

MINI PROYECTO

Equipo 12 – Deep Learning

Descripción breve

- El presente informe aborda la exploración de los datos para su entendimiento dentro del contexto organizacional, la preparación de los datos para poder utilizarlos como entrada para modelos predictivos, el análisis preliminar de selección de modelos relevantes para responder a la pregunta, el desarrollo y calibración de modelos y la visualización de resultados .

INTEGRANTES

-Christian Beraún Chamorro
- D'sharlie Sánchez Rozo
- William Alexander Morales Valera

El presente miniproyecto se basa en las características musicales de una canción y si estas influyen o son determinantes respecto del año de publicación o lanzamiento. En resumen, se busca responder a la pregunta ¿Existe una relación entre las características musicales de una canción y el año en que fue publicada/lanzada? Una respuesta positiva a esta pregunta revelaría una profunda comprensión de la naturaleza de una composición musical y, lo que es más importante, esta comprensión se demostraría matemáticamente.

En general, se deben reportar los resultados del modelamiento predictivo siguiendo los pasos que se muestran a continuación:

- ## DATASETS PROPORCIONADOS

- Disponibles en: <https://www.kaggle.com/competitions/music-year-prediction/data>

Para empezar, necesitamos mapear la descripción, cantidad de variables y calidad de los datos de ambos datasets, los cuales se realizó a continuación:

Tabla 1. Análisis Descriptivo de los datos

[illegible]

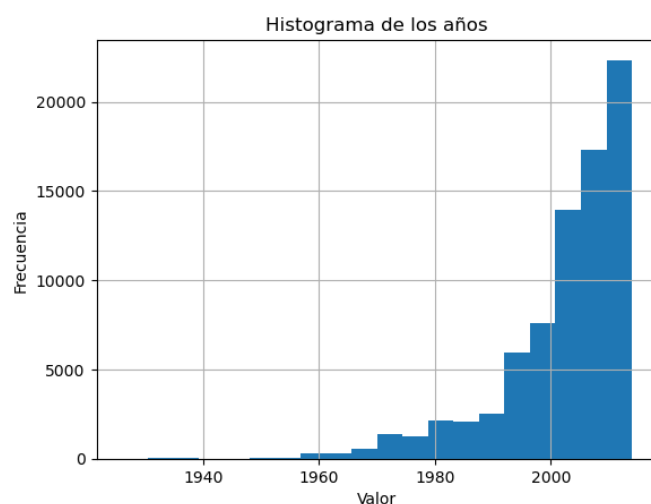
Se realiza un primer análisis descriptivos de las variables revisando la medida, mínimos, máximos y desviación estándar, a simple vista se puede evidenciar que la columna V3 esta presente en todos los registros pero su valor es cero, por otra lado en las primeras 12 variables se pueden evidenciar valores mas acordes según su desviación estándar, siendo que son 91 variables, se revisó en el notebook de una forma más amplia en histogramas, para efectos de observar el comportamiento y distribución de las variables.

Gráfico n.º 1 – Histogramas



Luego se realizó un análisis de la distribución de la variable a predecir:

Gráfico n.º 2 – Distribución de la variable a predecir



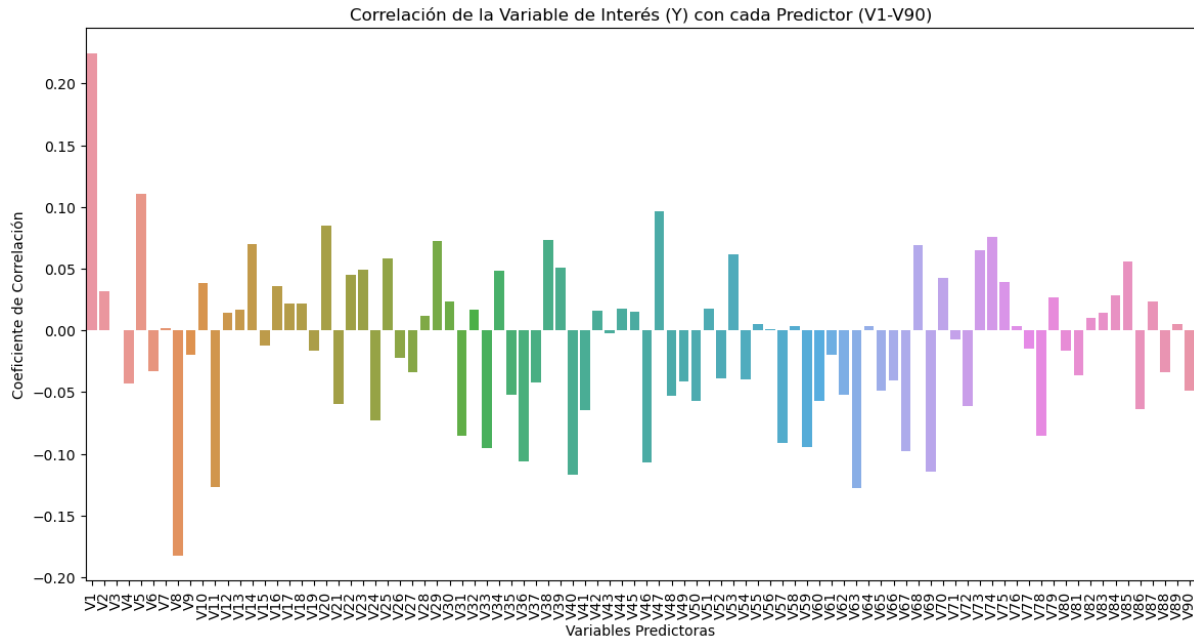
Como podemos observar, existe una distribución mayor para los años más cercanos a la fecha; ello podría deberse a que los datos se encuentran disponibles en diversas plataformas streaming de música siendo que estas surgieron en los últimos 20 años, tal como lo precisa en la página web de los datasets:

- SecondHandSongs dataset -> cover songs
- musiXmatch dataset -> lyrics
- Last.fm dataset -> song-level tags and similarity
- Taste Profile subset -> user data
- thisismyjam-to-MSD mapping -> more user data

- tagtraum genre annotations -> genre labels
- Top MAGD dataset -> more genre labels

Ahora bien, tengamos en cuenta este dato importante "los 12 primeros (V1 A V12) corresponden al timbre promedio y los 78 siguientes a la covarianza (V13 a V90)". Analizaremos la correlación que existe entre las variables predictoras con la variable de interés:

Gráfico n.º 3 – Matriz de Correlación

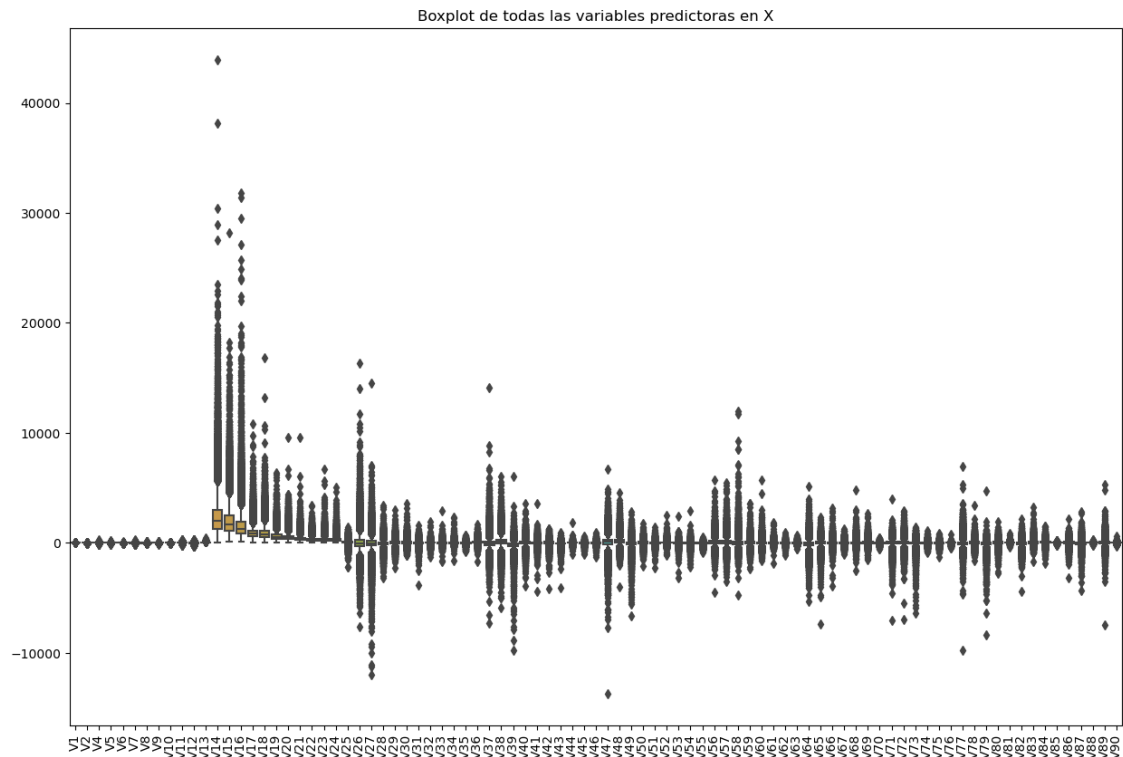


Realizando el análisis de correlación de las variables predictoras vs la variable objetivo, podemos evidenciar que hay una correlación de -20 a 20, lo que indica que hay muy poca correlación entre la variable objetivo y las variables predictoras.

3. PARTE II: PREPARACIÓN DE LOS DATOS PARA MODELOS PREDICTIVOS

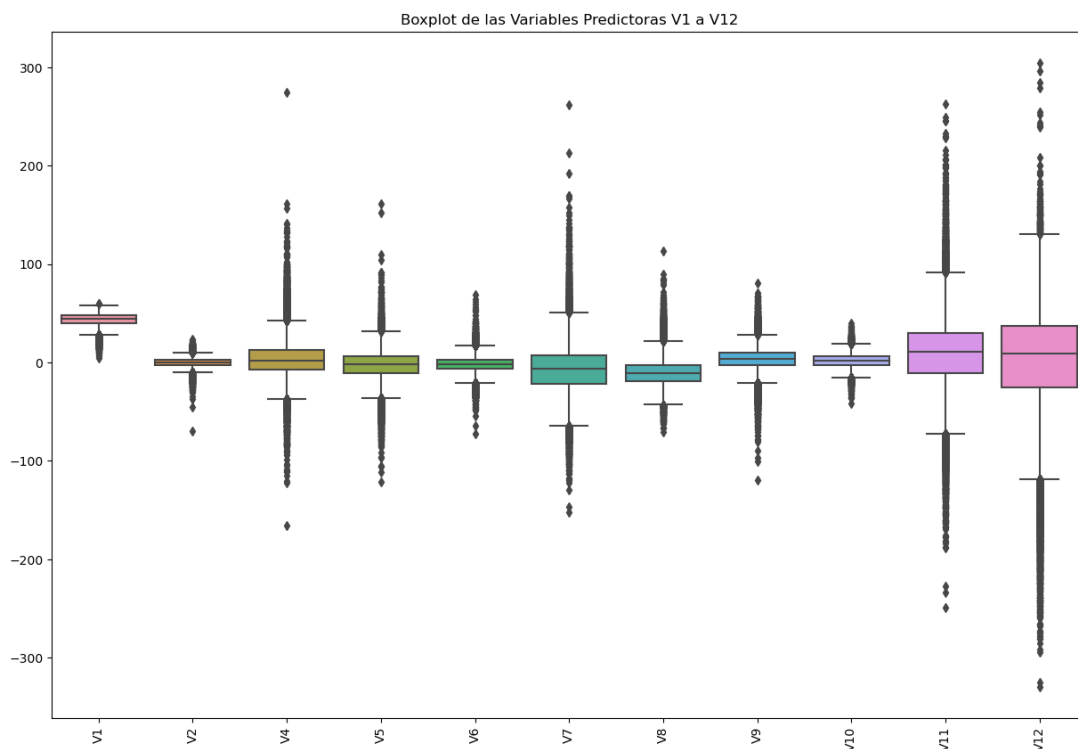
Como pudimos observar, las variables predictoras van de V1 a V90, pero se debe eliminar V3 y ID, siendo que no agregan valor al modelo, por ser valores vacíos o el índice de los registros. Ahora procedemos a analizar la dispersión de los datos

Gráfico n.º 4 – Determinación de outliers



Ahora bien, si observamos las variables V14 en adelante, podemos ver que hay mucha variabilidad. Mientras que de V1 a V12, no hay mucha.

Gráfico n.º 5 – Boxplot de las variables V1 a V12



Considerando todo lo anterior, se vio por conveniente crear variables que capturen los datos de V1 a V12 y de V13 a V90, usando métricas como suma, promedio, mediana, producto, entre otras. Conforme se detalla a continuación:

```

# Función para crear características
train_data = train_data.drop(columns=['ID','V3'])
test_data = test_data.drop(columns=['ID'])

def create_features(data):
    # Adjusted range to exclude 'V3'
    timbre_cols = [f"V{i}" for i in range(1, 13) if i != 3]

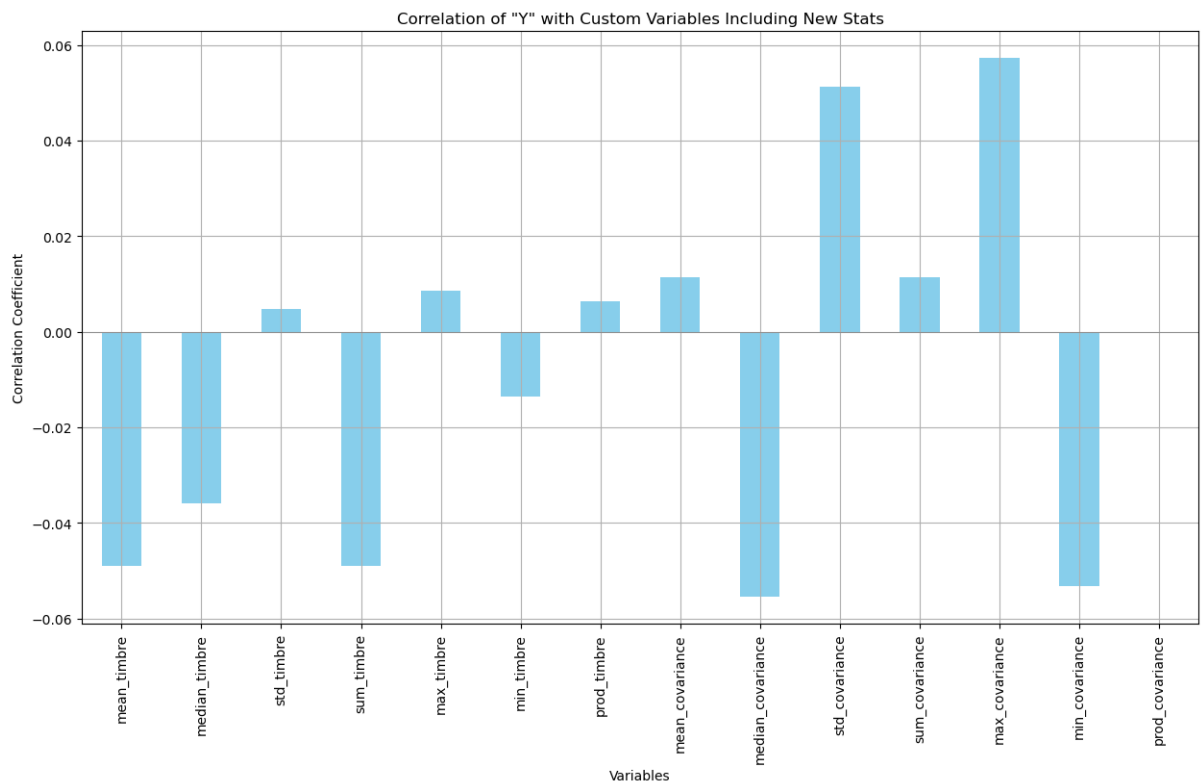
    # Check if all required timbre columns except V3 exist
    if all(col in data.columns for col in timbre_cols):
        data['mean_timbre'] = data[timbre_cols].mean(axis=1)
        data['median_timbre'] = data[timbre_cols].median(axis=1)
        data['std_timbre'] = data[timbre_cols].std(axis=1)
        data['sum_timbre'] = data[timbre_cols].sum(axis=1)
        data['max_timbre'] = data[timbre_cols].max(axis=1)
        data['min_timbre'] = data[timbre_cols].min(axis=1)
        data['prod_timbre'] = data[timbre_cols].prod(axis=1)
    else:
        print("Algunas columnas necesarias (V1, V2, V4 a V12) no existen en el DataFrame")
    # Check if all covariance columns (V13 to V90) exist
    if all(f"V{i}" in data.columns for i in range(13, 91)):
        covariance_cols = [f"V{i}" for i in range(13, 91)]
        data['mean_covariance'] = data[covariance_cols].mean(axis=1)
        data['median_covariance'] = data[covariance_cols].median(axis=1)
        data['std_covariance'] = data[covariance_cols].std(axis=1)
        data['sum_covariance'] = data[covariance_cols].sum(axis=1)
        data['max_covariance'] = data[covariance_cols].max(axis=1)
        data['min_covariance'] = data[covariance_cols].min(axis=1)
        # Using Logarithmic addition to prevent overflow in product
        data['prod_covariance'] = np.exp(data[covariance_cols].apply(np.log).sum(axis=1))
    else:
        print("Algunas columnas de V13 a V90 no existen en el DataFrame")

    return data
# Aplicar la creación de características
train_data = create_features(train_data)
test_data = create_features(test_data)

```

Como se puede evidenciar en el fragmento de código anterior, se valida la existencia de columnas en el set de datos con el objetivo de manejar posibles errores al momento de realizar la ejecución de proceso.

Gráfico n.º 6 – Correlación con “Y” de las nuevas variables



En la gráfico anterior podemos observar que fue buena decisión agregar estas características, especialmente aquellas que muestran una correlación más alta de 0.02. Estas no solo capturan los datos de los dos grupos mencionados, sino que también están más estrechamente relacionadas con la variable objetivo

4. PARTE III: ANÁLISIS PRELIMINAR DE SELECCIÓN DE MODELOS RELEVANTES

*Modelo 1: Una sola neurona *

```
In [21]: # Modelo 1: Simplificado con una sola neurona
model1 = Sequential([
    Dense(1, input_dim=X_train_scaled.shape[1], activation='linear') # Una sola salida para regresión
])

# Compilación del modelo con el optimizador Adam
optimizer = Adam(learning_rate=0.2)
model1.compile(optimizer=optimizer, loss='mean_squared_error')

# Callback de EarlyStopping para evitar el sobreajuste
early_stopping = EarlyStopping(monitor='val_loss', patience=5, verbose=1, restore_best_weights=True)

# Entrenamiento del modelo con EarlyStopping
model1.fit(X_train_scaled, y_train, epochs=10, batch_size=32, validation_data=(X_val_scaled, y_val), callbacks=[early_stopping])

# Evaluación del modelo
y_pred = model1.predict(X_val_scaled)
if np.isnan(y_pred).any():
    raise ValueError("Las predicciones contienen NaN")

# Cálculo del RMSE para evaluar el rendimiento del modelo
rmse = sqrt(mean_squared_error(y_val, y_pred))
print(f'RMSE: {rmse}')
```

Epoch 1/10
1945/1945 ————— 2s 640us/step - loss: 3646274.7500 - val_loss: 2653277.7500
Epoch 2/10
1945/1945 ————— 1s 602us/step - loss: 2378656.0000 - val_loss: 1631743.1250
Epoch 3/10
1945/1945 ————— 1s 597us/step - loss: 1429066.5000 - val_loss: 888608.5000
Epoch 4/10
1945/1945 ————— 1s 621us/step - loss: 748726.3125 - val_loss: 392420.5625
Epoch 5/10
1945/1945 ————— 1s 616us/step - loss: 309997.9375 - val_loss: 116953.0547
Epoch 6/10
1945/1945 ————— 1s 602us/step - loss: 81577.9453 - val_loss: 14191.3877
Epoch 7/10
1945/1945 ————— 1s 611us/step - loss: 7712.5869 - val_loss: 218.8185
Epoch 8/10
1945/1945 ————— 1s 613us/step - loss: 135.9781 - val_loss: 102.4660
Epoch 9/10
1945/1945 ————— 1s 599us/step - loss: 95.4779 - val_loss: 97.0490
Epoch 10/10
1945/1945 ————— 1s 605us/step - loss: 102.8793 - val_loss: 108.8134
Restoring model weights from the end of the best epoch: 9.
487/487 ————— 0s 491us/step
RMSE: 9.851345353532812

*Modelo 2: Múltiples Capas *

*Modelo 2: Múltiples Capas *

```
model2 = Sequential([
    Dense(200, input_dim=X_train_scaled.shape[1], activation='relu'), # Primera capa oculta con 100 neuronas
    Dense(50, activation='relu'), # Segunda capa oculta con 20 neuronas
    Dense(1, activation='linear') # Capa de salida para regresión
])

# Compilación del modelo con el optimizador Adam
optimizer = Adam(learning_rate=0.001)
model2.compile(optimizer=optimizer, loss='mean_squared_error')

# Callback de EarlyStopping para evitar el sobreajuste
early_stopping = EarlyStopping(monitor='val_loss', patience=5, verbose=1, restore_best_weights=True)

# Entrenamiento del modelo con EarlyStopping
model2.fit(X_train_scaled, y_train, epochs=10, batch_size=32, validation_data=(X_val_scaled, y_val), callbacks=[early_stopping])

# Evaluación del modelo
y_pred = model2.predict(X_val_scaled)
if np.isnan(y_pred).any():
    raise ValueError("Las predicciones contienen NaN")

# Cálculo del RMSE para evaluar el rendimiento del modelo
rmse = sqrt(mean_squared_error(y_val, y_pred))
print(f'RMSE: {rmse}')
```

Epoch 1/10

C:\Users\CHRISTIAN\anaconda3\lib\site-packages\keras\src\layers\core\dense.py:86: UserWarning: Do not pass an `input_shape` in 'input_dim' argument to a layer. When using Sequential models, prefer using an `Input(shape)` object as the first layer in the model instead.

super().__init__(activity_regularizer=activity_regularizer, **kwargs)

```
2066/2066 — 2s 754us/step - loss: 1203407.5000 - val_loss: 51299.4648
Epoch 2/10
2066/2066 — 1s 716us/step - loss: 37058.0742 - val_loss: 11845.1768
Epoch 3/10
2066/2066 — 2s 719us/step - loss: 7206.3481 - val_loss: 1883.6832
Epoch 4/10
2066/2066 — 2s 722us/step - loss: 1215.3699 - val_loss: 666.9042
Epoch 5/10
2066/2066 — 2s 737us/step - loss: 395.8275 - val_loss: 492.5071
Epoch 6/10
2066/2066 — 2s 725us/step - loss: 283.0558 - val_loss: 659.5401
Epoch 7/10
2066/2066 — 1s 713us/step - loss: 341.4096 - val_loss: 618.0717
Epoch 8/10
2066/2066 — 2s 727us/step - loss: 418.1860 - val_loss: 289.3042
Epoch 9/10
2066/2066 — 2s 722us/step - loss: 258.7532 - val_loss: 192.1150
Epoch 10/10
2066/2066 — 2s 721us/step - loss: 308.6641 - val_loss: 233.8360
Restoring model weights from the end of the best epoch: 9.
365/365 — 0s 606us/step
RMSE: 13.86055419151385
```

Cabe precisar que existirá un modelo 3, pero será resultante de la calibración del modelo 2 lo cual se observará en las siguientes páginas.

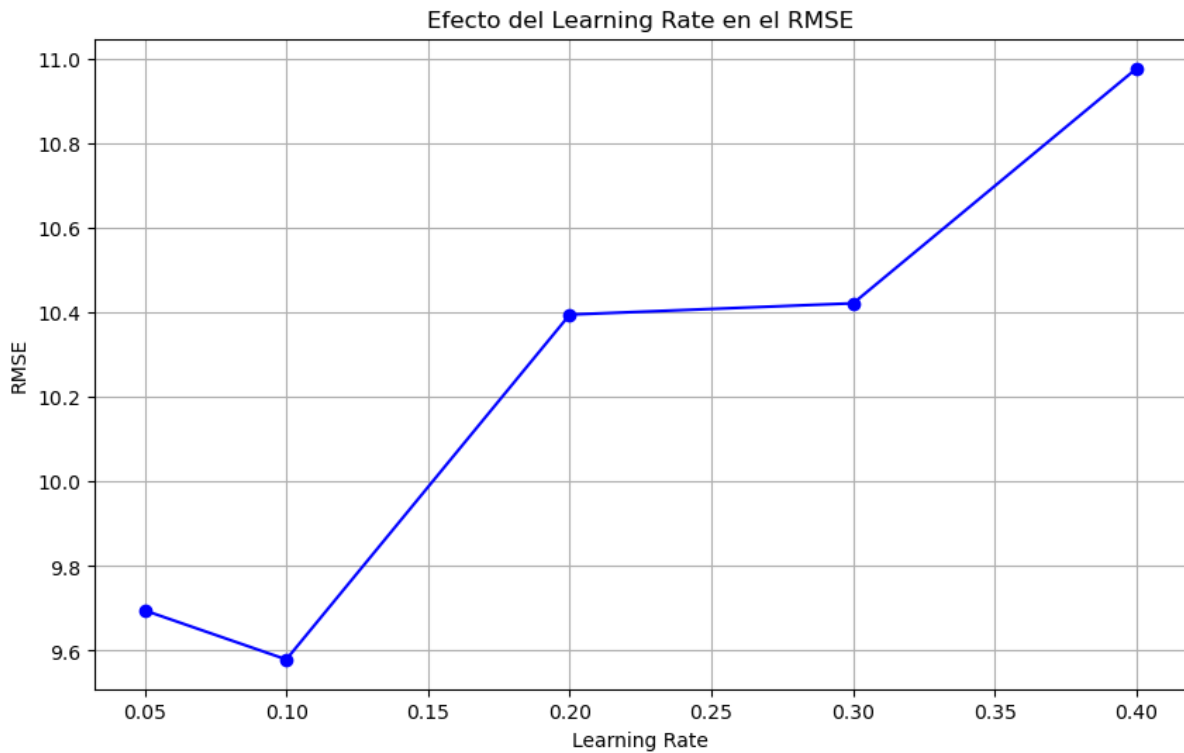
5. PARTE IV: DESARROLLO Y CALIBRACIÓN DE MODELOS

Se realizó la calibración de ambos modelos, por ejemplo a nivel de learning rate:

```
In [80]: learning_rates = [0.05, 0.1, 0.2, 0.3, 0.4]
rmse_results = []
for lr in learning_rates:
    print(f"Entrenando modelo con learning rate: {lr}")
    model1 = Sequential([
        Dense(1, input_dim=X_train_scaled.shape[1], activation='linear')
    ])
    model1.compile(optimizer=Adam(learning_rate=lr), loss='mean_squared_error')
    model1.fit(X_train_scaled, y_train, epochs=50, batch_size=32, validation_data=(X_val_scaled, y_val), verbose=0)
    y_pred = model1.predict(X_val_scaled)
    rmse = sqrt(mean_squared_error(y_val, y_pred))
    rmse_results.append(rmse)

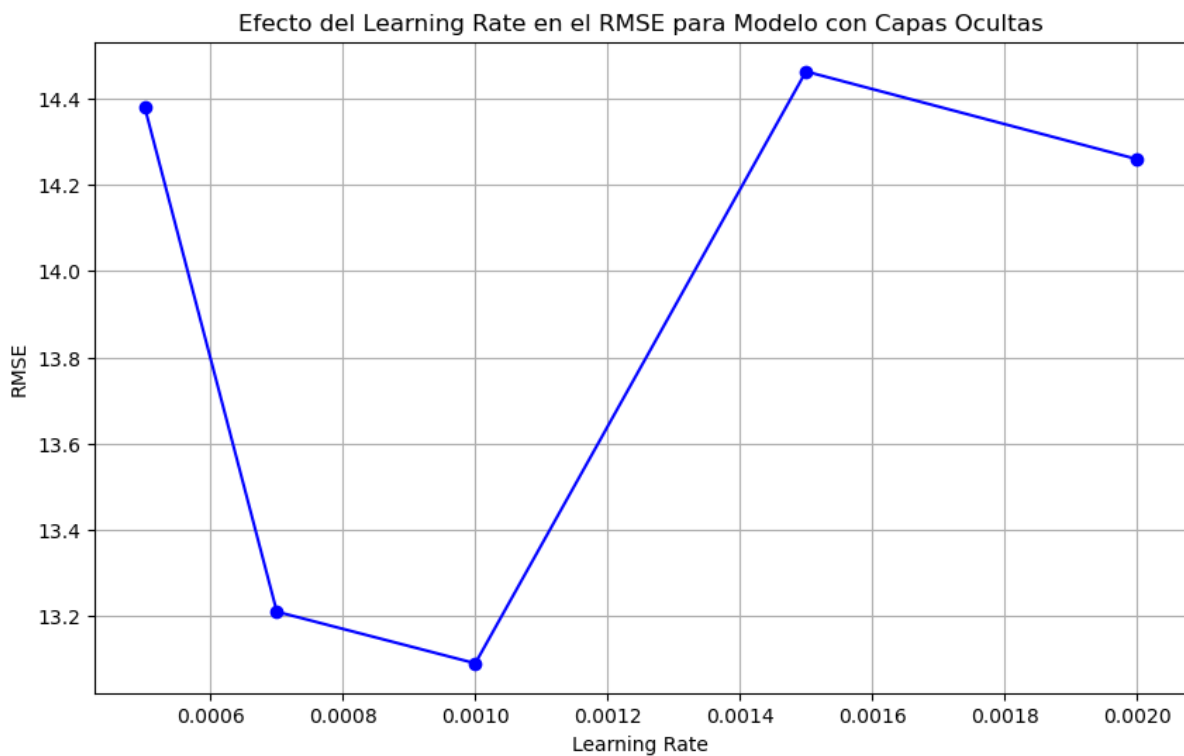
plt.figure(figsize=(10, 6))
plt.plot(learning_rates, rmse_results, marker='o', linestyle='-', color='b')
plt.title('Efecto del Learning Rate en el RMSE')
plt.xlabel('Learning Rate')
plt.ylabel('RMSE')
plt.grid(True)
plt.show()
```


Gráfico n.º 7 – Calibración Learning Rate Modelo 1



Y para el modelo 2 con capas ocultas (múltiples capas):

Gráfico n.º 78– Calibración Learning Rate Modelo 2



Como podemos observar, en ambos casos el óptimo del Learning Rate es 0.001; ahora procedemos a calibrar aún más el modelo 2:

```

# Configuraciones de neuronas para cada modelo
configurations = {
    'Config 1': [160, 80, 45, 25, 20], # Configuración original
    'Config 2': [165, 80, 45, 25, 20], # Ligeramente más grande
    'Config 3': [170, 80, 45, 25, 20], # Ligeramente más pequeño
    'Config 4': [175, 80, 45, 25, 20], # Configuración original
    'Config 5': [180, 80, 45, 25, 20], # Ligeramente más grande
    'Config 6': [160, 85, 45, 25, 20], # Ligeramente más pequeño
    'Config 7': [160, 75, 45, 25, 20], # Configuración original
    'Config 8': [160, 70, 45, 25, 20], # Ligeramente más grande
    'Config 9': [160, 80, 45, 20, 20], # Ligeramente más pequeño
    'Config 10': [160, 80, 45, 20, 15], # Configuración original
    'Config 11': [160, 80, 45, 25, 10], # Ligeramente más grande
    'Config 12': [160, 90, 45, 25, 20] # Ligeramente más pequeño
}

learning_rate = 0.001
results = []

# Pruebas de cada configuración
for config_name, layers in configurations.items():
    print(f"Probando modelo {config_name} con configuración de capas: {layers}")
    model = Sequential()
    model.add(Dense(layers[0], input_dim=X_train_scaled.shape[1], activation='elu'))
    model.add(BatchNormalization())

    for neurons in layers[1:-1]:
        model.add(Dense(neurons, activation='elu'))
        model.add(BatchNormalization())

    # Última capa antes de la salida
    model.add(Dense(layers[-1], activation='relu'))
    model.add(BatchNormalization())

    model.add(Dense(1, activation='linear')) # Salida para regresión

    optimizer = Adam(learning_rate=learning_rate)
    model.compile(optimizer=optimizer, loss='mean_squared_error')
    early_stopping = EarlyStopping(monitor='val_loss', patience=10, verbose=1, restore_best_weights=True)
    model.fit(X_train_scaled, y_train, epochs=50, batch_size=32, validation_data=(X_val_scaled, y_val), callbacks=[early_stopping])

    y_pred = model.predict(X_val_scaled)
    if np.isnan(y_pred).any():
        raise ValueError("Las predicciones contienen NaN")

    rmse = sqrt(mean_squared_error(y_val, y_pred))
    results.append((config_name, layers, rmse))
    print(f"RMSE para {config_name}: {rmse}")

# Encontrar la mejor configuración
best_result = min(results, key=lambda x: x[2])
print(f"Mejor configuración: {best_result[0]} con capas {best_result[1]}, RMSE: {best_result[2]}")

```

El código anterior fue diseñado para explorar la configuración óptima de una red neuronal multicapa. Una vez establecida la estructura básica, se procedió a calibrar ciertos parámetros y a incorporar diversas técnicas para mejorar el rendimiento y la generalización del modelo. Entre estas técnicas se incluyeron el uso de dropout y normalización por lotes (batch normalization) para regularizar el modelo, la implementación de paradas tempranas (early stopping) para evitar el sobreajuste, la optimización mediante el algoritmo Adam, y la selección cuidadosa de las funciones de activación. Estas estrategias fueron aplicadas con el objetivo de afinar el modelo y alcanzar los mejores resultados posibles en términos de precisión y eficiencia durante el entrenamiento.

En el modelo final calibrado (Modelo 3), tenemos lo siguiente:

```
In [30]: from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Dense, BatchNormalization, Dropout
from tensorflow.keras.optimizers import Adam
from tensorflow.keras.callbacks import EarlyStopping
from sklearn.metrics import mean_squared_error
from math import sqrt

# Modelo modificado
model = Sequential([
    Dense(360, input_dim=X_train_scaled.shape[1], activation='relu'),
    Dropout(0.3), # Agregar dropout para regularización
    BatchNormalization(),
    Dense(180, activation='elu'),
    Dropout(0.3), # Agregar dropout para regularización
    BatchNormalization(),
    Dense(60, activation='relu'),
    BatchNormalization(),
    Dense(25, activation='relu'),
    BatchNormalization(),
    Dense(20, activation='relu'),
    BatchNormalization(),
    Dense(1, activation='linear') # Salida para regresión
])

# Compilación con Adam
optimizer = Adam(learning_rate=0.001)
model.compile(optimizer=optimizer, loss='mean_squared_error')

# Callback de EarlyStopping modificado
early_stopping = EarlyStopping(monitor='val_loss', patience=10, verbose=1, restore_best_weights=True)

# Entrenamiento con EarlyStopping
model.fit(X_train_scaled, y_train, epochs=50, batch_size=50, validation_data=(X_val_scaled, y_val), callbacks=[early_stopping])

# Evaluación
y_pred = model.predict(X_val_scaled)
if np.isnan(y_pred).any():
    raise ValueError("Las predicciones contienen NaN")

rmse = sqrt(mean_squared_error(y_val, y_pred))
print(f'RMSE: {rmse}')
```

Sobre el modelo:

- **Dense:** Representa una capa densamente conectada (o completamente conectada). Cada neurona en una capa densa recibe entrada de todas las neuronas de la capa anterior, lo cual es una característica de las redes neuronales tradicionales.
- **La primera capa densa con 360 neuronas,** donde `input_dim` especifica el número de características de entrada del modelo. La función de activación 'relu' (Rectified Linear Unit) es comúnmente utilizada por su eficiencia y efectividad en redes neuronales profundas.
- **Dropout(0.3):** Una técnica de regularización donde aleatoriamente se "apagan" un porcentaje de las neuronas durante el entrenamiento para prevenir el sobreajuste. 0.3 significa que el 30% de las neuronas en la capa anterior se desactivarán aleatoriamente en cada paso durante el entrenamiento.
- **BatchNormalization():** Normaliza las activaciones de la capa anterior al restar la media del batch y dividir por la desviación estándar, lo que ayuda a mejorar la velocidad, rendimiento y estabilidad del entrenamiento.

Considerando que tiene normalización, aplicación técnicas y ajustes, se considero a este modelo 2 calibrado, como el modelo 3 “Múltiples capas Adam y calibración de hiperparámetros”

Siendo el mejor modelo obtenido, se generaron los predicts

```
In [31]: test_data2 = pd.read_csv('testReg.csv')
test_ids = test_data2['ID'].copy()
test_data = test_data[selected_features2]

# Asumiendo que 'scaler' y 'model' ya están definidos y entrenados
X_test_scaled = scaler.transform(test_data)
predictions = model.predict(X_test_scaled).flatten()

# Creación de un DataFrame para guardar las predicciones junto con los IDs 8.67698
output = pd.DataFrame({'ID': test_ids, 'Y': np.round(predictions).astype(int)})
output.to_csv('predictionsplease292.csv', index=False)
print("Predicciones guardadas en 'predictionsplease292.csv'.")

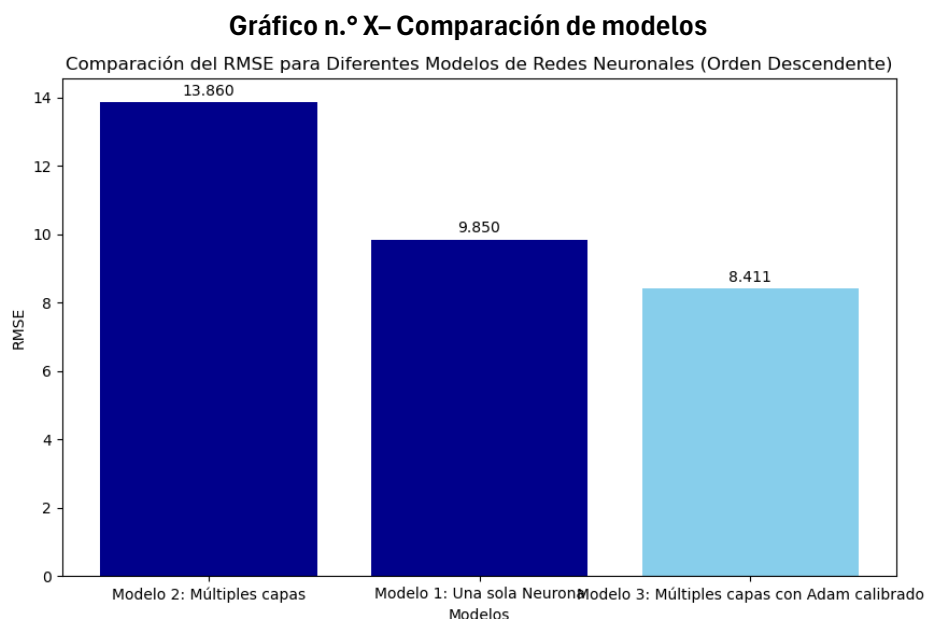
584/584 ————— 0s 741us/step
Predicciones guardadas en 'predictionsplease292.csv'.
```

Hito:

- RMSE EN LAB LOCAL - 8.41185
- RMSE EN KAGGLE - 8.64070
- RANKING EN EL PUBLIC LEADERBOARD 5 DE 21 EQUIPOS.
- DIFERENCIA CON EL RANKING 1: 8.64070 - 8.4901 = 0.1505

6. PARTE V: VISUALIZACIÓN DE RESULTADOS

A continuación se presentan los resultados obtenidos por cada tipo de modelo empleado:



Como podemos observar, inicialmente daba la impresión de que el modelo de una sola neurona sería mejor que de múltiples capas, pero a lo largo de las calibraciones realizadas, empleando batchnormalization, técnica de dropout y búsqueda intensiva de parámetros óptimos que incluye el learning rate, se obtuvieron los siguientes resultados:

1. Modelo 1: Una sola Neurona / RMSE OBTENIDO: 9.85
2. Modelo 2: Múltiples capas / RMSE OBTENIDO: 13.86
3. Modelo 3: Múltiples capas Adam y calibración de hiperparámetros / RMSE OBTENIDO: 8.411

En ese sentido, el modelo 3 al tener el mejor RMSE, fue utilizado para generar las predicciones y así obtener el puntaje de 8.64070, con lo cual nuestro equipo se colocó en el ranking 5 de la competencia.

Tabla 2. Modelos y los resultados obtenidos

Detalle del Modelo	RMSE JupyterNotebook	Puntuación Kaggle
Modelo 01: Red neuronal – Una sola neurona	9.85	10.11
Modelo 02: Red neuronal – Multicapa	13.86	15.86
Modelo 03: Red neuronal – Multicapa Adam calibrado	8.411	8.64

Recomendaríamos utilizar el modelo de redes neuronales Adam, calibrando los parámetros de número de capas, tasa de aprendizaje, neuronas por capa, tipo de función de activación, implementando estrategias de batchnormalization, dropout y earlystopping, a efectos de mejorar aún más el modelo.

Ahora bien, respondiendo a la pregunta de interés... sí, es posible siendo que las variables indican que existiría una relación entre las características musicales de una canción y el año en que fue publicada/lanzada.

Anexos.

[Anexo 1. Resultado en Competencia \(pdf\)](#)

[Anexo 2. HTML Notebook](#)

[Anexo 3. Notebook en Jupyter](#)