

# Proyecto Final de Machine Learning

## Sistema de Recomendación

Addison Amin Reyes Cedano, 2021-2026

Mi [proyecto](#) es el de un sistema de recomendación el cual dependiendo de las valoraciones del usuario a los anime que ya ha visto crea n recomendaciones al usuario de animes para ver, lo primero que hice fue buscar un dataset que cumpliera los requisitos del proyecto y que satisficiera mis necesidades. Este es el dataset que encontré:

<https://www.kaggle.com/datasets/CooperUnion/anime-recommendations-database?select=rating.csv>

Lo primero que realice fue crear el entorno virtual e importar todas las librerías que utilizare en el proyecto, aquí todas las librerías que utilice:

```
1 from sklearn.model_selection import train_test_split
2 from sklearn.preprocessing import LabelEncoder
3 from dataprep.eda import create_report
4 import matplotlib.pyplot as plt
5 import tensorflow as tf
6 import pandas as pd
7 import numpy as np
```

✓ 0.0s

Una vez importadas las librerías leo los dos archivos '.csv' que son "rating.csv" y "anime.csv"

```
1 ratings = pd.read_csv("rating.csv")
2 items = pd.read_csv("anime.csv")
```

✓ 1.4s

Ya con los dos archivos leídos inicia la preparación y preprocesamiento de la data, en donde junto los dos DataFrame en uno solo en donde tengo los usuarios, los id de los animes, la puntuación que le han dado al anime que han visto y el nombre del anime aparte del id, con este DataFrame completo hago ingeniería de características eliminando datos innecesarios, ajustando algunas columnas, revisando si hay filas vacías, etc..

```
1 ratings.head(6)
```

✓ 0.0s

	user_id	anime_id	rating
0	1	20	-1
1	1	24	-1
2	1	79	-1
3	1	226	-1
4	1	241	-1
5	1	355	-1

```
1 items.head(6)
```

✓ 0.0s

	anime_id	name	genre	type	episodes	rating	members
0	32281	Kimi no Na wa.	Drama, Romance, School, Supernatural	Movie	1	9.37	200630
1	5114	Fullmetal Alchemist: Brotherhood	Action, Adventure, Drama, Fantasy, Magic, Mili...	TV	64	9.26	793665
2	28977	Gintama°	Action, Comedy, Historical, Parody, Samurai, S...	TV	51	9.25	114262
3	9253	Steins;Gate	Sci-Fi, Thriller	TV	24	9.17	673572
4	9969	Gintama&#039;	Action, Comedy, Historical, Parody, Samurai, S...	TV	51	9.16	151266
5	32935	Haikyuu!!: Karasuno Koukou VS Shiratorizawa Ga...	Comedy, Drama, School, Shounen, Sports	TV	10	9.15	93351

```
1 items_titles = items[["anime_id", "name"]]  
2 items_titles.head(6)
```

✓ 0.0s

	anime_id	name
0	32281	Kimi no Na wa.
1	5114	Fullmetal Alchemist: Brotherhood
2	28977	Gintama°
3	9253	Steins;Gate
4	9969	Gintama&#039;
5	32935	Haikyuu!!: Karasuno Koukou VS Shiratorizawa Ga...

```
1 items_ratings = ratings.merge(items_titles, on="anime_id")
2 items_ratings.head(6)
```

✓ 1.0s

	user_id	anime_id	rating	name
0	1	20	-1	Naruto
1	3	20	8	Naruto
2	5	20	6	Naruto
3	6	20	-1	Naruto
4	10	20	-1	Naruto
5	21	20	8	Naruto

```
1 items_ratings.rename(columns={'anime_id': 'item_id'}, inplace = True)
```

✓ 0.0s

```
1 items_titles[items_titles["name"] == "Naruto"]
```

✓ 0.0s

	anime_id	name
841	20	Naruto

```
1 items_ratings.info()
```

✓ 0.0s

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 7813727 entries, 0 to 7813726
Data columns (total 4 columns):
#   Column  Dtype
---  -
0   user_id  int64
1   item_id  int64
2   rating   int64
3   name     object
dtypes: int64(3), object(1)
memory usage: 298.1+ MB
```

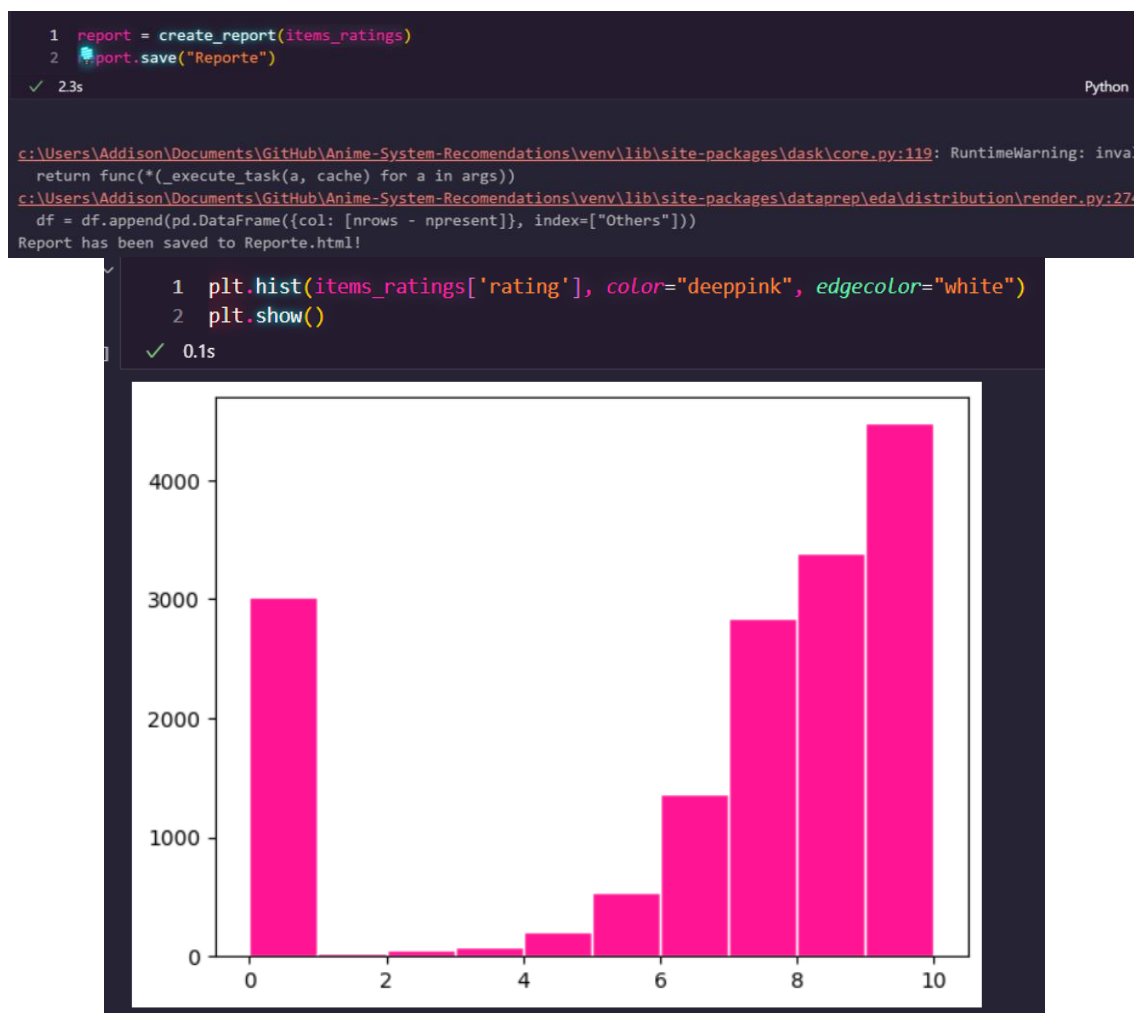
```
1 items_ratings.isna().any()
✓ 0.9s

user_id    False
item_id    False
rating     False
name       False
dtype: bool

1 items_ratings.loc[items_ratings['rating'] == -1, 'rating'] = 0
✓ 0.1s

1 items_ratings = items_ratings.sample(16000)
✓ 0.5s
```

Hice el análisis descriptivo de los datos con dataprep, aquí el [link](#) y una captura del histograma a los 'ratings':



Después de analizar la data paso la columna de 'user\_id' y 'item\_id' por un LabelEncoder para luego hacer un train test Split con Scikit learn implemento un modelo de recomendación utilizando el método de factorización de matrices con redes neuronales para predecir la calificación que un usuario podrá asignar a un anime que aún no ha visto.

```
1 encoder = LabelEncoder()
✓ 0.0s

1 encoder.fit(items_ratings['user_id'])
2 items_ratings.loc[:, 'user_id'] = sorted(encoder.transform(items_ratings['user_id']))
3
4 encoder.fit(items_ratings['item_id'])
5 items_ratings.loc[:, 'item_id'] = encoder.transform(items_ratings['item_id'])
✓ 0.1s

1 df_train, df_val = train_test_split(items_ratings, test_size=0.2, random_state=42)
```

Primero se crean dos matrices de ceros, train\_matrix y test\_matrix, con un número de filas igual al número de usuarios y un número de columnas igual al número de películas únicas en el conjunto de datos. La matriz train\_matrix se inicializa con las calificaciones de entrenamiento del conjunto de datos de entrenamiento, y la matriz test\_matrix se inicializa con las calificaciones de prueba del conjunto de datos de prueba.

```
1 n_users = len(items_ratings.user_id.unique())
2 n_items = len(items_ratings.item_id.unique())
3
4 train_matrix = np.zeros((n_users, n_items))
5
6 for row in df_train.itertuples():
7     train_matrix[row[1]-1, row[2]-1] = row[3]
✓ 0.1s
```

Luego, se definen dos capas de entrada para el modelo, input\_users e input\_items, que representan los identificadores de los usuarios y de las películas, respectivamente. Cada una de estas capas se conecta a una capa de embedding, embedding\_users y embedding\_items, que transforma los identificadores de los usuarios y las películas en vectores densos de baja dimensión. Estos vectores se aplanan en capas flatten\_users y flatten\_items, respectivamente.

```

1 input_users = tf.keras.layers.Input(shape=[1])
2 embedding_users = tf.keras.layers.Embedding(n_users, 100)(input_users)
3 flatten_users = tf.keras.layers.Flatten()(embedding_users)
4
5 input_items = tf.keras.layers.Input(shape=[1])
6 embedding_items = tf.keras.layers.Embedding(n_items, 100)(input_items)
7 flatten_items = tf.keras.layers.Flatten()(embedding_items)

```

✓ 0.1s

Después, se concatenan los vectores de usuario y película en una única capa de entrada concat, que se alimenta a dos capas densas con activación relu y una tasa de dropout del 20% en cada capa, dense1 y dense2. Esto ayuda a evitar el sobreajuste y mejorar la generalización del modelo.

```

1 concat = tf.keras.layers.Concatenate()([flatten_users, flatten_items])
2
3 dense1 = tf.keras.layers.Dense(128, activation='relu')(concat)
4 dropout1 = tf.keras.layers.Dropout(0.2)(dense1)
5
6 dense2 = tf.keras.layers.Dense(64, activation='relu')(dropout1)
7 dropout2 = tf.keras.layers.Dropout(0.2)(dense2)

```

✓ 0.0s

Finalmente, se utiliza una capa densa sin activación en la salida, output, que predice la calificación de la película para el usuario. El modelo se compila utilizando el error cuadrático medio (MSE) como función de pérdida y el optimizador Adam.

```

1 output = tf.keras.layers.Dense(1)(dropout2)
2
3 model = tf.keras.Model(inputs=[input_users, input_items], outputs=output)
4 model.compile(loss='mse', optimizer='adam')

```

✓ 0.0s

Una vez tenido el modelo entrenado utilizo la función .fit del modelo para ajustar los pesos dl modelo para minimizar la pérdida de un anime a otro, por ultimo evaluó el modelo con los datos que me había dado el train\_split\_train.

```

1 model.fit([df_train.user_id, df_train.item_id], df_train.rating, epochs=10, verbose=1, validation_split=0.2, batch_size=70)

```

✓ 28.6s

```

Epoch 1/10
147/147 [=====] - 3s 20ms/step - loss: 0.9447 - val_loss: 16.2858
Epoch 2/10
147/147 [=====] - 3s 20ms/step - loss: 0.8780 - val_loss: 16.0329
Epoch 3/10
147/147 [=====] - 3s 19ms/step - loss: 0.8570 - val_loss: 16.2854
Epoch 4/10
147/147 [=====] - 3s 19ms/step - loss: 0.8409 - val_loss: 16.2366

```



```

1 test_loss = model.evaluate([df_val.user_id, df_val.item_id], df_val.rating)
2 print('Test loss:', test_loss)
✓ 0.2s

100/100 [=====] - 0s 1ms/step - loss: 15.7229
Test loss: 15.722870826721191

```

Por ultimo en el proyecto creo una función la cual es la que se encarga de devolver las recomendaciones al usuario eliminando los animes visto por el usuario y luego prediciendo de los animes no vistos cuales tienen más probabilidad de ser una mejor recomendación para luego devolver las n recomendaciones en un DataFrame.

```

1 def recommend_items(user_id, model, n=2):
2     #Se obtiene items que ya ha visto el usuario
3     seen_movies = items_ratings[items_ratings['user_id'] == user_id]['item_id'].tolist()
4
5     #Se obtiene items que el usuario aún no ha visto
6     item_ids = list(set(items_ratings['item_id']) - set(seen_movies))
7
8     #Se predicen las puntuaciones para las items no vistas
9     ratings_pred = model.predict([np.array([user_id]*len(item_ids)), np.array(item_ids)])
10    ratings_pred = np.array([rating[0] for rating in ratings_pred])
11
12    #Obtener los mejores n items recomendadas
13    top_item_ids = (-ratings_pred).argsort()[0:n]
14    top_item_ids = [item_ids[i] for i in top_item_ids]
15    top_item_ids += [0]*(n - len(top_item_ids))
16    top_ratings_pred = ratings_pred[top_item_ids]
17
18    #Se crea un DataFrame con los mejores n items
19    items_pred = pd.DataFrame(columns=['anime_id', 'name', 'genre', 'type', 'episodes', 'rating', 'members'])
20
21    i = 0
22    for ids in top_item_ids:
23        items_pred.loc[i] = items.iloc[ids]
24        i += 1
25
26    items_pred['predicted_rating'] = top_ratings_pred
27
28    return items_pred
✓ 0.0s

```

## Recomendaciones

```

1 user_id = 16
2 recommendations = recommend_items(user_id, model, n=3)
3 recommendations
✓ 0.4s

104/104 [=====] - 0s 1ms/step

```

	anime_id	name	genre	type	episodes	rating	members	predicted_rating
0	6324	Omamori Himari	Action, Comedy, Demons, Ecchi, Harem, Romance,...	TV	12	7.19	119572	-0.018156
1	4186	Chrome Shelled Regios	Action, Adventure, Fantasy, School, Sci-Fi	TV	24	7.57	103615	6.951766
2	1434	Lupin III: Harimao no Zaihou wo Oe!!	Action, Adventure, Comedy, Shounen	Special	1	7.16	3057	0.009716

```

1 user_id = 69
2 recommendations = recommend_items(user_id, model, n=3)
3 recommendations
✓ 0.3s

104/104 [=====] - 0s 1ms/step

```

	anime_id	name	genre	type	episodes	rating	members	predicted_rating
0	6024	Chi&#039;s Sweet Home: Atarashii Ouchi	Comedy, Kids, Slice of Life	TV	104	7.9	16670	9.494140
1	6324	Omamori Himari	Action, Comedy, Demons, Ecchi, Harem, Romance,...	TV	12	7.19	119572	6.650553
2	3483	Maison Ikkoku: Utsuriyuku Kisetsu no Naka de	Comedy, Drama, Romance	OVA	1	7.15	1458	3.755569