



4.3 Avance de proyecto 2: Sistema de Recomendación

Análisis de grandes volúmenes de datos

TC4034 grupo 10 Equipo 26

Luis Arturo Dan Fong | A01650672

Eduardo Rodríguez Ramírez | A01794892

Felipe Enrique Vázquez Ruiz | A01638116

26 de Mayo de 2024

Sistema de Recomendación de Música Personalizado Basado en Listas de Reproducción de Spotify

1.- Descripción del algoritmo de recomendación avanzado elegido:

El algoritmo de recomendación elegido es un modelo de recomendación colaborativa basado en embeddings. En este modelo, se utilizan embeddings para representar tanto a los usuarios como a los artistas en un espacio vectorial de dimensionalidad reducida. Luego, se entrena un modelo de aprendizaje automático para predecir las reproducciones de un usuario para un artista específico. El modelo se optimiza para minimizar el error cuadrático medio entre las reproducciones reales y las predicciones del modelo.

2.- Identificación y justificación de métricas de evaluación utilizadas:

1. Precision@K: Mide la proporción de elementos relevantes recuperados en las primeras K recomendaciones. Es útil para evaluar la precisión de las recomendaciones en una lista corta.

2. Mean Average Precision (MAP): Calcula el promedio de los valores de precisión para cada usuario sobre el conjunto de recomendaciones. Proporciona una medida agregada de la precisión de las recomendaciones.

3. Normalized Discounted Cumulative Gain (NDCG): Evalúa la calidad de las recomendaciones considerando la relevancia y el orden de las mismas en una lista de recomendación. Es una métrica útil para clasificar las recomendaciones en función de su relevancia y posición.

Estas métricas se seleccionaron porque proporcionan una evaluación integral de la calidad de las recomendaciones, considerando tanto la precisión de las mismas como su orden relativo.

3.- Experimentación con al menos 1 algoritmo de recomendación:

El modelo propuesto para este sistema de recomendación consta de una red neuronal basada en embeddings que utiliza información sobre las preferencias de los usuarios y las características de los artistas para hacer predicciones sobre las reproducciones de los usuarios, la arquitectura del modelo consta de lo siguiente:

Entradas: El modelo recibe dos tipos de entradas: el índice del usuario y el índice del artista. Estos índices se utilizan para buscar los vectores de embedding correspondientes en las capas de embedding del modelo.

Capas de embedding: El modelo utiliza capas de embedding para representar a los usuarios y a los artistas en un espacio vectorial de dimensionalidad reducida. Cada usuario y cada artista están asociados con un vector de embedding único que captura sus características y preferencias.

Capas densas: Después de obtener los vectores de embedding para el usuario y el artista, se concatenan y se pasan a través de capas densas de redes neuronales. Estas capas aprenden las relaciones complejas entre los usuarios y los artistas a partir de los vectores de embedding.

Capa de salida: La capa de salida del modelo produce una única salida, que es la predicción de la reproducción del usuario para el artista específico. Esta predicción se calcula como un valor continuo que representa la probabilidad de que el usuario reproduzca el artista.

Resultados de la evaluación del modelo:

Después de entrenar el modelo, se evaluó su rendimiento utilizando las métricas mencionadas anteriormente. Se obtuvieron los siguientes resultados de evaluación:

Precision@10: 0.2463

Mean Average Precision (MAP)@10: 0.0308

Normalized Discounted Cumulative Gain (NDCG)@10: 0.1447

Una precisión del 24.63% en el top 10 significa que aproximadamente una cuarta parte de las recomendaciones en las primeras 10 posiciones son relevantes para los usuarios lo cual indica una calidad razonablemente buena en términos de precisión de las recomendaciones que el sistema genera.

Un MAP de 3.08% indica que, en promedio, el modelo tiene un rendimiento bajo en términos de la precisión del ranking de las recomendaciones, quiere decir que aunque algunas recomendaciones en las posiciones superiores pueden ser acertadas pero no para la mayoría del ranking de cada usuario.

NDCG de 14.47% sugiere que el modelo tiene un rendimiento moderado en términos de la relevancia y el orden de las recomendaciones en las primeras 10 posiciones. Aunque no es excepcionalmente alto, indica que las recomendaciones relevantes tienden a estar más arriba en la lista.

Referencias:

Schedl, M., Zamani, H., Chen, C., Deldjoo, Y., & Elahi, M. (2018). Current challenges and future directions in music recommender systems research. International Journal of Multimedia Information Retrieval, 7(4), 277-300.

Bonnin, G., & Jannach, D. (2015). Automated playlist continuation at scale. In Proceedings of the 9th ACM Conference on Recommender Systems (pp. 115-122).


```

1 import tensorflow as tf
2 from tensorflow.keras.models import Model
3 from tensorflow.keras.layers import Input, Embedding, Flatten, Concatenate, Dense
4 from petastorm.spark import SparkDatasetConverter, make_spark_converter
5 from petastorm.tf_utils import make_petastorm_dataset
6 from tensorflow.keras.callbacks import EarlyStopping

1 # Calculate steps per epoch
2 batch_size = 2048
3 steps_per_epoch = total_rows // batch_size
4 epochs = 10
5 embedding_size = 100

1 # Initialize SparkDatasetConverter
2 converter = make_spark_converter(train_data)
3
4
5 user_input = Input(shape=(1,), name='user_input')
6 user_embedding = Embedding(max_id_users, embedding_size, input_length=1, name='user_embedding')(user_input)
7 user_vec = Flatten(name='user_flatten')(user_embedding)
8
9 artist_input = Input(shape=(1,), name='artist_input')
10 artist_embedding = Embedding(max_id_artists, embedding_size, input_length=1, name='artist_embedding')(artist_input)
11 artist_vec = Flatten(name='artist_flatten')(artist_embedding)
12
13 concat = Concatenate()([user_vec, artist_vec])
14 dense1 = Dense(128, activation='relu')(concat)
15 dense2 = Dense(64, activation='relu')(dense1)
16 output = Dense(1)(dense2)
17
18 model = Model([user_input, artist_input], output)
19 model.compile(optimizer='adam', loss='mean_squared_error')
20
21
22 # Training with Petastorm dataset
23 def transform_row(row):
24     return {"user_input": row.user_id_index, "artist_input": row.artistname_index}, row.normalized_reproduction
25
26 # Create TensorFlow dataset
27 with converter.make_tf_dataset(batch_size=batch_size) as dataset:
28     dataset = dataset.map(lambda x: transform_row(x))
29
30     # Implement early stopping
31     early_stopping = EarlyStopping(monitor='val_loss', patience=3, restore_best_weights=True)
32
33     model.fit(dataset, epochs=epochs, steps_per_epoch=steps_per_epoch, validation_data=dataset, validation_steps=steps_per_epoch, callbacks=[early_stopping])

/usr/local/lib/python3.10/dist-packages/petastorm/fs_utils.py:88: FutureWarning: pyarrow.localfs is deprecated as of 2.0.0, please use pyarrow.fs.LocalFileSystem instead.
  self._filesystem = pyarrow.localfs
WARNING:petastorm:spark:spark_dataset_converter:Converting floating-point columns to float32
WARNING:petastorm:spark:spark_dataset_converter:The median size 472859 B (< 50 MB) of the parquet files is too small. Total size: 906227 B. Increase the median file size by calling df.repartition(n) or df.coalesce(n), which might be required.
/usr/local/lib/python3.10/dist-packages/petastorm/etl/dataset_metadata.py:402: FutureWarning: Passing 'use_legacy_dataset=True' to get the legacy behaviour is deprecated as of pyarrow 11.0.0, and the legacy implementation will be removed.
dataset = pq.ParquetDataset(path_or_paths, filesystem=fs, validate_schema=False, metadata_nthreads=10)
/usr/local/lib/python3.10/dist-packages/petastorm/etl/dataset_metadata.py:402: FutureWarning: Specifying the 'metadata_nthreads' argument is deprecated as of pyarrow 8.0.0, and the argument will be removed in a future version.
dataset = pq.ParquetDataset(path_or_paths, filesystem=fs, validate_schema=False, metadata_nthreads=10)
/usr/local/lib/python3.10/dist-packages/petastorm/etl/dataset_metadata.py:362: FutureWarning: 'ParquetDataset.common_metadata' attribute is deprecated as of pyarrow 5.0.0 and will be removed in a future version.
if not dataset.common_metadata:
/usr/local/lib/python3.10/dist-packages/petastorm/reader.py:420: FutureWarning: Passing 'use_legacy_dataset=True' to get the legacy behaviour is deprecated as of pyarrow 11.0.0, and the legacy implementation will be removed in a future version.
self.dataset = pq.ParquetDataset(dataset_path, filesystem=pyarrow.localfs,
/usr/local/lib/python3.10/dist-packages/petastorm/reader.py:420: FutureWarning: Specifying the 'metadata_nthreads' argument is deprecated as of pyarrow 8.0.0, and the argument will be removed in a future version.
self.dataset = pq.ParquetDataset(dataset_path, filesystem=pyarrow.localfs,
/usr/local/lib/python3.10/dist-packages/petastorm/unischema.py:317: FutureWarning: 'ParquetDataset.pieces' attribute is deprecated as of pyarrow 5.0.0 and will be removed in a future version. Specify 'use_legacy_dataset=False'.
meta = parquet_dataset.pieces[0].get_metadata()
/usr/local/lib/python3.10/dist-packages/petastorm/unischema.py:321: FutureWarning: 'ParquetDataset.partitions' attribute is deprecated as of pyarrow 5.0.0 and will be removed in a future version. Specify 'use_legacy_dataset=False'.
for partition in (parquet_dataset.partitions or []):
/usr/local/lib/python3.10/dist-packages/petastorm/etl/dataset_metadata.py:253: FutureWarning: 'ParquetDataset.metadata' attribute is deprecated as of pyarrow 5.0.0 and will be removed in a future version.
metadata = dataset.metadata
/usr/local/lib/python3.10/dist-packages/petastorm/etl/dataset_metadata.py:254: FutureWarning: 'ParquetDataset.common_metadata' attribute is deprecated as of pyarrow 5.0.0 and will be removed in a future version.
common_metadata = dataset.common_metadata
/usr/local/lib/python3.10/dist-packages/petastorm/etl/dataset_metadata.py:350: FutureWarning: 'ParquetDataset.pieces' attribute is deprecated as of pyarrow 5.0.0 and will be removed in a future version. Specify 'use_legacy_dataset=False'.
futures_list = [tread_pool.submit(_split_piece, piece, dataset, fs.open) for piece in dataset.pieces]
/usr/local/lib/python3.10/dist-packages/petastorm/etl/dataset_metadata.py:350: FutureWarning: 'ParquetDataset.fs' attribute is deprecated as of pyarrow 5.0.0 and will be removed in a future version. Specify 'use_legacy_dataset=False'.
futures_list = [tread_pool.submit(_split_piece, piece, dataset, fs.open) for piece in dataset.pieces]
/usr/local/lib/python3.10/dist-packages/petastorm/etl/dataset_metadata.py:334: FutureWarning: ParquetDatasetPiece is deprecated as of pyarrow 5.0.0 and will be removed in a future version.
return [pq.ParquetDatasetPiece(piece_path, open_file_func=f.open,
/usr/local/lib/python3.10/dist-packages/petastorm/arrow_reader_worker.py:132: FutureWarning: Passing 'use_legacy_dataset=True' to get the legacy behaviour is deprecated as of pyarrow 11.0.0, and the legacy implementation will be removed.
self._dataset = pq.ParquetDataset(
/usr/local/lib/python3.10/dist-packages/petastorm/arrow_reader_worker.py:140: FutureWarning: 'ParquetDataset.fs' attribute is deprecated as of pyarrow 5.0.0 and will be removed in a future version. Specify 'use_legacy_dataset=False'.
parquet_file = Parquetfile(self._dataset, fs.open(piece.path))
/usr/local/lib/python3.10/dist-packages/petastorm/arrow_reader_worker.py:288: FutureWarning: 'ParquetDataset.partitions' attribute is deprecated as of pyarrow 5.0.0 and will be removed in a future version. Specify 'use_legacy_dataset=False'.
partition_names = self._dataset.partitions.partition_names if self._dataset.partitions else set()
/usr/local/lib/python3.10/dist-packages/petastorm/arrow_reader_worker.py:291: FutureWarning: 'ParquetDataset.partitions' attribute is deprecated as of pyarrow 5.0.0 and will be removed in a future version. Specify 'use_legacy_dataset=False'.
table.piece.read(columns=column_names + partition_names, partitions=_self._dataset.partitions)

Epoch 1/10
101/101 [=====] - 87s 847ms/step - loss: 1.8578e-04 - val_loss: 7.0677e-05
Epoch 2/10
101/101 [=====] - 82s 815ms/step - loss: 6.2150e-05 - val_loss: 6.1484e-05
Epoch 3/10
101/101 [=====] - 86s 852ms/step - loss: 5.1474e-05 - val_loss: 4.9467e-05
Epoch 4/10
101/101 [=====] - 82s 817ms/step - loss: 4.5729e-05 - val_loss: 4.7544e-05
Epoch 5/10
101/101 [=====] - 84s 836ms/step - loss: 5.1317e-05 - val_loss: 4.0930e-05
Epoch 6/10
101/101 [=====] - 82s 814ms/step - loss: 3.9329e-05 - val_loss: 3.8764e-05
Epoch 7/10
101/101 [=====] - 84s 832ms/step - loss: 3.6409e-05 - val_loss: 4.5511e-05
Epoch 8/10
101/101 [=====] - 83s 819ms/step - loss: 2.5931e-05 - val_loss: 3.7847e-05
Epoch 9/10
101/101 [=====] - 83s 827ms/step - loss: 4.0968e-05 - val_loss: 3.2006e-05
Epoch 10/10
101/101 [=====] - 83s 819ms/step - loss: 3.1517e-05 - val_loss: 3.0189e-05

```

```

1 model_save_path = working_path + f"/model_{epochs}_{batch_size}"

1 model.save(model_save_path)

1 loaded_model = tf.keras.models.load_model(model_save_path)

```

Para evaluar este modelo se utilizaran 3 metricas:

- Presision@K
- Mean Average Precision (MAP)
- Normalized Discounted Cumulative Gain (NDCG)

```

1 def recommend_top_k_artists(user_id_index, model, data, k=10):
2     artist_ids = data['artistname_index'].unique()
3     user_ids = [user_id_index] * len(artist_ids)
4     user_ids_array = np.array(user_ids)
5     artist_ids_array = np.array(artist_ids)
6     predictions = model.predict([user_ids_array, artist_ids_array], verbose=0)
7     predicted_scores = list(zip(artist_ids, predictions))
8     top_k_artists = sorted(predicted_scores, key=lambda x: x[1], reverse=True)[:k]
9     top_k_artist_ids = [artist_id for artist_id, score in top_k_artists]
10    return top_k_artist_ids
11
12 #Precision@K
13
14 def precision_at_k(recommended_items, relevant_items, k):
15     relevant_set = set(relevant_items)
16     recommended_set = set(recommended_items)
17     return len(recommended_set & relevant_set) / k
18
19 #MAP
20
21 def average_precision(actual, predicted):
22     score = 0.0
23     num_hits = 0.0
24
25     for i, p in enumerate(predicted):
26         if p in actual and p not in predicted[:i]:
27             num_hits += 1.0
28             score += num_hits / (i + 1.0)
29
30     if not actual:
31         return 0.0
32     return score / min(len(actual), len(predicted))
33
34 def mean_average_precision(actual_items, predicted_items):
35     return np.mean([average_precision(actual, predicted) for actual, predicted in zip(actual_items, predicted_items)])
36
37 # NDCG
38
39 def dcg_at_k(r, k):
40     r = np.asarray(r)[:k]
41     if r.size:
42         return np.sum(np.divide(np.power(2, r) - 1, np.log2(np.arange(2, r.size + 2))))
43     return 0.0
44
45 def ndcg_at_k(r, k):
46     idcg = dcg_at_k(sorted(r, reverse=True), k)
47     if not idcg:
48         return 0.0
49     return dcg_at_k(r, k) / idcg
50
51 def normalized_discounted_cumulative_gain(actual, predicted, k):
52     r = [1 if p in actual else 0 for p in predicted]
53     return ndcg_at_k(r, k)
54
55 def evaluate_model(test_data, model, k=10):
56     user_ids = test_data['user_id_index'].unique()
57     avg_map = 0.0
58     avg_ndcg = 0.0
59     precision_scores = []
60
61     for user_id in user_ids:
62         relevant_items = test_data[test_data['user_id_index'] == user_id]['artistname_index'].tolist()
63
64         #Se ignoran usuarios que tengan pocos artistas (Arrancan en frio)
65         if len(relevant_items) < 25:
66             continue
67
68         recommended_items = recommend_top_k_artists(user_id, model, test_data, k)
69
70         #Calcular Precision@K
71         precision = precision_at_k(recommended_items, relevant_items, k)
72         precision_scores.append(precision)
73
74         # Calcular MAP para el usuario actual
75         map_score = average_precision(relevant_items, recommended_items[:k])
76         avg_map += map_score
77
78         # Calcular NDCG para el usuario actual
79         ndcg_score = normalized_discounted_cumulative_gain(relevant_items, recommended_items, k)
80         avg_ndcg += ndcg_score
81
82     # Calcular promedio de MAP y NDCG para todos los usuarios
83     avg_precisionAtK = np.mean(precision_scores)
84     avg_map /= len(user_ids)
85     avg_ndcg /= len(user_ids)
86
87     return avg_precisionAtK, avg_map, avg_ndcg
88
89 # Evaluar el modelo
90 avg_precisionAtK, avg_map, avg_ndcg = evaluate_model(test_data.toPandas(), model, k=10)
91 print(f"Precision@10: {avg_precisionAtK:.4f}")
92 print(f"Mean Average Precision (MAP)@10: {avg_map:.4f}")
93 print(f"Normalized Discounted Cumulative Gain (NDCG)@10: {avg_ndcg:.4f}")

```

→ Precision@10: 0.2463
 Mean Average Precision (MAP)@10: 0.0308
 Normalized Discounted Cumulative Gain (NDCG)@10: 0.1447

✓ Dando Recomendaciones

```

1 test_data.show()

```

user_id_index	artistname_index	reproductions	normalized_reproduction
355	258175	58	0.017040358744394617
8608	258175	42	0.012257100149476832
12719	258175	138	0.040956651718956
550	258175	183	0.05440956651718983
4052	258175	89	0.026507922727047833
11906	258175	98	0.026666875934230195
12659	258175	126	0.03736920777279522
6366	258175	35	0.016164424514260299
3222	258175	17	0.004783258594917788
313	258175	108	0.019388041853512
13393	258175	47	0.013751868469388639
2248	258175	5	0.00119581464879447
8162	258175	56	0.016442451420029897
13626	258175	113	0.03348281016442452
2830	258175	14	0.003886397698377...
14930	258175	12	0.003288490284005979
3102	258175	2	2.389536621823617...
6336	258175	47	0.013751868469388639

```

| 261 | 258175 | 12 | 0.003288490284005979 |
| 13235 | 258175 | 90 | 0.026606875934230195 |
+-----+
only showing top 20 rows

```

```

1 user_id_index=261
2 #13235
3
4 usr_reps = test_data.filter(fn.col("user_id_index") == user_id_index).distinct()\
5     .join(dfs['artistname'],'artistname_index').orderBy('reproductions', ascending=False).select('artistname','reproductions')
6
7 n_artist_usr=usr_reps.count()
8 usr_reps.show(n=n_artist_usr,truncate=False)

```

```

⇒ +-----+
| artistname | reproductions |
+-----+
| Grateful Dead | 42 |
| The Smashing Pumpkins | 40 |
| Iron Maiden | 31 |
| Muse | 16 |
| Blur | 16 |
| Radiohead | 14 |
| Coldplay | 13 |
| U2 | 13 |
| Arctic Monkeys | 12 |
| Foo Fighters | 12 |
| The Smiths | 12 |
| The Rolling Stones | 12 |
| Kings Of Leon | 12 |
| David Bowie | 12 |
| The Killers | 11 |
| The Clash | 10 |
| The Who | 9 |
| Green Day | 9 |
| Nirvana | 8 |
| The Cure | 7 |
| The Strokes | 7 |
| R.E.M. | 6 |
| Weezer | 6 |
| Nine Inch Nails | 5 |
| Guns N' Roses | 5 |
| Beastie Boys | 5 |
| Florence + The Machine | 4 |
| Daft Punk | 4 |
| Linkin Park | 4 |
| Pearl Jam | 3 |
| Depeche Mode | 3 |
| Arcade Fire | 3 |
| The Doors | 3 |
| Led Zeppelin | 3 |
| Queens Of The Stone Age | 3 |
| Vampire Weekend | 2 |
| Mumford & Sons | 2 |
| Massive Attack | 2 |
| The White Stripes | 2 |
| Gorillaz | 2 |
| blink-182 | 2 |
| Beck | 2 |
| Fall Out Boy | 2 |
| Red Hot Chili Peppers | 1 |
| Johnny Cash | 1 |
| Various Artists | 1 |
+-----+

```

```

1 #Generando 10 recomendaciones
2 from pyspark.sql.types import StructType, StructField, StringType, IntegerType
3
4
5 rec_items = recommend_top_k_artists(user_id_index, loaded_model, train_data.toPandas(), 10)
6
7 rec_items = [(int(item),) for item in rec_items]
8
9 rec_artist_names = spark.createDataFrame(rec_items, StructType([StructField("artistname_index", IntegerType(), True)])).join(dfs['artistname'],'artistname_index')
10
11 print(f'Recommended artists for user {user_id_index}')
12 rec_artist_names.show(n=rec_artist_names.count(), truncate=False)

```

```

⇒ Recommended artists for user 261
+-----+
| artistname_index | artistname |
+-----+
| 171812 | Miles Davis |
| 168699 | Metallica |
| 274942 | Vitamin String Quartet |
| 127924 | John Williams |
| 258175 | The Smiths |
| 100812 | Grateful Dead |
| 269792 | U2 |
| 100085 | Gorillaz |
| 194927 | Pearl Jam |
| 136060 | Kendrick Lamar |
+-----+

```

```

1 #Similitudes:
2 usr_reps.join(rec_artist_names,'artistname').select('artistname').show(n=n_artist_usr)

```

```

⇒ +-----+
| artistname |
+-----+
| Gorillaz |
| Grateful Dead |
| Pearl Jam |
| The Smiths |
| U2 |
+-----+

```

```

1 precision = usr_reps.join(rec_artist_names,'artistname').count()/10
2
3 print(f'Precisión de las recomendaciones para el usuario {user_id_index}: {precision}')

```

```

⇒ Precisión de las recomendaciones para el usuario 261: 0.5

```

✓ Ejercicio complementario Arranque en frío

Instalación de paquetes no incluidas en Google colab

```
1 pip install opendatasets
2
3 ➔ Requirement already satisfied: opendatasets in /usr/local/lib/python3.10/dist-packages (0.1.22)
Requirement already satisfied: tqdm in /usr/local/lib/python3.10/dist-packages (from opendatasets) (4.66.4)
Requirement already satisfied: kaggle in /usr/local/lib/python3.10/dist-packages (from opendatasets) (1.6.14)
Requirement already satisfied: click in /usr/local/lib/python3.10/dist-packages (from opendatasets) (8.1.7)
Requirement already satisfied: six>=1.10 in /usr/local/lib/python3.10/dist-packages (from kaggle->opendatasets) (1.16.0)
Requirement already satisfied: certifi>=2023.7.22 in /usr/local/lib/python3.10/dist-packages (from kaggle->opendatasets) (2024.2.2)
Requirement already satisfied: python-dateutil in /usr/local/lib/python3.10/dist-packages (from kaggle->opendatasets) (2.8.2)
Requirement already satisfied: requests in /usr/local/lib/python3.10/dist-packages (from kaggle->opendatasets) (2.31.0)
Requirement already satisfied: python-slugify in /usr/local/lib/python3.10/dist-packages (from kaggle->opendatasets) (8.0.4)
Requirement already satisfied: urllib3 in /usr/local/lib/python3.10/dist-packages (from kaggle->opendatasets) (2.0.7)
Requirement already satisfied: bleach in /usr/local/lib/python3.10/dist-packages (from kaggle->opendatasets) (6.1.0)
Requirement already satisfied: webencodings in /usr/local/lib/python3.10/dist-packages (from bleach->kaggle->opendatasets) (0.5.1)
Requirement already satisfied: text-unidecode>1.3 in /usr/local/lib/python3.10/dist-packages (from python-slugify->kaggle->opendatasets) (1.3)
Requirement already satisfied: charset-normalizer<4,>=2 in /usr/local/lib/python3.10/dist-packages (from requests->kaggle->opendatasets) (3.3.2)
Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.10/dist-packages (from requests->kaggle->opendatasets) (3.7)
```

```
1 pip install fuzzywuzzy
2
3 ➔ Requirement already satisfied: fuzzywuzzy in /usr/local/lib/python3.10/dist-packages (0.18.0)
```

Importar las librerías necesarias

```
1 import opendatasets as od
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 from sklearn.metrics.pairwise import cosine_similarity
```

Descargar el dataset

```
1 od.download(
2     "https://www.kaggle.com/datasets/andrewmvd/spotify-playlists")
3
4 ➔ Skipping, found downloaded files in "./spotify-playlists" (use force=True to force download)
```

Lectura del dataset y eliminación de observaciones no útiles

```
1 #leer el dataset
2 df = pd.read_csv("spotify-playlists//spotify_dataset.csv",on_bad_lines='skip')
3
4 #Renombrar las columnas
5 df.columns = ['user', 'artist', 'song', 'playlist']
6
7 #Eliminar todos aquellos valores NA
8 df.dropna(inplace=True)
9
10 #Eliminar todos los valores duplicados
11 df.drop_duplicates(inplace=True)
12
13 #Eliminar las observaciones que presenten los siguientes valores en la columna de playlist
14 #Estas observaciones se eliminan porque no son playlist públicas
15 df.drop(df[df['playlist']=='Starred'].index,inplace=True)
16 df.drop(df[df['playlist']=='Liked from Radio'].index,inplace=True)
17 df.drop(df[df['playlist']=='Favoritas de la radio'].index,inplace=True)
18 df.drop(df[df['playlist']=='New Playlist'].index,inplace=True)
19 df.drop(df[df['playlist']=='My Shazam Tracks'].index,inplace=True)
20 df.drop(df[df['playlist']=='All Live Files'].index,inplace=True)
21 df.drop(df[df['playlist']=='iPhone'].index,inplace=True)
22
23
24 #obtener un Id para cada usuario
25 user_id = {}
26 cnt = 1
27 for user in df['user']:
28     if user not in user_id:
29         user_id[user] = cnt
30         cnt+=1
31
32 #obtener el conteo de las playlist
33 playlists = {}
34 for playlist in df['playlist']:
35     if playlist in playlists:
36         playlists[playlist] += 1
37     else:
38         playlists[playlist] = 1
```

Generación de columnas útiles

```
1 #Añadir la dataset las columnas de ID y de reproducciones por playlist
2 df['User_Id'] = df['user'].apply(lambda user: user_id[user])
3 df['Repos'] = df['playlist'].apply(lambda playlist: playlists[playlist])
4
5
```

```
1 df.head()
2
3 ➔
4      user        artist          song    playlist  User_Id  Repos
5  0  9cc0cf4d7d7885102480dd99e7a90d6  Elvis Costello  (The Angels Wanna Wear My) Red Shoes  HARD ROCK 2010    1    67
6  1  9cc0cf4d7d7885102480dd99e7a90d6  Elvis Costello & The Attractions  (What's So Funny 'Bout) Peace, Love And Unders...  HARD ROCK 2010    1    67
7  2  9cc0cf4d7d7885102480dd99e7a90d6  Tiffany Page  7 Years Too Late  HARD ROCK 2010    1    67
8  3  9cc0cf4d7d7885102480dd99e7a90d6  Elvis Costello & The Attractions  Accidents Will Happen  HARD ROCK 2010    1    67
9  4  9cc0cf4d7d7885102480dd99e7a90d6  Elvis Costello  Alison  HARD ROCK 2010    1    67
```

Generar dataset con el número de reproducciones que tiene playlist

```

1 #Generar un dataset que solamente contenga el titulo de playlist y su numero de reproducciones
2 playlist_df = pd.DataFrame.from_dict(playlists,orient='index')
3 playlist_df.columns = ['Repros']
4 playlist_df.reset_index(inplace=True)
5 playlist_df.columns = ['Playlist','Repros']
6 playlist_df.sort_values(by='Repros',ascending=False,inplace=True)
7 playlist_df.reset_index(inplace = True,drop=True)

```

1 print(f'Existen {len(playlist_df)} playlists')

2 Existen 157313 playlists

Obtención del top 20 playlist con más reproducciones

```

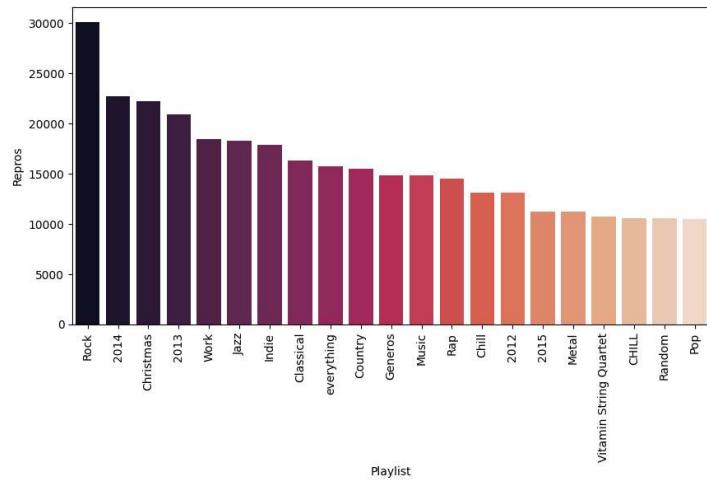
1 plt.figure(figsize=(10,5))
2 sns.barplot(x='Playlist', y='Repros', data=playlist_df.loc[0:20], palette='rocket')
3 plt.xticks(rotation=90)
4 plt.show()

```

1 <ipython-input-10-b488ca9cfa02>:2: FutureWarning:

Passing 'palette' without assigning 'hue' is deprecated and will be removed in v0.14.0. Assign the 'x' variable to 'hue' and set 'legend=False' for the same effect.

```
sns.barplot(x='Playlist', y='Repros', data=playlist_df.loc[0:20], palette='rocket')
```



Obtener solamente las 20 playlist con mas reproducciones en el dataset

```

1 toMatain = playlist_df.loc[:20]
2 del playlist_df

```

1 toMatain

	Playlist	Repros
0	Rock	30107
1	2014	22674
2	Christmas	22236
3	2013	20870
4	Work	18408
5	Jazz	18266
6	Indie	17858
7	Classical	16328
8	everything	15705
9	Country	15503
10	Generos	14854
11	Music	14843
12	Rap	14500
13	Chill	13113
14	2012	13089
15	2015	11185
16	Metal	11178
17	Vitamin String Quartet	10717
18	CHILL	10544
19	Random	10533
20	Pop	10454

Eliminar todas aquellas observaciones que no posean alguna playlist de las mostradas anteriormente.

```
1 df = df[df['playlist'].isin(toMatain['Playlist'])]
```

Realizar el mismo proceso que se lleva a cabo con las playlist, pero ahora con los artistas.

```

1 #Obtener las reproducciones por artista
2 ReprosArtistas = {}
3 for artist in df['artist']:
4     if artist in ReprosArtistas:
5         ReprosArtistas[artist]+=1
6     else:
7         ReprosArtistas[artist]=1

```

```
1 artists_df = pd.DataFrame.from_dict(ReprosArtistas,orient='index')
```



```

1 #Encontrar la matriz de similaridad
2 cosine_sim = cosine_similarity(dummy, dummy)
3 print(f"Las dimensiones de similaridad coseno de las características de nuestra matriz de similitud son: {cosine_sim.shape}")
↳ Las dimensiones de similaridad coseno de las características de nuestra matriz de similitud son: (28588, 28588)

1 from fuzzywuzzy import process
2
3 song_idx = dict(zip(df['song'], list(df.index)))
4 def encuentra_cancion(title):
5     all_titles = df['song'].tolist()
6     closest_match = process.extractOne(title,all_titles)
7     return closest_match[0]
8
9 def obtener_recomendaciones_basadas_contenido(title_string, numero_recomendaciones=10):
10    title = encuentra_cancion(title_string)
11    idx = song_idx.get(title)
12    sim_scores = list(enumerate(cosine_sim[idx]))
13    sim_scores = sorted(sim_scores, key=lambda x: x[1], reverse=True)
14    sim_scores = sim_scores[1:(numero_recomendaciones+1)]
15    similar_movies = [i[0] for i in sim_scores]
16    print(f"Las películas recomendadas con base en la película {title} son:")
17    print(df['song'].iloc[similar_movies])

/usr/local/lib/python3.10/dist-packages/fuzzywuzzy/fuzz.py:11: UserWarning: Using slow pure-python SequenceMatcher. Install python-Levenshtein to remove this warning
warnings.warn("Using slow pure-python SequenceMatcher. Install python-Levenshtein to remove this warning")

1 obtener_recomendaciones_basadas_contenido('Basin Street Blues',10)
↳ Las películas recomendadas con base en la película Basin Street Blues son:
1          Fall in Love Too Easily
2          My Funny Valentine
576          All Blues
582          Blue in Green
587          Flamenco Sketches
588          Freddie Freeloader
618          So What
814          Enigma
815          Feijo
820  It Never Entered My Mind (Higher & Higher)
Name: song, dtype: object

```