25000
positive    25000
Name: sentiment, dtype: int64

100000 I thought this movie did a down right good job positive
1 df['sentiment'].unique()

['positive', 'negative']
Categories (2, object): ['positive', 'negative']

45550 I'm going to have to disagree with the previous... negative
1 sns.countplot(x='sentiment', data=df)

```



## ▼ 2. Preprocesamiento básico

Pre-procesar el texto de manera básica:

- a. Eliminar tags html: por ejemplo <br/>.
- b. Eliminar puntuaciones.
- c. Eliminar stopwords.

```

1 import nltk
2 import re

```

Tags HTML.

```

1 tag_re = re.compile(r'<[>]+>')
2 df['review'] = df['review'].apply(lambda x: tag_re.sub('',x) )

```

```
3 df.head()
```

	<b>review</b>	<b>sentiment</b>
0	One of the other reviewers has mentioned that ...	positive
1	A wonderful little production. The filming tec...	positive
2	I thought this was a wonderful way to spend ti...	positive
3	Basically there's a family where a little boy ...	negative
4	Petter Mattei's "Love in the Time of Money" is...	positive

## Puntuación

```
1 df['review'] = df['review'].apply(lambda x: re.sub('[^a-zA-Z]', ' ', x))
2 df.head()
```

	<b>review</b>	<b>sentiment</b>
0	One of the other reviewers has mentioned that ...	positive
1	A wonderful little production The filming tec...	positive
2	I thought this was a wonderful way to spend ti...	positive
3	Basically there s a family where a little boy ...	negative
4	Petter Mattei s Love in the Time of Money is...	positive

## Stop Words

```
1 nltk.download('stopwords')
2 stopwords_re = re.compile(r'\b(' + r'|'.join(nltk.corpus.stopwords.words('engli
3 df['review'] = df['review'].apply(lambda x: stopwords_re.sub('',x) )
4 df.head()
```

```
[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data] Package stopwords is already up-to-date!
```

	<b>review</b>	<b>sentiment</b>
0	One reviewers mentioned watching Oz episode ho...	positive
1	A wonderful little production The filming tec...	positive
2	I thought wonderful way spend time hot summer ...	positive
3	Basically family little boy Jake thinks zomb...	negative
4	Petter Mattei Love Time Money visually stunni...	positive

## Lowercase

```
1 df['review'] = df['review'].apply(lambda x: " ".join(x.lower() for x in x.split))
2 df.head()
```

	review	sentiment
0	one reviewers mentioned watching oz episode ho...	positive
1	a wonderful little production the filming tech...	positive
2	i thought wonderful way spend time hot summer ...	positive
3	basically family little boy jake thinks zombie...	negative
4	petter mattei love time money visually stunnin...	positive

Etiquetas.

```
1 df["sentiment"] = df["sentiment"].astype('category')
2 df["target"] = df["sentiment"].cat.codes
3 df.head()
```

	review	sentiment	target
0	one reviewers mentioned watching oz episode ho...	positive	1
1	a wonderful little production the filming tech...	positive	1
2	i thought wonderful way spend time hot summer ...	positive	1
3	basically family little boy jake thinks zombie...	negative	0
4	petter mattei love time money visually stunnin...	positive	1

### ▼ 3. Clasificador BOW con Red Neuronal

Entrenar un clasificador utilizando BOW.

- a. Definir el tamaño del vocabulario.
- b. (opcional) Agregar pesos utilizando TFIDF.
- c. Transformar los textos en vectores.
- d. Crear los datasets de train, validation y test.
- e. (opcional) Aplicar PCA para reducir la dimensión de los vectores.
- f. Entrenar una red neuronal (seleccionar arquitectura, loss y optimizador).
- g. Medir AUC.

```
1 X = df.review
2 y = df.target
```

```
1 from sklearn.model_selection import train_test_split
2
3 TRAIN_TEST_SPLIT = 0.2
4 TRAIN_VAL_SPLIT = 0.25
```

```

5
6 X_train_set, X_test, y_train_set, y_test = train_test_split(X, y, test_size=TRA
7 X_train, X_val, y_train, y_val = train_test_split(X_train_set,y_train_set,test_
8
9 print("Test:",X_test.shape)
10 print("Train:",X_train.shape)
11 print("Val:",X_val.shape)

```

```

    Test: (10000,)
    Train: (30000,)
    Val: (10000,)

```

```

1 from sklearn.feature_extraction.text import CountVectorizer
2 from sklearn.feature_extraction.text import TfidfTransformer
3
4 VOCAB_SIZE = 200000
5 count_vectorizer = CountVectorizer(stop_words="english",max_features=VOCAB_SIZE

```

```

1 X_train = count_vectorizer.fit_transform(X_train)
2 X_val = count_vectorizer.transform(X_val)

```

```

1 tfidf_transformer = TfidfTransformer(smooth_idf=False)
2
3 X_train = tfidf_transformer.fit_transform(X_train)
4 X_val = tfidf_transformer.transform(X_val)

```

```

1 N_DIM_REDUCE = 512
2
3 from sklearn.decomposition import TruncatedSVD
4
5 svd = TruncatedSVD(N_DIM_REDUCE)
6
7 X_train = svd.fit_transform(X_train)
8 X_val = svd.transform(X_val)
9
10 X_train.shape,X_val.shape

```

```

    ((30000, 512), (10000, 512))

```

```

1 from sklearn.preprocessing import StandardScaler
2
3 scaler = StandardScaler()
4
5 X_train = scaler.fit_transform(X_train)
6 X_val = scaler.transform(X_val)

```

```

1 import keras
2 from keras.models import Sequential
3 from keras.layers import Dense,Dropout
4 from keras.utils.vis_utils import plot_model
5

```

```
6
7 model = Sequential()
8 model.add(Dense(50, activation = "relu", input_shape=(N_DIM_REduc, )))
9 model.add(Dropout(0.3, noise_shape=None, seed=None))
10 model.add(Dense(50, activation = "relu"))
11 model.add(Dropout(0.2, noise_shape=None, seed=None))
12 model.add(Dense(50, activation = "relu"))
13 model.add(Dense(1, activation = "sigmoid"))
14 model.summary()
15
16 plot_model(model, show_shapes=True, show_layer_names=True)
```

Model: "sequential\_4"

Layer (type)	Output Shape	Param #
dense_12 (Dense)	(None, 50)	25650
dropout_6 (Dropout)	(None, 50)	0
dense_13 (Dense)	(None, 50)	2550
dropout_7 (Dropout)	(None, 50)	0
dense_14 (Dense)	(None, 50)	2550
dense_15 (Dense)	(None, 1)	51

```

1 model.compile(
2     optimizer = "adam",
3     loss = "binary_crossentropy",
4     metrics = [keras.metrics.AUC(name="auc")]
5 )
6
7 NUM_EPOCHS = 5
8
9 history = model.fit( X_train, y_train,
10                     epochs= NUM_EPOCHS,
11                     batch_size = 32,
12                     validation_data = (X_val, y_val),
13                     verbose = True )

```

```

Epoch 1/5
938/938 [=====] - 7s 5ms/step - loss: 0.5964 - auc:
Epoch 2/5
938/938 [=====] - 4s 4ms/step - loss: 0.3231 - auc:
Epoch 3/5
938/938 [=====] - 4s 4ms/step - loss: 0.2839 - auc:
Epoch 4/5
938/938 [=====] - 4s 4ms/step - loss: 0.2537 - auc:
Epoch 5/5
938/938 [=====] - 4s 4ms/step - loss: 0.2393 - auc:

```

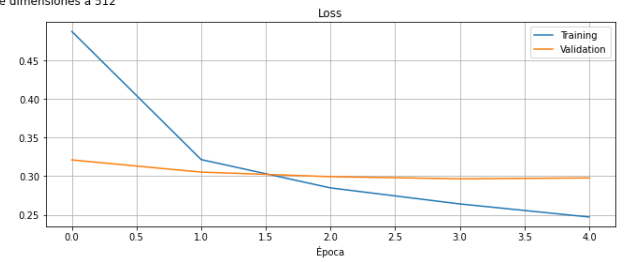
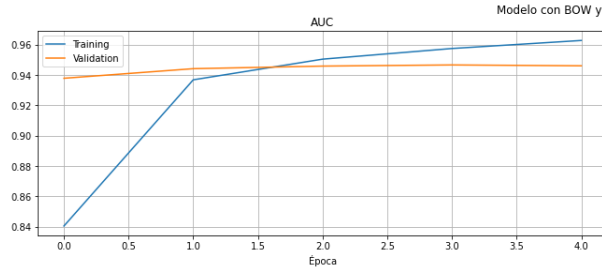
```

1 fig, axes = plt.subplots(1, 2, figsize=(24, 4))
2 plt.suptitle("Modelo con BOW y Reducción de dimensiones a %d" % N_DIM_REDUC)
3 axes[0].set_title("AUC")
4 axes[0].plot(np.arange(NUM_EPOCHS), history.history['auc'])
5 axes[0].plot(np.arange(NUM_EPOCHS), history.history['val_auc'])
6 axes[0].legend(["Training", "Validation"])
7 axes[0].grid(which="Both")
8 axes[0].set_xlabel("Época")
9 axes[1].set_title("Loss")
10 axes[1].plot(np.arange(NUM_EPOCHS), history.history['loss'])
11 axes[1].plot(np.arange(NUM_EPOCHS), history.history['val_loss'])
12 axes[1].legend(["Training", "Validation"])
13 axes[1].grid(which="Both")
14 axes[1].set_xlabel("Época")
15 plt.show()

```



10 print(fpr, tpr, thresholds)

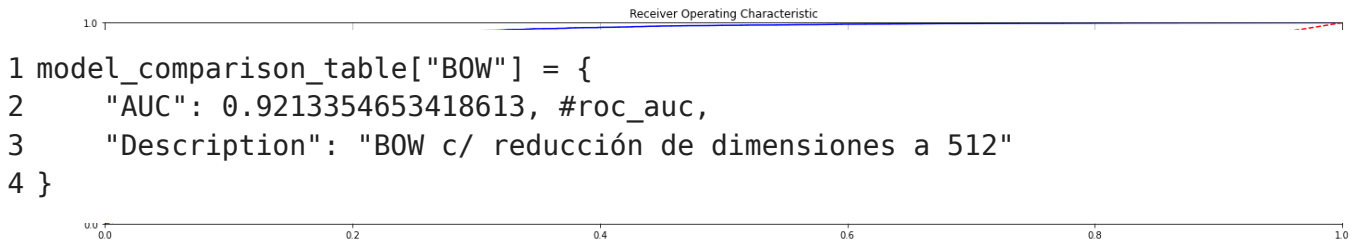


```

1 x_test_tx = count_vectorizer.transform(X_test)
2 x_test_tx = svd.transform(x_test_tx)
3 x_test_tx = scaler.transform(x_test_tx)
4 y_test_pred = model.predict(x_test_tx)

1 import sklearn
2 from sklearn.metrics import auc
3
4 fpr, tpr, thresholds = sklearn.metrics.roc_curve(y_test.values, y_test_pred)
5 roc_auc = auc(fpr, tpr)
6
7 plt.figure(figsize=(24,4))
8 plt.title('Receiver Operating Characteristic')
9 plt.plot(fpr, tpr, 'b', label = 'AUC = %0.2f' % roc_auc)
10 plt.legend(loc = 'lower right')
11 plt.plot([0, 1], [0, 1], 'r--')
12 plt.xlim([0, 1])
13 plt.ylim([0, 1])
14 plt.grid(which="Both")
15 plt.ylabel('True Positive Rate')
16 plt.xlabel('False Positive Rate')
17 plt.show()
18
19 print("AUC:", roc_auc)

```



## ▼ 4. Clasificador con Word Embeddings

- Entrenar un clasificador utilizando words embeddings (probar con GloVe y con FastText).
  - a. Calcular el embedding de cada texto como el promedio de los embeddings de las palabras.
  - b. Crear los datasets de train, validation y test.
  - c. Entrenar una red neuronal (seleccionar arquitectura, loss y optimizador).
  - d. Medir AUC.
  - e. Tratar de mejorar los embeddings de los textos con:
    - Agregar TFIDF como pesos a las palabras.
    - En lugar de tomar el promedio probar tomando el max o el min de cada componente de los embeddings de las palabras.
    - Probar diferentes tamaños de embeddings.

Descargar GloVe y FastText.

```

1 !wget http://nlp.stanford.edu/data/glove.twitter.27B.zip
2 !unzip glove.twitter.27B.zip

--2021-01-01 12:00:59-- http://nlp.stanford.edu/data/glove.twitter.27B.zip
Resolving nlp.stanford.edu (nlp.stanford.edu)... 171.64.67.140
Connecting to nlp.stanford.edu (nlp.stanford.edu)|171.64.67.140|:80... connect
HTTP request sent, awaiting response... 302 Found
Location: https://nlp.stanford.edu/data/glove.twitter.27B.zip [following]
--2021-01-01 12:00:59-- https://nlp.stanford.edu/data/glove.twitter.27B.zip
Connecting to nlp.stanford.edu (nlp.stanford.edu)|171.64.67.140|:443... connect
HTTP request sent, awaiting response... 301 Moved Permanently
Location: http://downloads.cs.stanford.edu/nlp/data/glove.twitter.27B.zip [following]
--2021-01-01 12:00:59-- http://downloads.cs.stanford.edu/nlp/data/glove.twitter.27B.zip
Resolving downloads.cs.stanford.edu (downloads.cs.stanford.edu)... 171.64.64.
Connecting to downloads.cs.stanford.edu (downloads.cs.stanford.edu)|171.64.64.
HTTP request sent, awaiting response... 200 OK
Length: 1520408563 (1.4G) [application/zip]
Saving to: 'glove.twitter.27B.zip'

glove.twitter.27B.z 100%[=====>] 1.42G 1.97MB/s in 11m 42s

2021-01-01 12:12:42 (2.06 MB/s) - 'glove.twitter.27B.zip' saved [1520408563/1520408563]

Archive: glove.twitter.27B.zip
  inflating: glove.twitter.27B.25d.txt
  inflating: glove.twitter.27B.50d.txt
  inflating: glove.twitter.27B.100d.txt
  inflating: glove.twitter.27B.200d.txt

```

```

1 !wget https://dl.fbaipublicfiles.com/fasttext/vectors-crawl/cc.en.300.vec.gz
2 !gunzip -k cc.en.300.vec.gz

--2021-01-01 22:21:09-- https://dl.fbaipublicfiles.com/fasttext/vectors-craw
Resolving dl.fbaipublicfiles.com (dl.fbaipublicfiles.com)... 104.22.74.142, 1
Connecting to dl.fbaipublicfiles.com (dl.fbaipublicfiles.com)|104.22.74.142|:
HTTP request sent, awaiting response... 200 OK
Length: 1325960915 (1.2G) [binary/octet-stream]
Saving to: 'cc.en.300.vec.gz'

cc.en.300.vec.gz      100%[=====>]    1.23G  50.6MB/s   in 26s

2021-01-01 22:21:35 (48.6 MB/s) - 'cc.en.300.vec.gz' saved [1325960915/132596

```

El siguiente código fue copiado del ejemplo de:

<https://github.com/ejesposito/ceai/blob/master/nlp/word2vec/pretrained.py>

```

1 import logging
2 import os
3 from pathlib import Path
4 from io import StringIO
5 import pickle
6
7 logging.basicConfig(level=logging.INFO)
8
9 class WordsEmbeddings(object):
10     logger = logging.getLogger(__name__)
11
12     def __init__(self):
13         # load the embeddings
14         words_embedding_pkl = Path(self.PKL_PATH)
15         if not words_embedding_pkl.is_file():
16             words_embedding_txt = Path(self.WORD_TO_VEC_MODEL_TXT_PATH)
17             assert words_embedding_txt.is_file(), 'Words embedding not available'
18             embeddings = self.convert_model_to_pickle()
19         else:
20             embeddings = self.load_model_from_pickle()
21         self.embeddings = embeddings
22         # build the vocabulary hashmap
23         index = np.arange(self.embeddings.shape[0])
24         self.word2idx = dict(zip(self.embeddings['word'], index))
25         self.idx2word = dict(zip(index, self.embeddings['word']))
26
27     def get_words_embeddings(self, words):
28         words_idx = self.words2idxs(words)
29         return self.embeddings[words_idx]['embedding']
30
31     def words2idxs(self, words):
32         return np.array([self.word2idx.get(word, -1) for word in words])
33
34     def idx2words(self, idxs):

```

```

34 def idx2word(self, idxs):
35     return np.array([self.idx2word.get(idx, '-1') for idx in idxs])
36
37 def load_model_from_pickle(self):
38     self.logger.debug(
39         'loading words embeddings from pickle {}'.format(
40             self.PKL_PATH
41         )
42     )
43     max_bytes = 2**28 - 1 # 256MB
44     bytes_in = bytearray(0)
45     input_size = os.path.getsize(self.PKL_PATH)
46     with open(self.PKL_PATH, 'rb') as f_in:
47         for _ in range(0, input_size, max_bytes):
48             bytes_in += f_in.read(max_bytes)
49     embeddings = pickle.loads(bytes_in)
50     self.logger.debug('words embeddings loaded')
51     return embeddings
52
53 def convert_model_to_pickle(self):
54     # create a numpy structured array:
55     # word      embedding
56     # U50       np.float32[]
57     # word_1    a, b, c
58     # word_2    d, e, f
59     # ...
60     # word_n    g, h, i
61     self.logger.debug(
62         'converting and loading words embeddings from text file {}'.format(
63             self.WORD_TO_VEC_MODEL_TXT_PATH
64         )
65     )
66     structure = [('word', np.dtype('U' + str(self.WORD_MAX_SIZE))),
67                 ('embedding', np.float32, (self.N_FEATURES,))]
68     structure = np.dtype(structure)
69     # load numpy array from disk using a generator
70     with open(self.WORD_TO_VEC_MODEL_TXT_PATH, encoding="utf8") as words_em
71         embeddings_gen = (
72             (line.split()[0], line.split()[1:]) for line in words_embedding
73             if len(line.split()[1:]) == self.N_FEATURES
74         )
75     embeddings = np.fromiter(embeddings_gen, structure)
76     # add a null embedding
77     null_embedding = np.array(
78         [('null_embedding', np.zeros((self.N_FEATURES,)), dtype=np.float32))
79         dtype=structure
80     )
81     embeddings = np.concatenate([embeddings, null_embedding])
82     # dump numpy array to disk using pickle
83     max_bytes = 2**28 - 1 # 256MB
84     bytes_out = pickle.dumps(embeddings, protocol=pickle.HIGHEST_PROTOCOL)
85     with open(self.PKL_PATH, 'wb') as f_out:
86         for idx in range(0, len(bytes_out), max_bytes):
87             f_out.write(bytes_out[idx:idx+max_bytes])
88     self.logger.debug('words embeddings loaded')
89     return embeddings

```

```

89         return embeddings
90
91
92 class GloveEmbeddings(WordsEmbeddings):
93
94     WORD_TO_VEC_MODEL_TXT_PATH = 'glove.twitter.27B.50d.txt'
95     PKL_PATH = 'gloveembedding.pkl'
96     N_FEATURES = 50
97     WORD_MAX_SIZE = 60
98
99
100 class FasttextEmbeddings(WordsEmbeddings):
101
102     WORD_TO_VEC_MODEL_TXT_PATH = 'cc.en.300.vec'
103     PKL_PATH = 'fasttext.pkl'
104     N_FEATURES = 300
105     WORD_MAX_SIZE = 60

```

a. Calcular el embedding de cada texto como el promedio de los embeddings de las palabras.

```

1 X = df.review
2 y = df.target

1 nltk.download('punkt')
2 from nltk.tokenize import word_tokenize

[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data]   Package punkt is already up-to-date!

1 embedder_model = FasttextEmbeddings()

1 def embed_corpus_average(corpus, embedder_model):
2     X_emb = np.zeros(shape=(len(corpus), embedder_model.N_FEATURES))
3     for i in range(len(corpus)):
4         X_emb[i] = np.average(
5             embedder_model.get_words_embeddings(
6                 word_tokenize(corpus[i])
7             ), axis=0)
8     return X_emb

1 X_emb = embed_corpus_average(X, embedder_model)

1 X_emb.shape

(50000, 300)

```

b. Crear los datasets de train, validation y test.

```

1 def split_dataset(X, y):
2     TRAIN_TEST_SPLIT = 0.2

```

```

2 TRAIN_TEST_SPLIT = 0.2
3 TRAIN_VAL_SPLIT = 0.25
4 X_train_set, X_test, y_train_set, y_test = train_test_split(X, y, test_size=T
5 X_train, X_val, y_train, y_val = train_test_split(X_train_set,y_train_set,tes
6 print("Test:",X_test.shape)
7 print("Train:",X_train.shape)
8 print("Val:",X_val.shape)
9 return X_train_set, X_test, y_train_set, y_test,X_train, X_val, y_train, y_va
10
11 X_train_set, X_test, y_train_set, y_test,X_train, X_val, y_train, y_val = split

    Test: (10000, 300)
    Train: (30000, 300)
    Val: (10000, 300)

```

- c. Entrenar una red neuronal (seleccionar arquitectura, loss y optimizador).
- d. Obtener AUC.

```

1 import sklearn
2 from sklearn.metrics import auc
3
4 def create_model(n_features):
5     model = Sequential()
6     model.add(Dense(50, activation = "relu", input_shape=(embedder_model.N_FEATUR
7     model.add(Dropout(0.3, noise_shape=None, seed=None))
8     model.add(Dense(50, activation = "relu"))
9     model.add(Dropout(0.2, noise_shape=None, seed=None))
10    model.add(Dense(50, activation = "relu"))
11    model.add(Dense(1, activation = "sigmoid"))
12    model.summary()
13    plot_model(model, show_shapes=True, show_layer_names=True)
14    return model
15
16 def train_and_evaluate_model(model):
17     model_metrics = [
18         keras.metrics.AUC(name="auc")
19     ]
20
21     # compiling the model
22     model.compile(
23         optimizer = "adam",
24         loss = "binary_crossentropy",
25         metrics = model_metrics
26     )
27
28     NUM_EPOCHS = 30
29
30     history = model.fit( X_train, y_train,
31                         epochs= NUM_EPOCHS,
32                         batch_size = 32,
33                         validation_data = (X_val, y_val),
34                         verbose = True )
35
36     fig,axes = plt.subplots(1,2,figsize=(24,4))

```

```
37 plt.suptitle("Entrenamiento")
38 axes[0].set_title("AUC")
39 axes[0].plot(np.arange(NUM_EPOCHS),history.history['auc'])
40 axes[0].plot(np.arange(NUM_EPOCHS),history.history['val_auc'])
41 axes[0].legend(["Training","Validation"])
42 axes[0].grid(which="Both")
43 axes[0].set_xlabel("Época")
44 axes[1].set_title("Loss")
45 axes[1].plot(np.arange(NUM_EPOCHS),history.history['loss'])
46 axes[1].plot(np.arange(NUM_EPOCHS),history.history['val_loss'])
47 axes[1].legend(["Training","Validation"])
48 axes[1].grid(which="Both")
49 axes[1].set_xlabel("Época")
50 plt.show()
51
52 y_test_pred = model.predict(X_test)
53
54 fpr, tpr, thresholds = sklearn.metrics.roc_curve(y_test.values, y_test_pred)
55 roc_auc = auc(fpr, tpr)
56
57 plt.figure(figsize=(24,4))
58 plt.title('Receiver Operating Characteristic')
59 plt.plot(fpr, tpr, 'b', label = 'AUC = %0.2f' % roc_auc)
60 plt.legend(loc = 'lower right')
61 plt.plot([0, 1], [0, 1], 'r--')
62 plt.xlim([0, 1])
63 plt.ylim([0, 1])
64 plt.grid(which="Both")
65 plt.ylabel('True Positive Rate')
66 plt.xlabel('False Positive Rate')
67 plt.show()
68
69 print("AUC:",roc_auc)
70 return roc_auc

1 model = create_model(n_features=embedder_model.N_FEATURES)
2 roc_auc = train_and_evaluate_model(model)
```

Model: "sequential\_1"

Layer (type)	Output Shape	Param #
dense_4 (Dense)	(None, 50)	15050
dropout_2 (Dropout)	(None, 50)	0
dense_5 (Dense)	(None, 50)	2550
dropout_3 (Dropout)	(None, 50)	0
dense_6 (Dense)	(None, 50)	2550
dense_7 (Dense)	(None, 1)	51

Total params: 20,201

Trainable params: 20,201

Non-trainable params: 0

Epoch 1/30

938/938 [=====] - 5s 4ms/step - loss: 0.5453 - auc:

Epoch 2/30

938/938 [=====] - 4s 4ms/step - loss: 0.3921 - auc:

Epoch 3/30

938/938 [=====] - 4s 4ms/step - loss: 0.3893 - auc:

Epoch 4/30

938/938 [=====] - 4s 5ms/step - loss: 0.3785 - auc:

Epoch 5/30

938/938 [=====] - 4s 4ms/step - loss: 0.3750 - auc:

Epoch 6/30

938/938 [=====] - 4s 4ms/step - loss: 0.3692 - auc:

Epoch 7/30

938/938 [=====] - 4s 4ms/step - loss: 0.3649 - auc:

Epoch 8/30

938/938 [=====] - 4s 4ms/step - loss: 0.3623 - auc:

Epoch 9/30

938/938 [=====] - 4s 4ms/step - loss: 0.3596 - auc:

Epoch 10/30

938/938 [=====] - 4s 4ms/step - loss: 0.3528 - auc:

Epoch 11/30

938/938 [=====] - 4s 4ms/step - loss: 0.3569 - auc:

Epoch 12/30

938/938 [=====] - 4s 4ms/step - loss: 0.3529 - auc:

Epoch 13/30

938/938 [=====] - 4s 4ms/step - loss: 0.3540 - auc:

Epoch 14/30

938/938 [=====] - 4s 4ms/step - loss: 0.3530 - auc:

Epoch 15/30

938/938 [=====] - 4s 4ms/step - loss: 0.3506 - auc:

Epoch 16/30

938/938 [=====] - 4s 4ms/step - loss: 0.3428 - auc:

Epoch 17/30

938/938 [=====] - 4s 5ms/step - loss: 0.3390 - auc:

Epoch 18/30

938/938 [=====] - 4s 4ms/step - loss: 0.3434 - auc:

Epoch 19/30

938/938 [=====] - 4s 4ms/step - loss: 0.3510 - auc:

Epoch 20/30

938/938 [=====] - 4s 4ms/step - loss: 0.3411 - auc:

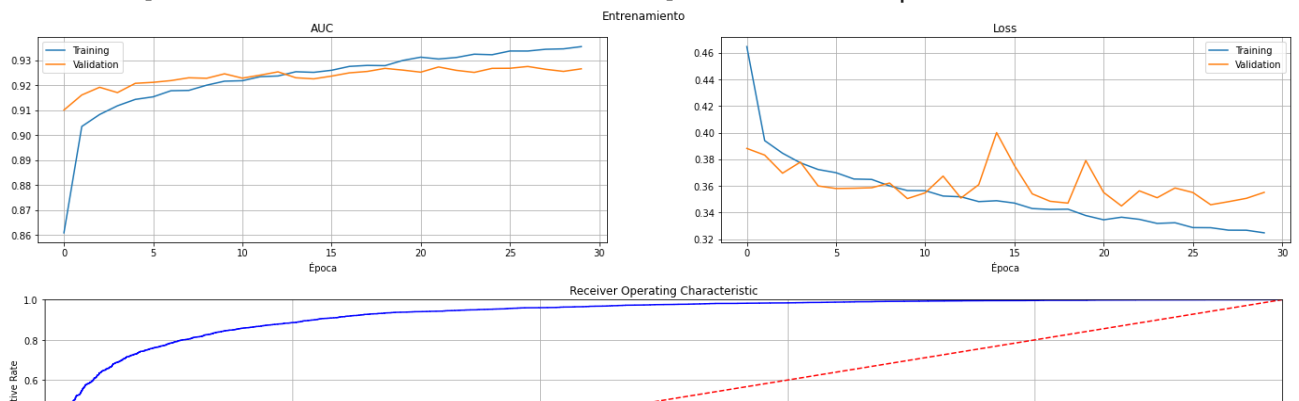
Epoch 21/30



```

938/938 [=====] - 4s 4ms/step - loss: 0.3390 - auc:
Epoch 22/30
938/938 [=====] - 4s 4ms/step - loss: 0.3348 - auc:
Epoch 23/30
938/938 [=====] - 4s 4ms/step - loss: 0.3346 - auc:
Epoch 24/30
938/938 [=====] - 4s 4ms/step - loss: 0.3277 - auc:
Epoch 25/30
938/938 [=====] - 4s 4ms/step - loss: 0.3318 - auc:
Epoch 26/30
938/938 [=====] - 4s 4ms/step - loss: 0.3302 - auc:
Epoch 27/30
938/938 [=====] - 4s 4ms/step - loss: 0.3286 - auc:
Epoch 28/30
938/938 [=====] - 4s 4ms/step - loss: 0.3250 - auc:
Epoch 29/30
938/938 [=====] - 4s 5ms/step - loss: 0.3259 - auc:
Epoch 30/30
938/938 [=====] - 4s 4ms/step - loss: 0.3241 - auc:

```



```

1 model_comparison_table["FastText (Average)"] = {
2     "AUC": 0.9237549829902969, #roc_auc,
3     "Description": "FastText 300 (cc.en.300.vec.gz)"
4 }

```

e. Tratar de mejorar los embeddings de los textos con:

- Agregar TFIDF como pesos a las palabras.

```

1 X = df.review.values
2 y = df.target.values
3 X_train_set, X_test, y_train_set, y_test, X_train, X_val, y_train, y_val = spli

Test: (10000,)
Train: (30000,)
Val: (10000,)

```

```

1 from sklearn.feature_extraction.text import TfidfVectorizer
2 tfidf_vectorizer = TfidfVectorizer(tokenizer=word_tokenize)

```

Entrenar el TF-IDF Vectorizer con el Train set y obtener los pesos para el corpus.

```

1 X_train_tfidf = tfidf_vectorizer.fit_transform(X_train)

```

```

2 X_val_tfidf = tfidf_vectorizer.transform(X_val)
3 X_test_tfidf = tfidf_vectorizer.transform(X_test)
4 feature_names = tfidf_vectorizer.get_feature_names()
5 inverse_feature_names = { v:k for k, v in enumerate(feature_names)}
6 len(feature_names)

```

82157

La siguiente función devuelve los pesos de las palabras de un documento.

```

1 def doc_to_weights(doc_idx,doc_matrix, tfidf_doc_matrix):
2     """
3     doc_idx: índice del documento en ambas matrices
4     doc_matrix: matriz con documentos (texto)
5     tfidf_doc_matrix: matriz TF-IDF
6     """
7     tokenized_doc = word_tokenize(doc_matrix[doc_idx])
8     indexed_terms = []
9     for term in tokenized_doc:
10         indexed_terms.append(inverse_feature_names[term] if term in inverse_feature
11 doc_weights = [tfidf_doc_matrix[doc_idx].toarray()[:,token_index][0] for toke
12 return np.array(doc_weights)
13
14 x0_w = doc_to_weights(1,X_train,X_train_tfidf).reshape(-1,1)
15 x0_w.shape

```

(43, 1)

```

1 def embed_corpus_tfidf_weight(corpus,corpus_tfidf,embedder_model):
2     X_emb = np.zeros(shape=(len(corpus),embedder_model.N_FEATURES))
3     for i in range(len(corpus)):
4         x_w = doc_to_weights(i,corpus,corpus_tfidf).reshape(-1,1)
5         x_e = embedder_model.get_words_embeddings( word_tokenize(corpus[i]))
6         X_emb[i] = np.sum(x_e*x_w,axis=0)/x_w.shape[0]
7     return X_emb
8
9 X_train = embed_corpus_tfidf_weight(X_train,X_train_tfidf,embedder_model)
10 X_val = embed_corpus_tfidf_weight(X_val,X_val_tfidf,embedder_model)
11 X_test = embed_corpus_tfidf_weight(X_test,X_test_tfidf,embedder_model)

```

```

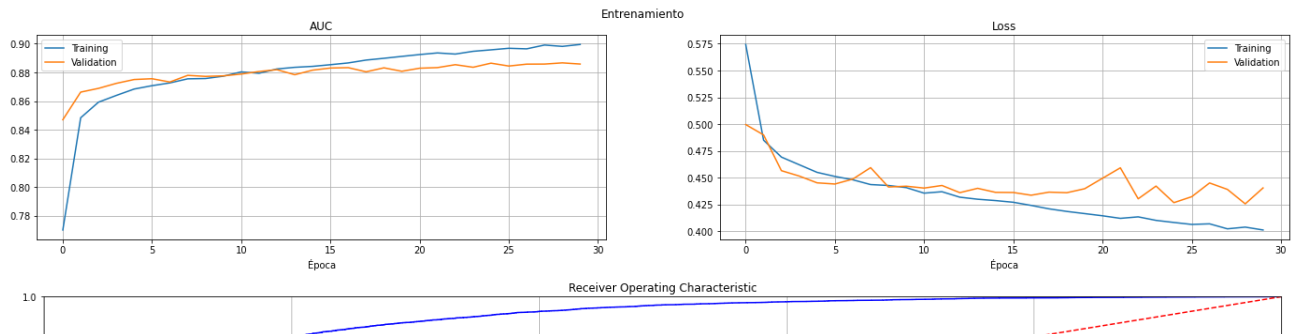
1 model = create_model(n_features=embedder_model.N_FEATURES)
2 model.compile(
3     optimizer = "adam",
4     loss = "binary_crossentropy",
5     metrics = [ keras.metrics.AUC(name="auc")])
6 )
7
8 NUM_EPOCHS = 30
9
10 history = model.fit( X_train, y_train,
11                     epochs= NUM_EPOCHS,
12                     batch_size = 32,
13                     validation_data = (X_val, y_val))

```

```
validation_data = (x_val, y_val),
verbose = True )
```

```
938/938 [=====] - 5s 5ms/step - loss: 0.6355 - auc
Epoch 2/30
938/938 [=====] - 4s 4ms/step - loss: 0.4897 - auc
Epoch 3/30
938/938 [=====] - 4s 4ms/step - loss: 0.4710 - auc
Epoch 4/30
938/938 [=====] - 4s 4ms/step - loss: 0.4639 - auc
Epoch 5/30
938/938 [=====] - 4s 4ms/step - loss: 0.4564 - auc
Epoch 6/30
938/938 [=====] - 4s 4ms/step - loss: 0.4546 - auc
Epoch 7/30
938/938 [=====] - 4s 4ms/step - loss: 0.4499 - auc
Epoch 8/30
938/938 [=====] - 4s 4ms/step - loss: 0.4443 - auc
Epoch 9/30
938/938 [=====] - 4s 5ms/step - loss: 0.4441 - auc
Epoch 10/30
938/938 [=====] - 4s 5ms/step - loss: 0.4412 - auc
Epoch 11/30
938/938 [=====] - 4s 4ms/step - loss: 0.4328 - auc
Epoch 12/30
938/938 [=====] - 4s 5ms/step - loss: 0.4401 - auc
Epoch 13/30
938/938 [=====] - 4s 4ms/step - loss: 0.4285 - auc
Epoch 14/30
938/938 [=====] - 4s 4ms/step - loss: 0.4331 - auc
Epoch 15/30
938/938 [=====] - 4s 4ms/step - loss: 0.4313 - auc
Epoch 16/30
938/938 [=====] - 4s 4ms/step - loss: 0.4255 - auc
Epoch 17/30
938/938 [=====] - 4s 4ms/step - loss: 0.4210 - auc
Epoch 18/30
938/938 [=====] - 4s 4ms/step - loss: 0.4187 - auc
Epoch 19/30
938/938 [=====] - 4s 4ms/step - loss: 0.4203 - auc
Epoch 20/30
938/938 [=====] - 4s 5ms/step - loss: 0.4114 - auc
Epoch 21/30
938/938 [=====] - 4s 4ms/step - loss: 0.4093 - auc
Epoch 22/30
938/938 [=====] - 4s 4ms/step - loss: 0.4164 - auc
Epoch 23/30
938/938 [=====] - 4s 4ms/step - loss: 0.4140 - auc
Epoch 24/30
938/938 [=====] - 4s 4ms/step - loss: 0.4097 - auc
Epoch 25/30
938/938 [=====] - 4s 4ms/step - loss: 0.4058 - auc
Epoch 26/30
938/938 [=====] - 4s 4ms/step - loss: 0.4032 - auc
Epoch 27/30
938/938 [=====] - 4s 5ms/step - loss: 0.4059 - auc
Epoch 28/30
938/938 [=====] - 4s 4ms/step - loss: 0.4036 - auc
Epoch 29/30
938/938 [=====] - 4s 4ms/step - loss: 0.4043 - auc
Epoch 30/30
```

```
1 fig, axes = plt.subplots(1, 2, figsize=(24, 4))
2 plt.suptitle("Entrenamiento")
3 axes[0].set_title("AUC")
4 axes[0].plot(np.arange(NUM_EPOCHS), history.history['auc'])
5 axes[0].plot(np.arange(NUM_EPOCHS), history.history['val_auc'])
6 axes[0].legend(["Training", "Validation"])
7 axes[0].grid(which="Both")
8 axes[0].set_xlabel("Época")
9 axes[1].set_title("Loss")
10 axes[1].plot(np.arange(NUM_EPOCHS), history.history['loss'])
11 axes[1].plot(np.arange(NUM_EPOCHS), history.history['val_loss'])
12 axes[1].legend(["Training", "Validation"])
13 axes[1].grid(which="Both")
14 axes[1].set_xlabel("Época")
15 plt.show()
16
17 y_test_pred = model.predict(X_test)
18
19 fpr, tpr, thresholds = sklearn.metrics.roc_curve(y_test, y_test_pred)
20 roc_auc = auc(fpr, tpr)
21
22 plt.figure(figsize=(24, 4))
23 plt.title('Receiver Operating Characteristic')
24 plt.plot(fpr, tpr, 'b', label = 'AUC = %0.2f' % roc_auc)
25 plt.legend(loc = 'lower right')
26 plt.plot([0, 1], [0, 1], 'r--')
27 plt.xlim([0, 1])
28 plt.ylim([0, 1])
29 plt.grid(which="Both")
30 plt.ylabel('True Positive Rate')
31 plt.xlabel('False Positive Rate')
32 plt.show()
33
34 print("AUC:", roc_auc)
```



```
1 model_comparison_table["FastText (TF-IDF)"] = {
2     "AUC": 0.8834895915067472, # roc_auc,
3     "Description": "FastText 300 (cc.en.300.vec.gz)"
4 }
```

```
AUC: 0.8834895915067472
```

- En lugar de tomar el promedio probar tomando el max o el min de cada componente de los embeddings de las palabras.

```
1 def embed_corpus_max(corpus, embedder_model):
2     X_emb = np.zeros(shape=(len(corpus), embedder_model.N_FEATURES))
3     for i in range(len(corpus)):
4         X_emb[i] = np.max(
5             embedder_model.get_words_embeddings(
6                 word_tokenize(corpus[i])
7             ), axis=0)
8     return X_emb
9
10 def embed_corpus_min(corpus, embedder_model):
11     X_emb = np.zeros(shape=(len(corpus), embedder_model.N_FEATURES))
12     for i in range(len(corpus)):
13         X_emb[i] = np.min(
14             embedder_model.get_words_embeddings(
15                 word_tokenize(corpus[i])
16             ), axis=0)
17     return X_emb

1 X_emb = embed_corpus_max(X, embedder_model)
2 X_train_set, X_test, y_train_set, y_test, X_train, X_val, y_train, y_val = spli
3 model = create_model(n_features=embedder_model.N_FEATURES)
4 roc_auc = train_and_evaluate_model(model)
```

Test: (10000, 300)  
 Train: (30000, 300)  
 Val: (10000, 300)  
 Model: "sequential\_2"

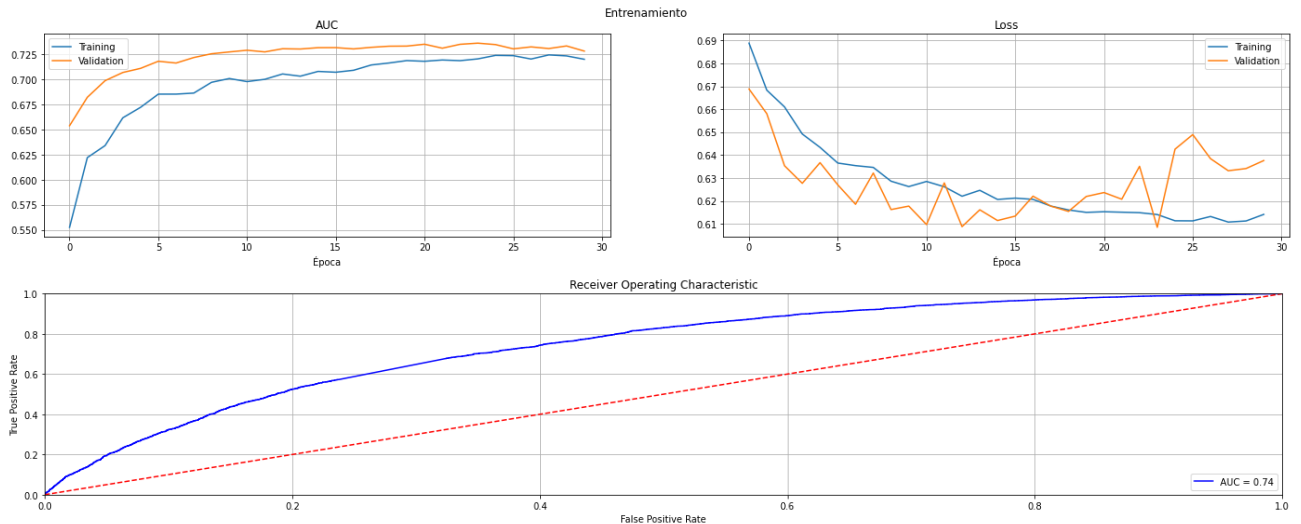
Layer (type)	Output Shape	Param #
dense_8 (Dense)	(None, 50)	15050
dropout_4 (Dropout)	(None, 50)	0
dense_9 (Dense)	(None, 50)	2550
dropout_5 (Dropout)	(None, 50)	0
dense_10 (Dense)	(None, 50)	2550
dense_11 (Dense)	(None, 1)	51
Total params: 20,201		
Trainable params: 20,201		
Non-trainable params: 0		

```
Epoch 1/30
938/938 [=====] - 5s 4ms/step - loss: 0.6926 - auc:
Epoch 2/30
938/938 [=====] - 4s 4ms/step - loss: 0.6733 - auc:
Epoch 3/30
938/938 [=====] - 4s 4ms/step - loss: 0.6614 - auc:
Epoch 4/30
938/938 [=====] - 4s 4ms/step - loss: 0.6516 - auc:
Epoch 5/30
938/938 [=====] - 4s 4ms/step - loss: 0.6438 - auc:
Epoch 6/30
938/938 [=====] - 4s 4ms/step - loss: 0.6348 - auc:
Epoch 7/30
938/938 [=====] - 4s 4ms/step - loss: 0.6365 - auc:
Epoch 8/30
938/938 [=====] - 4s 4ms/step - loss: 0.6339 - auc:
Epoch 9/30
938/938 [=====] - 4s 4ms/step - loss: 0.6324 - auc:
Epoch 10/30
938/938 [=====] - 4s 4ms/step - loss: 0.6303 - auc:
Epoch 11/30
938/938 [=====] - 4s 4ms/step - loss: 0.6281 - auc:
Epoch 12/30
938/938 [=====] - 4s 4ms/step - loss: 0.6221 - auc:
Epoch 13/30
938/938 [=====] - 4s 4ms/step - loss: 0.6189 - auc:
Epoch 14/30
938/938 [=====] - 4s 4ms/step - loss: 0.6259 - auc:
Epoch 15/30
938/938 [=====] - 4s 4ms/step - loss: 0.6185 - auc:
Epoch 16/30
938/938 [=====] - 4s 4ms/step - loss: 0.6221 - auc:
Epoch 17/30
938/938 [=====] - 4s 4ms/step - loss: 0.6212 - auc:
Epoch 18/30
938/938 [=====] - 4s 4ms/step - loss: 0.6210 - auc:
Epoch 19/30
938/938 [=====] - 4s 4ms/step - loss: 0.6175 - auc:
```

```

Epoch 20/30
938/938 [=====] - 4s 4ms/step - loss: 0.6143 - auc:
Epoch 21/30
938/938 [=====] - 4s 4ms/step - loss: 0.6138 - auc:
Epoch 22/30
938/938 [=====] - 4s 4ms/step - loss: 0.6148 - auc:
Epoch 23/30
938/938 [=====] - 4s 4ms/step - loss: 0.6150 - auc:
Epoch 24/30
938/938 [=====] - 4s 4ms/step - loss: 0.6142 - auc:
Epoch 25/30
938/938 [=====] - 4s 4ms/step - loss: 0.6085 - auc:
Epoch 26/30
938/938 [=====] - 4s 4ms/step - loss: 0.6114 - auc:
Epoch 27/30
938/938 [=====] - 4s 4ms/step - loss: 0.6143 - auc:
Epoch 28/30
938/938 [=====] - 4s 4ms/step - loss: 0.6110 - auc:
Epoch 29/30
938/938 [=====] - 4s 4ms/step - loss: 0.6116 - auc:
Epoch 30/30
938/938 [=====] - 4s 4ms/step - loss: 0.6161 - auc:

```



AUC: 0.7392006923537724

```

1 model_comparison_table["FastText (Max)"] = {
2     "AUC": 0.7392006923537724, #roc_auc
3     "Description": "FastText 300 (cc.en.300.vec.gz)"
4 }

```

```

1 X_emb = embed_corpus_min(X, embedder_model)
2 X_train_set, X_test, y_train_set, y_test, X_train, X_val, y_train, y_val = spli
3 model = create_model(n_features=embedder_model.N_FEATURES)
4 roc_auc = train_and_evaluate_model(model)

```

Test: (10000, 300)  
 Train: (30000, 300)  
 Val: (10000, 300)  
 Model: "sequential\_3"

Layer (type)	Output Shape	Param #
dense_12 (Dense)	(None, 50)	15050
dropout_6 (Dropout)	(None, 50)	0
dense_13 (Dense)	(None, 50)	2550
dropout_7 (Dropout)	(None, 50)	0
dense_14 (Dense)	(None, 50)	2550
dense_15 (Dense)	(None, 1)	51
Total params: 20,201		
Trainable params: 20,201		
Non-trainable params: 0		

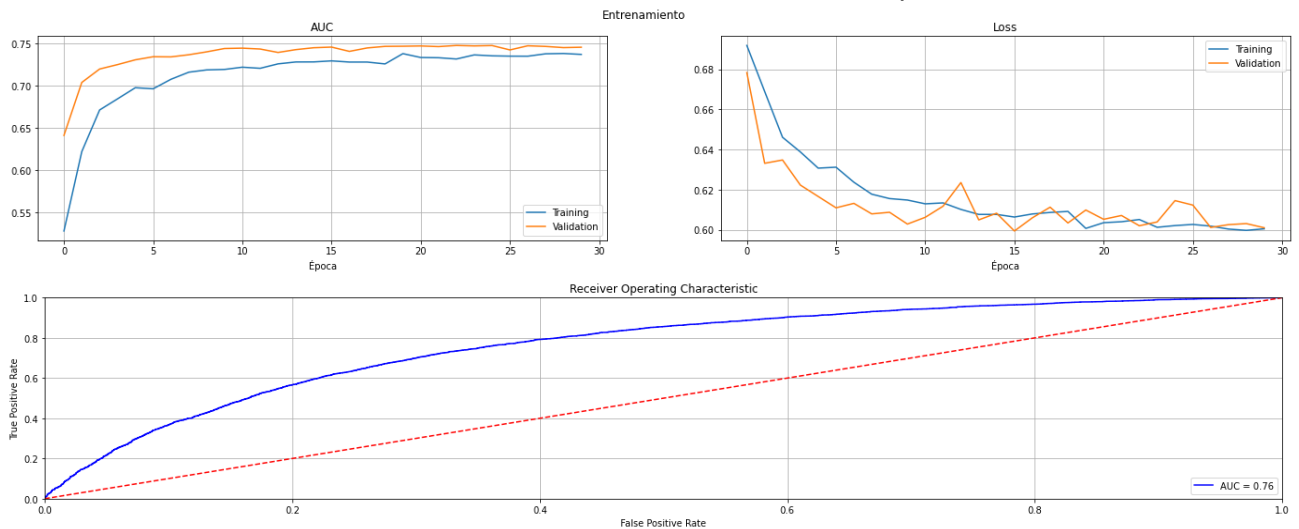
```
Epoch 1/30
938/938 [=====] - 5s 4ms/step - loss: 0.6940 - auc:
Epoch 2/30
938/938 [=====] - 4s 4ms/step - loss: 0.6767 - auc:
Epoch 3/30
938/938 [=====] - 4s 4ms/step - loss: 0.6466 - auc:
Epoch 4/30
938/938 [=====] - 4s 4ms/step - loss: 0.6416 - auc:
Epoch 5/30
938/938 [=====] - 4s 4ms/step - loss: 0.6305 - auc:
Epoch 6/30
938/938 [=====] - 4s 4ms/step - loss: 0.6355 - auc:
Epoch 7/30
938/938 [=====] - 4s 4ms/step - loss: 0.6295 - auc:
Epoch 8/30
938/938 [=====] - 4s 4ms/step - loss: 0.6162 - auc:
Epoch 9/30
938/938 [=====] - 4s 4ms/step - loss: 0.6158 - auc:
Epoch 10/30
938/938 [=====] - 4s 4ms/step - loss: 0.6161 - auc:
Epoch 11/30
938/938 [=====] - 4s 4ms/step - loss: 0.6126 - auc:
Epoch 12/30
938/938 [=====] - 4s 4ms/step - loss: 0.6105 - auc:
Epoch 13/30
938/938 [=====] - 4s 4ms/step - loss: 0.6082 - auc:
Epoch 14/30
938/938 [=====] - 4s 4ms/step - loss: 0.6082 - auc:
Epoch 15/30
938/938 [=====] - 4s 4ms/step - loss: 0.6071 - auc:
Epoch 16/30
938/938 [=====] - 4s 4ms/step - loss: 0.6040 - auc:
Epoch 17/30
938/938 [=====] - 4s 4ms/step - loss: 0.6105 - auc:
Epoch 18/30
938/938 [=====] - 4s 4ms/step - loss: 0.6125 - auc:
Epoch 19/30
938/938 [=====] - 4s 4ms/step - loss: 0.6151 - auc:
```



```

Epoch 20/30
938/938 [=====] - 4s 4ms/step - loss: 0.5981 - auc:
Epoch 21/30
938/938 [=====] - 4s 4ms/step - loss: 0.6060 - auc:
Epoch 22/30
938/938 [=====] - 4s 4ms/step - loss: 0.6057 - auc:
Epoch 23/30
938/938 [=====] - 4s 4ms/step - loss: 0.6066 - auc:
Epoch 24/30
938/938 [=====] - 4s 4ms/step - loss: 0.6012 - auc:
Epoch 25/30
938/938 [=====] - 4s 4ms/step - loss: 0.6052 - auc:
Epoch 26/30
938/938 [=====] - 4s 4ms/step - loss: 0.6002 - auc:
Epoch 27/30
938/938 [=====] - 4s 4ms/step - loss: 0.5980 - auc:
Epoch 28/30
938/938 [=====] - 4s 4ms/step - loss: 0.6036 - auc:
Epoch 29/30
938/938 [=====] - 4s 4ms/step - loss: 0.6026 - auc:
Epoch 30/30
938/938 [=====] - 4s 4ms/step - loss: 0.6023 - auc:

```



AUC: 0.7621117737390766

```

1 model_comparison_table["FastText (Min)"] = {
2     "AUC": 0.7621117737390766, #roc_auc,
3     "Description": "FastText 300 (cc.en.300.vec.gz)"
4 }

```

- Probar diferentes tamaños de embeddings.

```

1 embedder_model = GloveEmbeddings()
2 X_emb = embed_corpus_average(X, embedder_model)
3 X_train_set, X_test, y_train_set, y_test, X_train, X_val, y_train, y_val = spli
4 model = create_model(n_features=embedder_model.N_FEATURES)
5 roc_auc = train_and_evaluate_model(model)

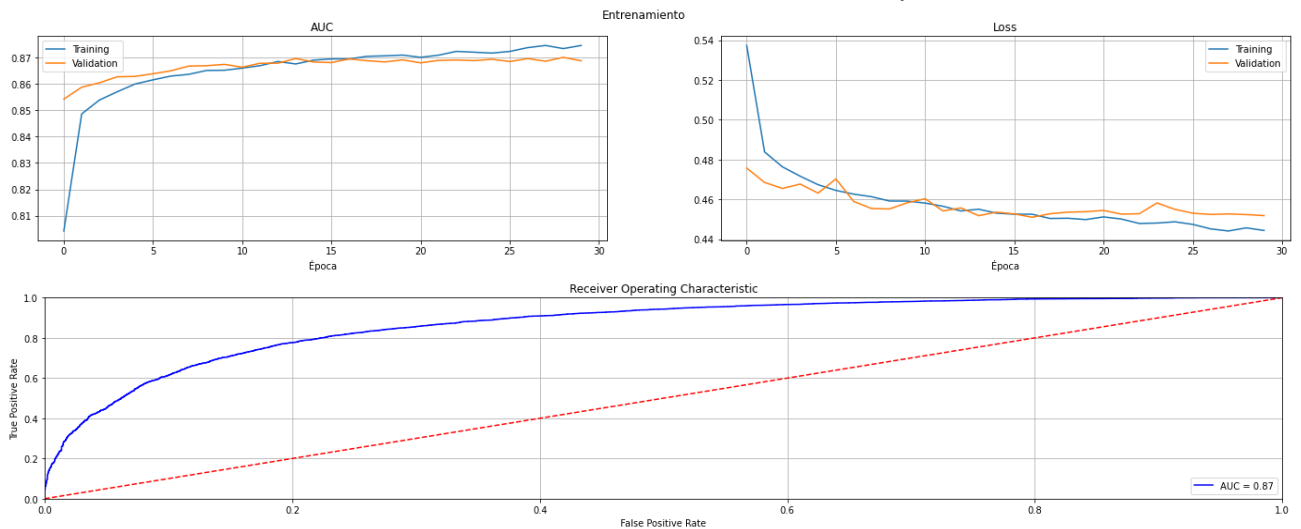
```

Test: (10000, 50)  
 Train: (30000, 50)  
 Val: (10000, 50)  
 Model: "sequential\_4"

Layer (type)	Output Shape	Param #
dense_16 (Dense)	(None, 50)	2550
dropout_8 (Dropout)	(None, 50)	0
dense_17 (Dense)	(None, 50)	2550
dropout_9 (Dropout)	(None, 50)	0
dense_18 (Dense)	(None, 50)	2550
dense_19 (Dense)	(None, 1)	51
Total params: 7,701		
Trainable params: 7,701		
Non-trainable params: 0		

```
Epoch 1/30
938/938 [=====] - 5s 5ms/step - loss: 0.5966 - auc:
Epoch 2/30
938/938 [=====] - 4s 4ms/step - loss: 0.4832 - auc:
Epoch 3/30
938/938 [=====] - 4s 4ms/step - loss: 0.4755 - auc:
Epoch 4/30
938/938 [=====] - 4s 4ms/step - loss: 0.4736 - auc:
Epoch 5/30
938/938 [=====] - 4s 4ms/step - loss: 0.4654 - auc:
Epoch 6/30
938/938 [=====] - 4s 4ms/step - loss: 0.4645 - auc:
Epoch 7/30
938/938 [=====] - 4s 4ms/step - loss: 0.4577 - auc:
Epoch 8/30
938/938 [=====] - 4s 4ms/step - loss: 0.4599 - auc:
Epoch 9/30
938/938 [=====] - 4s 5ms/step - loss: 0.4550 - auc:
Epoch 10/30
938/938 [=====] - 4s 4ms/step - loss: 0.4577 - auc:
Epoch 11/30
938/938 [=====] - 4s 5ms/step - loss: 0.4530 - auc:
Epoch 12/30
938/938 [=====] - 4s 4ms/step - loss: 0.4597 - auc:
Epoch 13/30
938/938 [=====] - 4s 4ms/step - loss: 0.4492 - auc:
Epoch 14/30
938/938 [=====] - 4s 4ms/step - loss: 0.4550 - auc:
Epoch 15/30
938/938 [=====] - 4s 4ms/step - loss: 0.4536 - auc:
Epoch 16/30
938/938 [=====] - 4s 4ms/step - loss: 0.4516 - auc:
Epoch 17/30
938/938 [=====] - 4s 4ms/step - loss: 0.4560 - auc:
Epoch 18/30
938/938 [=====] - 4s 4ms/step - loss: 0.4532 - auc:
Epoch 19/30
938/938 [=====] - 4s 4ms/step - loss: 0.4567 - auc:
```

Epoch 20/30  
 938/938 [=====] - 4s 4ms/step - loss: 0.4485 - auc:  
 Epoch 21/30  
 938/938 [=====] - 4s 4ms/step - loss: 0.4478 - auc:  
 Epoch 22/30  
 938/938 [=====] - 4s 4ms/step - loss: 0.4498 - auc:  
 Epoch 23/30  
 938/938 [=====] - 4s 4ms/step - loss: 0.4515 - auc:  
 Epoch 24/30  
 938/938 [=====] - 4s 4ms/step - loss: 0.4514 - auc:  
 Epoch 25/30  
 938/938 [=====] - 4s 4ms/step - loss: 0.4468 - auc:  
 Epoch 26/30  
 938/938 [=====] - 4s 4ms/step - loss: 0.4421 - auc:  
 Epoch 27/30  
 938/938 [=====] - 4s 4ms/step - loss: 0.4455 - auc:  
 Epoch 28/30  
 938/938 [=====] - 4s 4ms/step - loss: 0.4500 - auc:  
 Epoch 29/30  
 938/938 [=====] - 4s 4ms/step - loss: 0.4464 - auc:  
 Epoch 30/30  
 938/938 [=====] - 4s 4ms/step - loss: 0.4459 - auc:



AUC: 0.8682643656506908

```
1 model_comparison_table["GloVE"] = {
2     "AUC": 0.8682643656506908, #roc_auc,
3     "Description": "GloVE 50 (glove.twitter.27B.zip)"
4 }
```

## ▼ 5. Clasificador con LSTM

Entrenar un clasificador utilizando celdas LSTM.

- a. Convertir secuencias de palabras a secuencias de números (indexer).

- b. Agregar padding para que cada elemento de entrenamiento tenga la misma longitud.
- c. Armar una red LSTM con las siguientes capas:
  - Capa de embedding que transforma un número (index) en un embedding.
  - Agregar uno dos layers LSTM para obtener el embedding de la secuencias.
  - Agregar un layer denso para entrenar el clasificador.

```
1 X = df.review
2 y = df.target
```

```
1 from sklearn.model_selection import train_test_split
2
3 TRAIN_TEST_SPLIT = 0.2
4 TRAIN_VAL_SPLIT = 0.25
5
6 X_train_set, X_test, y_train_set, y_test = train_test_split(X, y, test_size=TRA
7 X_train, X_val, y_train, y_val = train_test_split(X_train_set,y_train_set,test_
8
9 print("Test:",X_test.shape)
10 print("Train:",X_train.shape)
11 print("Val:",X_val.shape)

    Test: (10000,)
    Train: (30000,)
    Val: (10000,)
```

```
1 NUM_WORDS = 20000
2 EMBEDDING_DIM = 300
3 SEQ_LENGTH = 200
4 PADDING_TYPE='post'
5 TRUNC_TYPE='post'
```

```
1 from keras.preprocessing.text import Tokenizer
2 from keras.preprocessing.sequence import pad_sequences
3
4 tokenizer = Tokenizer(num_words=NUM_WORDS)
5 tokenizer.fit_on_texts(X_train)
```

```
1 X_train_seq = tokenizer.texts_to_sequences(X_train.values)
2 len(X_train_seq)
```

```
30000
```

```
1 X_val_seq = tokenizer.texts_to_sequences(X_val.values)
2 len(X_val_seq)
```

```
10000
```

```
1 X_train_padded = pad_sequences(X_train_seq,maxlen=SEQ_LENGTH,padding=PADDING_TY
2 X_train_padded[0]
```

```
array([ 2, 1457, 45, 20, 1849, 8669, 108, 19, 20,
       975, 20, 1133, 1596, 20, 14595, 20, 2817, 107,
       432, 3899, 12386, 103, 730, 1894, 2921, 1023, 1495,
       4065, 2870, 4213, 2165, 10568, 3617, 3175, 55, 13,
       571, 370, 2234, 45, 690, 1067, 9428, 4647, 55,
       13, 2294, 1270, 1, 6, 14, 1397, 1825, 2119,
       1379, 190, 5526, 1379, 1849, 3110, 8669, 306, 519,
       255, 5031, 367, 1468, 64, 89, 1240, 838, 6,
       1002, 3155, 1044, 2772, 109, 679, 11510, 1497, 6063,
       0, 0, 0, 0, 0, 0, 0, 0, 0,
       0, 0, 0, 0, 0, 0, 0, 0, 0,
       0, 0, 0, 0, 0, 0, 0, 0, 0,
       0, 0, 0, 0, 0, 0, 0, 0, 0,
       0, 0, 0, 0, 0, 0, 0, 0, 0,
       0, 0, 0, 0, 0, 0, 0, 0, 0,
       0, 0, 0, 0, 0, 0, 0, 0, 0,
       0, 0, 0, 0, 0, 0, 0, 0, 0,
       0, 0, 0, 0, 0, 0, 0, 0, 0,
       0, 0, 0, 0, 0, 0, 0, 0, 0,
       0, 0, 0, 0, 0, 0, 0, 0, 0,
       0, 0], dtype=int32)
```

```
1 X_val_padded = pad_sequences(X_val_seq,maxlen=SEQ_LENGTH,padding=PADDING_TYPE,
2 X_val_padded[0]
```

```
array([ 1, 237, 5663, 5361, 292, 2382, 62, 1394, 35,
       83, 4, 624, 148, 23, 10, 5361, 4153, 4,
       5088, 1085, 1150, 1, 23, 120, 14, 717, 5361,
       4065, 2870, 1676, 772, 6718, 1922, 916, 2, 5,
       158, 409, 3617, 120, 2927, 15458, 1, 8782, 5676,
       10622, 3468, 541, 845, 707, 314, 23, 714, 592,
       739, 16495, 420, 1597, 2709, 135, 34, 7040, 314,
       714, 83, 6040, 138, 38, 4678, 2615, 34, 221,
       1005, 120, 1286, 49, 558, 34, 24, 13, 20,
       1259, 8057, 4866, 4647, 1258, 10622, 115, 30, 4678,
       2988, 16495, 1625, 1686, 2140, 5937, 572, 2688, 297,
       253, 62, 15244, 2775, 615, 100, 3124, 245, 4409,
       253, 2080, 115, 669, 605, 3674, 1205, 1491, 7765,
       236, 2878, 9673, 2937, 32, 2, 174, 132, 16635,
       2, 4678, 8913, 2614, 355, 8039, 1952, 3843, 572,
       1113, 187, 212, 381, 1, 716, 825, 1492, 11700,
       6192, 5, 10, 1, 1777, 2, 3416, 9, 1,
       716, 1, 330, 51, 551, 2, 1, 284, 147,
       14, 2, 1306, 481, 19059, 39, 1, 5514, 2384,
       0, 0, 0, 0, 0, 0, 0, 0, 0,
       0, 0, 0, 0, 0, 0, 0, 0, 0,
       0, 0, 0, 0, 0, 0, 0, 0, 0,
       0, 0], dtype=int32)
```

```
1 from tensorflow.keras.models import Sequential
2 from tensorflow.keras.layers import Dense, Flatten, LSTM, Dropout, Activation,
3
4 def create_model(vocab_size,embedding_dim):
5     model = Sequential()
6
7     model.add(Embedding(vocab_size,embedding_dim))
```

```

8  model.add(Dropout(0.5))
9  model.add(Bidirectional(LSTM(embedding_dim)))
10 model.add(Dense(1,activation='sigmoid'))
11 return model
12
13
14 model = create_model(NUM_WORDS,EMBEDDING_DIM)
15 model.summary()

```

Model: "sequential\_5"

Layer (type)	Output Shape	Param #
embedding (Embedding)	(None, None, 300)	6000000
dropout_10 (Dropout)	(None, None, 300)	0
bidirectional (Bidirectional)	(None, 600)	1442400
dense_20 (Dense)	(None, 1)	601
Total params: 7,443,001		
Trainable params: 7,443,001		
Non-trainable params: 0		

```

1 from tensorflow.keras.optimizers import Adam
2 from tensorflow.keras.metrics import AUC
3
4 model.compile(
5     loss='binary_crossentropy',
6     optimizer=Adam(lr=0.001, decay=1e-6),
7     metrics=[AUC(name="auc")],
8 )

```

```

1 NUM_EPOCHS = 10
2 history = model.fit(
3     X_train_padded,
4     y_train,
5     epochs=NUM_EPOCHS,
6     validation_data=(X_val_padded, y_val),
7     verbose=True)

```

```

Epoch 1/10
938/938 [=====] - 121s 120ms/step - loss: 0.6030 - a
Epoch 2/10
938/938 [=====] - 112s 120ms/step - loss: 0.3128 - a
Epoch 3/10
938/938 [=====] - 112s 120ms/step - loss: 0.1767 - a
Epoch 4/10
938/938 [=====] - 112s 120ms/step - loss: 0.1043 - a
Epoch 5/10
938/938 [=====] - 112s 119ms/step - loss: 0.0632 - a
Epoch 6/10
938/938 [=====] - 112s 119ms/step - loss: 0.0357 - a
Epoch 7/10

```

```

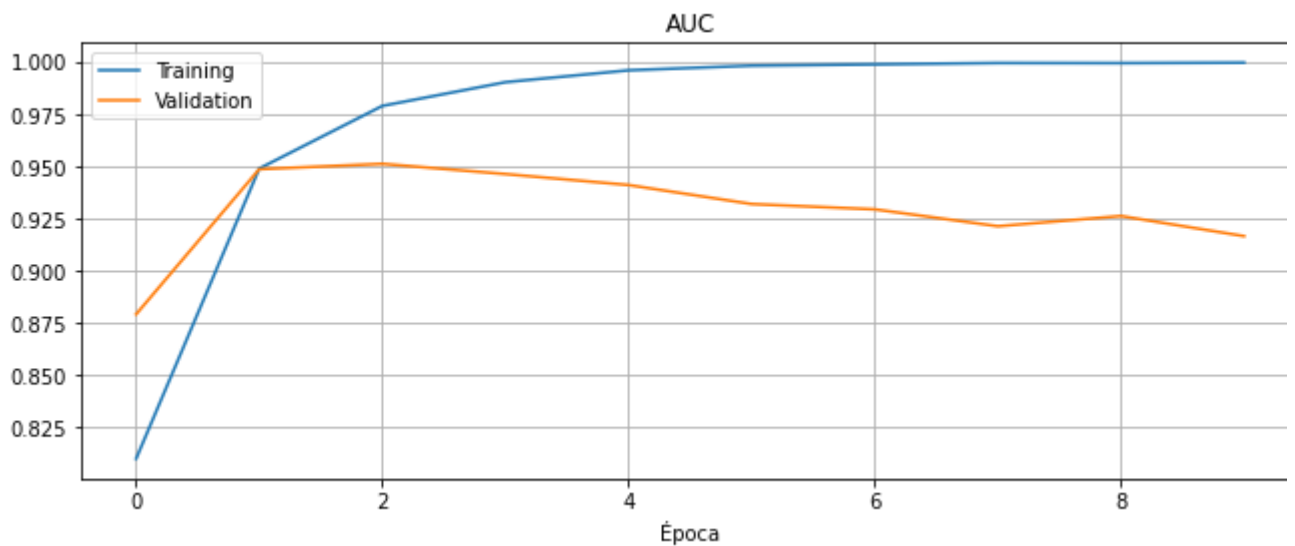
938/938 [=====] - 113s 120ms/step - loss: 0.0272 - a
Epoch 8/10
938/938 [=====] - 112s 120ms/step - loss: 0.0154 - a
Epoch 9/10
938/938 [=====] - 112s 120ms/step - loss: 0.0111 - a
Epoch 10/10
938/938 [=====] - 112s 119ms/step - loss: 0.0077 - a

```

```

1  fig,axes = plt.subplots(1,2,figsize=(24,4))
2  plt.suptitle("LSTM")
3  axes[0].set_title("AUC")
4  axes[0].plot(np.arange(NUM_EPOCHS),history.history['auc'])
5  axes[0].plot(np.arange(NUM_EPOCHS),history.history['val_auc'])
6  axes[0].legend(["Training","Validation"])
7  axes[0].grid(which="Both")
8  axes[0].set_xlabel("Época")
9  axes[1].set_title("Loss")
10 axes[1].plot(np.arange(NUM_EPOCHS),history.history['loss'])
11 axes[1].plot(np.arange(NUM_EPOCHS),history.history['val_loss'])
12 axes[1].legend(["Training","Validation"])
13 axes[1].grid(which="Both")
14 axes[1].set_xlabel("Época")
15 plt.show()

```



Se reentrena con una época para obtener el mejor modelo.

```

1 NUM_EPOCHS = 1
2 model = create_model(NUM_WORDS,EMBEDDING_DIM)
3 model.compile(loss='binary_crossentropy',optimizer=Adam(lr=0.001, decay=1e-6),
4               metrics=[AUC(name="auc")])
5 history = model.fit(X_train_padded, y_train, epochs=NUM_EPOCHS,
6                    validation_data=(X_val_padded, y_val), verbose=True)

938/938 [=====] - 95s 99ms/step - loss: 0.5153 - auc

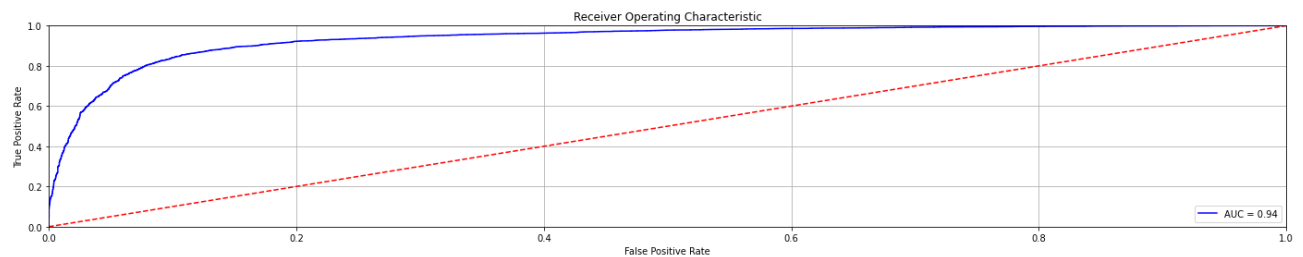
```

```

1 X_test_seq = tokenizer.texts_to_sequences(X_test.values)
2 X_test_padded = pad_sequences(X_test_seq,maxlen=SEQ_LENGTH,padding=PADDING_TYPE)
3 y_test_pred = model.predict(X_test_padded)

1 import sklearn
2 from sklearn.metrics import auc
3
4 fpr, tpr, thresholds = sklearn.metrics.roc_curve(y_test.values, y_test_pred)
5 roc_auc = auc(fpr, tpr)
6
7 plt.figure(figsize=(24,4))
8 plt.title('Receiver Operating Characteristic')
9 plt.plot(fpr, tpr, 'b', label = 'AUC = %0.2f' % roc_auc)
10 plt.legend(loc = 'lower right')
11 plt.plot([0, 1], [0, 1], 'r--')
12 plt.xlim([0, 1])
13 plt.ylim([0, 1])
14 plt.grid(which="Both")
15 plt.ylabel('True Positive Rate')
16 plt.xlabel('False Positive Rate')
17 plt.show()
18
19 print("AUC:", roc_auc)

```



AUC: 0.9353843632321274

```

1 model_comparison_table["LSTM"] = {
2     "AUC": 0.9353843632321274, # roc_auc,
3     "Description": "LSTM c/ Keras Tokenizer"
4 }

```

## ▼ 6. Clasificador con BERT

Entrenar un clasificador utilizando BERT.

- a.Utilizar BERT sin fine-tuning.
- b.Utilizar BERT con fine-tuning.

Descargar BERT.



```
1 !pip install transformers > /dev/null 2>&1
```

## Preparación de dataset para BERT

```
1 from sklearn.model_selection import train_test_split
2 TRAIN_TEST_SPLIT = 0.2
3 TRAIN_VAL_SPLIT = 0.25
4
5 X = df.review.values
6 y = df.target.values
7
8 X_train_set, X_test, y_train_set, y_test = train_test_split(X, y, test_size=TRA
9 X_train, X_val, y_train, y_val = train_test_split(X_train_set,y_train_set,test_

1 MAX_SEQ_LENGTH = 64

1 from transformers import BertTokenizer
2
3 PRETRAINED_MODEL_NAME = 'bert-base-uncased'
4 bert_tokenizer = BertTokenizer.from_pretrained(PRETRAINED_MODEL_NAME)

1 def batch_encode(X, tokenizer):
2     return tokenizer.batch_encode_plus(
3         X,
4         max_length=MAX_SEQ_LENGTH,
5         add_special_tokens=True, # [CLS] y [SEP] tokens
6         return_attention_mask=True,
7         return_token_type_ids=False, # no es necesario para clasificador
8         pad_to_max_length=True,
9         return_tensors='tf'
10    )

1 X_train_bert = batch_encode(X_train,bert_tokenizer)
2 X_val_bert = batch_encode(X_val,bert_tokenizer)
3 X_test_bert = batch_encode(X_test,bert_tokenizer)
```

Truncation was not explicitly activated but `max\_length` is provided a speci  
 /usr/local/lib/python3.6/dist-packages/transformers/tokenization\_utils\_base.py  
 FutureWarning,

## ▼ BERT sin fine-tuning

```
1 import tensorflow as tf
2 from tensorflow.keras.metrics import AUC
3 from transformers import TFBertForSequenceClassification
4
5 def create_bert_model_no_ft():
```

```
6 bert_model = TFBertForSequenceClassification.from_pretrained(
7     PRETRAINED_MODEL_NAME, num_labels=2)
8 # Inputs:
9
10 # Layer de input con cadenas de token Ids
11 input_ids = tf.keras.layers.Input(shape=(MAX_SEQ_LENGTH,), dtype=tf.int32, name='input_ids')
12
13 # Attention Mask. Máscara binaria para saber a qué tokens prestar atención y
14 attention_mask = tf.keras.layers.Input((MAX_SEQ_LENGTH,), dtype=tf.int32, name='attention_mask')
15
16 # Conectar salidas anteriores con entradas de BERT:
17 output = bert_model([input_ids, attention_mask])[0]
18
19 # Clasificación binaria
20 output = tf.keras.layers.Dense(1, activation='sigmoid')(output)
21
22 model = tf.keras.models.Model(inputs=[input_ids, attention_mask], outputs=output)
23 model.compile(
24     optimizer=tf.keras.optimizers.Adam(learning_rate=3e-5),
25     loss='binary_crossentropy',
26     metrics=[AUC(name="auc")]
27 )
28 return model
29
30 model = create_bert_model_no_ft()
31 model.summary()
32 tf.keras.utils.plot_model(model)
```

All model checkpoint layers were used when initializing TFBertForSequenceClass

Some layers of TFBertForSequenceClassification were not initialized from the  
 You should probably TRAIN this model on a down-stream task to be able to use  
 The parameters `output\_attentions`, `output\_hidden\_states` and `use\_cache` ca  
 The parameter `return\_dict` cannot be set in graph mode and will always be se  
 Model: "model\_12"

Layer (type)	Output Shape	Param #	Connected to
=====	=====	=====	=====

Entrenamiento BERT sin Fine Tuning.

```
attention_mask = InputLayer, [None, 512], 0
```

Como los modelos con BERT tardan mucho en entrenar en Google Colab, se agregan los siguientes callbacks para detener el entrenamiento y obtener el mejor modelo cuando empieza a ocurrir overfitting.

```
Total params: 100 402 701
```

```
1 from tensorflow.keras.callbacks import EarlyStopping, ModelCheckpoint, ReduceLR0
2 model_fit_callbacks = [
3     EarlyStopping(monitor='val_loss', patience=4, verbose=0, mode='min'),
4     #ModelCheckpoint('bert_nofinetuning.hdf5', save_best_only=True, monitor='val_
5     ReduceLR0nPlateau(monitor='val_loss', factor=0.1, patience=2, verbose=1, min_
6 ]
```

Lamentablemente ModelCheckpoint genera un error cuando se llama a save\_model.  
 TFBertForSequenceClassification no tiene implementado get\_config()?.

I

Para reducir los tiempos de entrenamiento y poder aumentar la cantidad de épocas se trabajará con una partición de train/val reducida.

```
1 REDUX_SPLIT = 0.4
2 TRAIN_REDUX_SIZE = int(X_train.shape[0]*REDUX_SPLIT)
3 VAL_REDUX_SIZE = int(X_val.shape[0]*REDUX_SPLIT)
4 TRAIN_REDUX_SIZE, VAL_REDUX_SIZE

(12000, 4000)

1 train_redux_idx = np.random.choice(np.arange(0,X_train.shape[0]),size=TRAIN_RED
2 X_train_redux, y_train_redux = X_train[train_redux_idx],y_train[train_redux_idx
3
4 val_redux_idx = np.random.choice(np.arange(0,X_val.shape[0]),size=VAL_REDUX_SIZ
5 X_val_redux,y_val_redux = X_val[val_redux_idx],y_val[val_redux_idx]
6
7 X_train_redux_bert = batch_encode(X_train_redux,bert_tokenizer)
8 X_val_redux_bert = batch_encode(X_val_redux,bert_tokenizer)

/usr/local/lib/python3.6/dist-packages/transformers/tokenization_utils_base.py
FutureWarning,
```

```

1 BATCH_SIZE=16
2 NUM_EPOCHS=10
3
4 model = create_bert_model_no_ft()
5 history = model.fit(
6     x=X_train_redux_bert.values(),
7     y=y_train_redux,
8     validation_data=(X_val_redux_bert.values(), y_val_redux),
9     callbacks=model_fit_callbacks,
10    epochs=NUM_EPOCHS,
11    batch_size=BATCH_SIZE
12 )

```

The parameters `output\_attentions`, `output\_hidden\_states` and `use\_cache` cannot be used in graph mode. Epoch 1/10

The parameter `return\_dict` cannot be set in graph mode and will always be set to True. The parameters `output\_attentions`, `output\_hidden\_states` and `use\_cache` cannot be used in graph mode. Epoch 2/10

The parameter `return\_dict` cannot be set in graph mode and will always be set to True. 750/750 [=====] - ETA: 0s - loss: 0.5682 - auc: 0.73

The parameter `return\_dict` cannot be set in graph mode and will always be set to True. 750/750 [=====] - 193s 245ms/step - loss: 0.5681 - auc: 0.73

Epoch 2/10

750/750 [=====] - 193s 258ms/step - loss: 0.2686 - auc: 0.73

Epoch 3/10

750/750 [=====] - 196s 261ms/step - loss: 0.1350 - auc: 0.73

Epoch 00003: ReduceLROnPlateau reducing learning rate to 2.9999999242136257e-05.

Epoch 4/10

750/750 [=====] - 197s 262ms/step - loss: 0.0482 - auc: 0.73

Epoch 5/10

750/750 [=====] - 197s 262ms/step - loss: 0.0181 - auc: 0.73

Epoch 00005: ReduceLROnPlateau reducing learning rate to 2.9999998787388907e-06.

```

1 NUM_EPOCHS_TRAINED = len(history.history["auc"])
2 fig, axes = plt.subplots(1, 2, figsize=(24, 4))
3 plt.suptitle("BERT con Fine Tuning")
4 axes[0].set_title("AUC")
5 axes[0].plot(np.arange(NUM_EPOCHS_TRAINED), history.history['auc'])
6 axes[0].plot(np.arange(NUM_EPOCHS_TRAINED), history.history['val_auc'])
7 axes[0].legend(["Training", "Validation"])
8 axes[0].grid(which="Both")
9 axes[0].set_xlabel("Época")
10 axes[1].set_title("Loss")
11 axes[1].plot(np.arange(NUM_EPOCHS_TRAINED), history.history['loss'])
12 axes[1].plot(np.arange(NUM_EPOCHS_TRAINED), history.history['val_loss'])
13 axes[1].legend(["Training", "Validation"])
14 axes[1].grid(which="Both")
15 axes[1].set_xlabel("Época")
16 plt.show()

```

Se reentrena con una época porque se vé que no hay una mejora luego. Esta vez se usa el

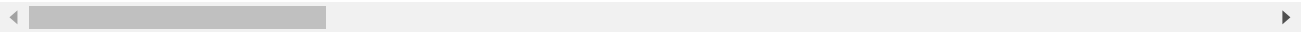
```

1 NUM_EPOCHS = 1
2 BATCH_SIZE=16
3 model = create_bert_model_no_ft()
4 history = model.fit(
5     #x=X_train_redux_bert.values(),
6     #y=y_train_redux,
7     x=X_train_bert.values(),
8     y=y_train,
9     validation_data=(X_val_bert.values(), y_val),
10    #validation_data=(X_val_redux_bert.values(), y_val_redux),
11    #callbacks=model_fit_callbacks,
12    epochs=NUM_EPOCHS,
13    batch_size=BATCH_SIZE
14 )

```

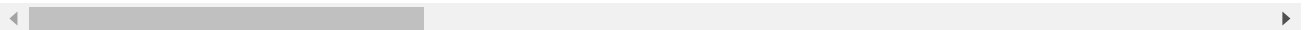
All model checkpoint layers were used when initializing TFBertForSequenceClass

Some layers of TFBertForSequenceClassification were not initialized from the  
 You should probably TRAIN this model on a down-stream task to be able to use  
 The parameters `output\_attentions`, `output\_hidden\_states` and `use\_cache` ca  
 The parameter `return\_dict` cannot be set in graph mode and will always be se  
 The parameters `output\_attentions`, `output\_hidden\_states` and `use\_cache` ca  
 The parameter `return\_dict` cannot be set in graph mode and will always be se  
 The parameters `output\_attentions`, `output\_hidden\_states` and `use\_cache` ca  
 The parameter `return\_dict` cannot be set in graph mode and will always be se  
 1875/1875 [=====] - ETA: 0s - loss: 0.4464 - auc: 0.  
 The parameter `return\_dict` cannot be set in graph mode and will always be se  
 1875/1875 [=====] - 469s 245ms/step - loss: 0.4464 -



```
1 y_test_pred = model.predict([X_test_bert.input_ids,X_test_bert.attention_mask])
```

The parameters `output\_attentions`, `output\_hidden\_states` and `use\_cache` ca  
 The parameter `return\_dict` cannot be set in graph mode and will always be se



```

1 import sklearn
2 from sklearn.metrics import auc
3
4 fpr, tpr, thresholds = sklearn.metrics.roc_curve(y_test, y_test_pred)
5 roc_auc = auc(fpr, tpr)
6
7 plt.figure(figsize=(24,4))
8 plt.title('Receiver Operating Characteristic')
9 plt.plot(fpr, tpr, 'b', label = 'AUC = %0.2f' % roc_auc)
10 plt.legend(loc = 'lower right')
11 plt.plot([0, 1], [0, 1], 'r--')
12 plt.xlim([0, 1])
13 plt.ylim([0, 1])
14 plt.grid(which="Both")
15 plt.ylabel('True Positive Rate')
16 plt.xlabel('False Positive Rate')
17 plt.show()
18

```

```
10
19 print("AUC:", roc_auc)
```



AUC: 0.93906602808635

```
1 model_comparison_table["BERT-NoFT"] = {
2     "AUC": 0.93906602808635, #roc_auc, #0.9257985149861581, # roc_auc,
3     "Description": "BERT sin Fine Tuning"
4 }
```

### ▼ BERT con fine-tuning

Ahora se entrenará la misma arquitectura del caso anterior pero con una capa agrega de Dropout y una capa densa.

```
1 def create_bert_model_ft():
2     bert_model = TFBertForSequenceClassification.from_pretrained(PRETRAINED_MODEL
3                                                                num_labels=2)
4     # Inputs:
5
6     # Layer de input con cadenas de token Ids
7     input_ids = tf.keras.layers.Input(shape=(MAX_SEQ_LENGTH,), dtype=tf.int32, na
8
9     # Attention Mask. Máscara binaria para saber a qué tokens prestar atención y
10    attention_mask = tf.keras.layers.Input((MAX_SEQ_LENGTH,), dtype=tf.int32, nam
11
12    # Conectar salidas anteriores con entradas de BERT:
13    output = bert_model([input_ids, attention_mask])[0]
14
15    # Capas custom
16    output = tf.keras.layers.Dropout(rate=0.15)(output)
17    output = tf.keras.layers.Dense(64)(output)
18
19    # Clasificación binaria
20    output = tf.keras.layers.Dense(1, activation='sigmoid')(output)
21
22    model = tf.keras.models.Model(inputs=[input_ids, attention_mask], outputs=out
23
24    model.compile(
25        optimizer=tf.keras.optimizers.Adam(learning_rate=3e-5),
26        loss='binary_crossentropy',
```

```

27     metrics=[AUC(name="auc")]
28 )
29 return model

1 model = create_bert_model_ft()
2 model.summary()
3 tf.keras.utils.plot_model( model )
4
5 BATCH_SIZE=16
6 NUM_EPOCHS=10
7
8 history = model.fit(
9     x=X_train_redux_bert.values(),
10    y=y_train_redux,
11    validation_data=(X_val_redux_bert.values(), y_val_redux),
12    callbacks=model_fit_callbacks,
13    epochs=NUM_EPOCHS,
14    batch_size=BATCH_SIZE
15 )

```

All model checkpoint layers were used when initializing TFBertForSequenceClass

Some layers of TFBertForSequenceClassification were not initialized from the  
 You should probably TRAIN this model on a down-stream task to be able to use  
 The parameters `output\_attentions`, `output\_hidden\_states` and `use\_cache` ca  
 The parameter `return\_dict` cannot be set in graph mode and will always be se  
 Model: "model\_2"

Layer (type)	Output Shape	Param #	Connected to
input_ids (InputLayer)	[(None, 64)]	0	
attention_mask (InputLayer)	[(None, 64)]	0	
tf_bert_for_sequence_classifica	TFSequenceClassifier	109483778	input_ids[0] attention_ma
dropout_115 (Dropout)	(None, 2)	0	tf_bert_for_
dense_3 (Dense)	(None, 64)	192	dropout_115[
dense_4 (Dense)	(None, 1)	65	dense_3[0][6
Total params: 109,484,035			
Trainable params: 109,484,035			
Non-trainable params: 0			

The parameters `output\_attentions`, `output\_hidden\_states` and `use\_cache` ca  
 The parameter `return\_dict` cannot be set in graph mode and will always be se  
 Epoch 1/10

The parameters `output\_attentions`, `output\_hidden\_states` and `use\_cache` ca  
 The parameter `return\_dict` cannot be set in graph mode and will always be se  
 750/750 [=====] - ETA: 0s - loss: 0.5619 - auc: 0.75  
 The parameter `return\_dict` cannot be set in graph mode and will always be se  
 750/750 [=====] - 185s 235ms/step - loss: 0.5618 - a  
 Epoch 2/10  
 750/750 [=====] - 180s 240ms/step - loss: 0.3153 - a  
 Epoch 3/10

750/750 [=====] - 181s 241ms/step - loss: 0.1737 - a  
Epoch 4/10

750/750 [=====] - 182s 243ms/step - loss: 0.1116 - a

Epoch 00004: ReduceLROnPlateau reducing learning rate to 2.9999999242136257e-  
Epoch 5/10

750/750 [=====] - 183s 244ms/step - loss: 0.0595 - a

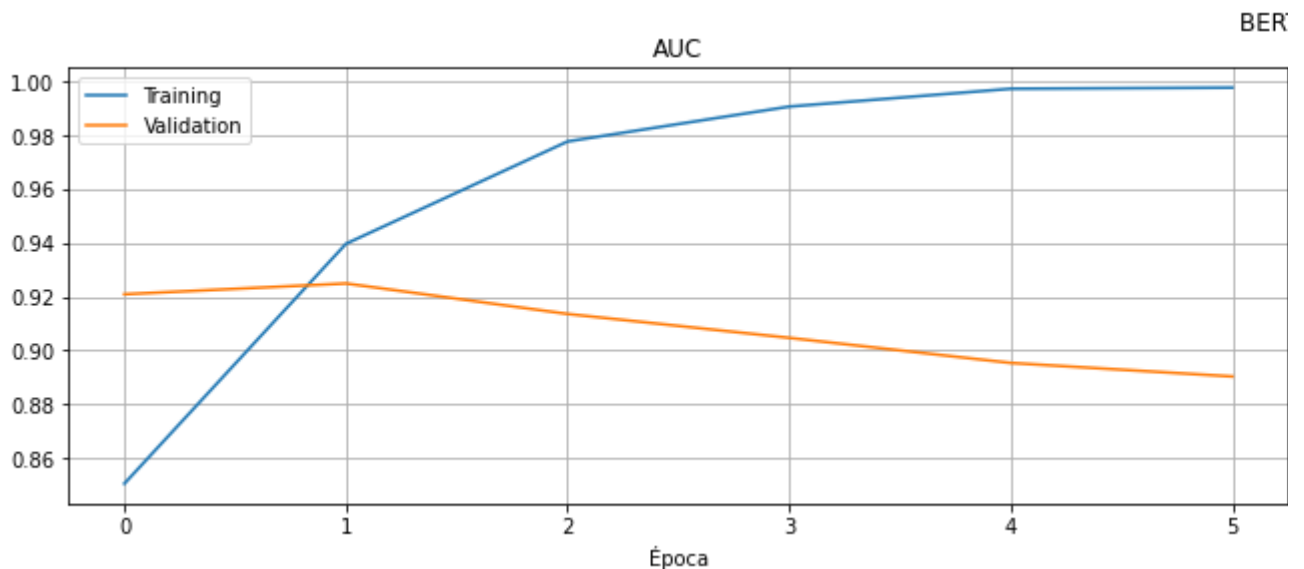
Epoch 6/10

750/750 [=====] - 183s 243ms/step - loss: 0.0443 - a

Epoch 00006: ReduceLROnPlateau reducing learning rate to 2.9999998787388907e-



```
1 NUM_EPOCHS_TRAINED = len(history.history["auc"])
2 fig, axes = plt.subplots(1, 2, figsize=(24, 4))
3 plt.suptitle("BERT con Fine Tuning")
4 axes[0].set_title("AUC")
5 axes[0].plot(np.arange(NUM_EPOCHS_TRAINED), history.history['auc'])
6 axes[0].plot(np.arange(NUM_EPOCHS_TRAINED), history.history['val_auc'])
7 axes[0].legend(["Training", "Validation"])
8 axes[0].grid(which="Both")
9 axes[0].set_xlabel("Época")
10 axes[1].set_title("Loss")
11 axes[1].plot(np.arange(NUM_EPOCHS_TRAINED), history.history['loss'])
12 axes[1].plot(np.arange(NUM_EPOCHS_TRAINED), history.history['val_loss'])
13 axes[1].legend(["Training", "Validation"])
14 axes[1].grid(which="Both")
15 axes[1].set_xlabel("Época")
16 plt.show()
```



```
1 model = create_bert_model_ft()
2 BATCH_SIZE=16
3 NUM_EPOCHS=1
4 history = model.fit(
5     #x=X_train_redux_bert.values(),
6     #y=y_train_redux,
7     x=X_train_bert.values(),
8     ...)
```