

Universidad Nacional Autónoma de México

Por mi raza hablara el espiritu



Facultad de Estudios Superior de Acatlan
Matematicas Aplicadas Y Computación

IDENTIFICACIÓN DE PATRONES DE ÉXITO MUSICAL EN SPOTIFY MEDIANTE ANÁLISIS DE DATOS Y CLASIFICACIÓN PREDICTIVA

Solis Vilchis Roberto Atonatiuh
Ortiz Garcia Miguel Angel

Índice general

1.	Resumen Ejecutivo (Abstract)	2
2.	Introducción	2
3.	Preparación de Datos (Data Preparation)	3
4.	Análisis de Datos (Data Analysis)	9
5.	Reducción de Dimensiones (Dimensionality Reduction)	15
6.	Clustering	16
7.	Modelación de Clasificación (Classification Modeling)	19
8.	Evaluación del Modelo (Model Evaluation)	23
9.	Conclusiones	27

1. Resumen Ejecutivo (Abstract)

Este proyecto presenta un análisis y modelación de datos sobre el éxito de canciones en plataformas de streaming como Spotify, Apple Music, Shazam y Deezer. El objetivo principal es identificar las características que influyen en el éxito de las canciones y predecir su popularidad utilizando diversas técnicas de análisis y aprendizaje supervisado.

Inicialmente, se realiza una exploración de datos para entender mejor las características musicales como el BPM (beats per minute), la danzabilidad y otros atributos relevantes. A continuación, se aplica el Análisis de Componentes Principales (PCA) para la reducción de dimensiones, permitiendo simplificar los datos sin perder información crítica.

Se implementan técnicas de clustering para categorizar las canciones en grupos según su potencial de éxito. Posteriormente, se utilizan varios modelos de clasificación supervisada, incluyendo Regresión Logística, Random Forest, Redes Neuronales, Análisis Discriminante y Support Vector Machine, para predecir la categoría de éxito de cada canción.

Para evaluar el rendimiento de los modelos, se emplean métricas como la precisión, la recall, el F1-score y el área bajo la curva ROC (ROC-AUC). Estos análisis permiten determinar cuál es el modelo más eficaz para la clasificación de canciones.

Los resultados de este estudio pueden ayudar a identificar canciones con alto potencial de éxito en las diferentes plataformas de streaming, ofreciendo valiosas recomendaciones tanto para productores musicales como para plataformas de streaming en la toma de decisiones estratégicas.

2. Introducción

Vivimos en una época donde la música se ha vuelto cada vez más relevante. En un mundo caracterizado por altos niveles de estrés, la música ofrece un medio para relajarse y aliviar esos momentos de tensión. Por esta razón, las plataformas de streaming han proliferado, convirtiéndose en medios cruciales para la difusión de canciones.

Dada la importancia de la música, surge la pregunta de cómo saber si una canción será relevante o, como ocurre con la mayoría, pasará desapercibida sin alcanzar el éxito. Este proyecto busca abordar esta cuestión utilizando técnicas de machine learning para predecir la posibilidad de que una canción

tenga éxito en alguna plataforma de streaming.

El objetivo principal de este proyecto es desarrollar un modelo que pueda predecir el éxito de una canción, proporcionando así una herramienta útil para discográficas y estrategias de marketing en plataformas como Spotify, Apple Music y otras. Además, el análisis de los datos permitirá observar cómo varían los rankings de popularidad entre diferentes plataformas.

Para alcanzar este objetivo, es fundamental seguir un camino metodológico claro y riguroso. Este proyecto se divide en varias etapas: exploración de datos, análisis, reducción de dimensiones, clustering, modelación con algoritmos de clasificación y evaluación de los modelos creados. Cada una de estas etapas será explicada en detalle a lo largo del documento.

Este documento está organizado de la siguiente manera: primero, se describe el análisis de los datos y las características musicales utilizadas. Luego, se detalla la metodología empleada para la reducción de dimensiones y el clustering. Posteriormente, se presentan los modelos de clasificación utilizados y sus evaluaciones. Finalmente, se discuten los resultados y se presentan las conclusiones y posibles trabajos futuros.

3. Preparación de Datos (Data Preparation)

Nuestro conjunto de datos proviene de la plataforma Kaggle y contiene información sobre las canciones en Spotify. El dataset se puede encontrar en Kaggle. Incluye datos sobre las canciones más populares del año 2023, así como de años anteriores como 2022 y 2021. Además, contiene el número de reproducciones y la cantidad de listas de reproducción en otras plataformas donde se encuentra cada canción.

El conjunto de datos se carga utilizando la biblioteca pandas y se realiza una descripción inicial de las variables.

```
1 import pandas as pd
2
3 data = pd.read_csv('./spotify-2023.csv', encoding='ISO-8859-1')
4
5 print(data.info())
6 print(data.describe())
```

Listing 1: Carga de datos y descripción inicial

Para entender mejor las variables presentes en nuestro conjunto de datos, se ha creado un diccionario de datos que describe el significado y el tipo de cada columna. La siguiente tabla muestra las variables, sus significados y sus categorías.

Variable	Descripción	Tipo
track_name	Nombre de la canción	Categorica
artist(s)_name	Nombre del/los artista(s) de la canción	Categorica
artist_count	Número de artistas que contribuyen a la canción	Numérica
released_year	Año de lanzamiento de la canción	Numérica
released_month	Mes de lanzamiento de la canción	Numérica
released_day	Día del mes de lanzamiento de la canción	Numérica
in_spotify_playlists	Número de listas de reproducción de Spotify en las que está incluida la canción	Numérica
in_spotify_charts	Presencia y rango de la canción en las listas de Spotify	Numérica
streams	Número total de reproducciones en Spotify	Numérica
in_apple_playlists	Número de listas de reproducción de Apple Music en las que está incluida la canción	Numérica
in_apple_charts	Presencia y rango de la canción en las listas de Apple Music	Numérica
in_deezer_playlists	Número de listas de reproducción de Deezer en las que está incluida la canción	Numérica
in_deezer_charts	Presencia y rango de la canción en las listas de Deezer	Numérica
in_shazam_charts	Presencia y rango de la canción en las listas de Shazam	Numérica
bpm	Beats por minuto, una medida del tempo de la canción	Numérica
key	Clave de la canción	Categorica
mode	Modo de la canción (mayor o menor)	Categorica
danceability %	Porcentaje que indica cuán adecuada es la canción para bailar	Numérica
valence %	Positividad del contenido musical de la canción	Numérica
energy %	Nivel de energía percibido de la canción	Numérica
acousticness %	Cantidad de sonido acústico en la canción	Numérica
instrumentalness %	Cantidad de contenido instrumental en la canción	Numérica
liveness %	Presencia de elementos de actuación en vivo	Numérica
speechiness %	Cantidad de palabras habladas en la canción	Numérica

Cuadro 1: Diccionario de datos

La limpieza de datos es crucial para garantizar la calidad y precisión de los análisis posteriores. Incluye la identificación y corrección de valores faltantes, incorrectos y atípicos.

1. Valores Incorrectos Se identificó un error en la columna streams, donde debería haber un valor numérico pero se encontró la cadena:

'BPM110KeyAModeMajorDanceability53Valence75Energy69Acousticness7Instrumentalness0Liveness17Speechiness3'. Este valor se cambió a NaN (valor faltante). Tambien se identificaron errores en el registro de datos en ciertas columnas, como el agregar coma "," en valores que deberian ser solo numericos.

```

1      vi = spo[spo['streams'] == '
          BPM110KeyAModeMajorDanceability53Valence75
          Energy69Acousticness7Instrumentalness0
          Liveness17Speechiness3'].index[0]
2
3      spo.loc[vi, ['streams']] = np.nan
4
5      spo['streams'] = spo['streams'].astype('Int64')
6
7      spo['in_deezer_playlists'] = spo['
          in_deezer_playlists'].astype(str)
8
9      spo['in_deezer_playlists'] = spo['
          in_deezer_playlists'].apply(lambda x: int(x.
          replace(',', '')))
10
11     spo['in_deezer_playlists'] = spo['
          in_deezer_playlists'].astype(int)
12
13     spo['in_shazam_charts'] = spo['in_shazam_charts'].
          astype(str)
14
15     spo['in_shazam_charts'].replace('nan', np.nan,
          inplace=True)
16
17     spo['in_shazam_charts'] = spo['in_shazam_charts'].
          apply(lambda x: int(str(x).replace(',', '')) if
          pd.notnull(x) else np.nan )
18
19     spo['in_shazam_charts'] = spo['in_shazam_charts'].
          astype('Int64')

```

Listing 2: Valores incorrectos

2. Valores Faltantes Se encontraron valores faltantes en las siguientes columnas:

```

1      nulos = spo.isnull().sum()
2      print(f'Columnas con al menos un nulo: \n{nulos[
          nulos > 0]}')

```

Listing 3: Valores nulos

```

Columnas con al menos un nulo:
streams          1

```

```
in_shazam_charts    50
key                 95
dtype: int64
```

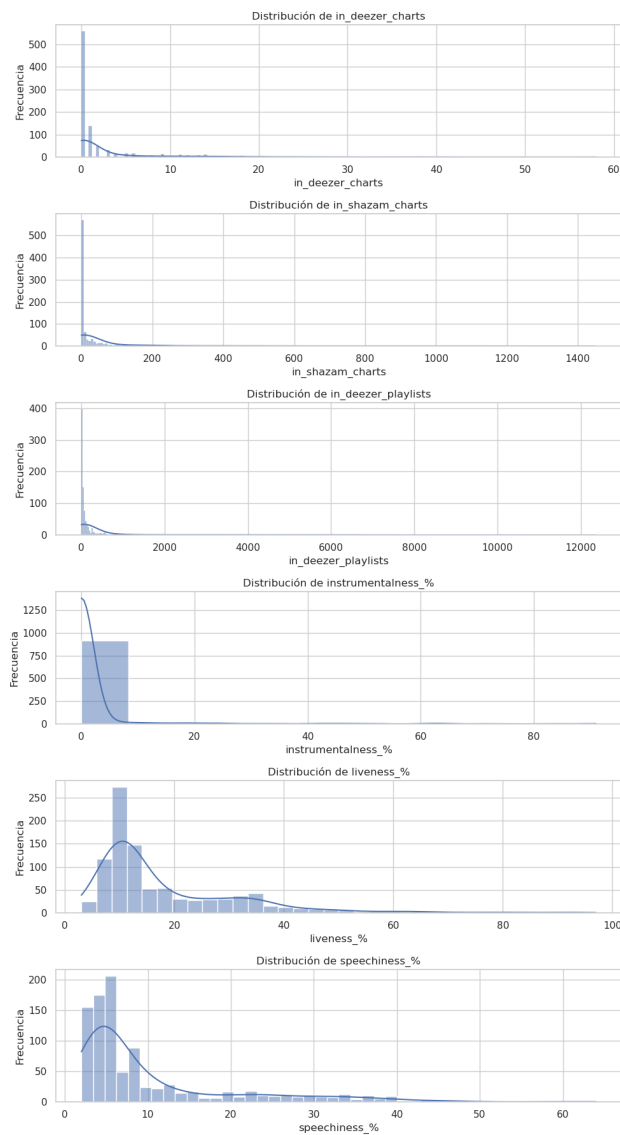
Para imputar estos valores faltantes, se utilizaron las siguientes técnicas:

- streams: Mediana
- in_shazam_charts: Media
- key: Moda

```
1      imputar_moda_cat = SimpleImputer(strategy='
2          most_frequent')
3
4      spo['key'] = imputar_moda_cat.fit_transform(spo[['
5          key']]).ravel()
6
7      spo['key'] = spo['key'].astype(str)
8
9      imput_streams = SimpleImputer(strategy='mean')
10
11     spo['streams'] = imput_streams.fit_transform(spo[['
12         streams']])
13
14     imput_shazam_charts = SimpleImputer(strategy='
15         median')
16
17     spo['in_shazam_charts'] = imput_shazam_charts.
18         fit_transform(spo[['in_shazam_charts']])
```

Listing 4: Metodos de imputacion

3. Valores Atípicos (Outliers) Para identificar valores atípicos, se utilizaron gráficos de distribución. Las columnas con valores potencialmente atípicos fueron tratadas con técnicas de Capping y Flooring.



Estas técnicas se aplicaron a las siguientes columnas:

- in_deezer_charts
- in_shazam_charts
- in_deezer_playlists
- instrumentalness_
- liveness_
- speechiness_


```

1 change_column = ['in_deezer_charts', '
2 in_shazam_charts', 'in_deezer_playlists', '
3 instrumentalness_%', 'liveness_%', '
4 speechiness_%']
5
6 def cap_floor(df, column, lower_quantile=0.01,
7 upper_quantile=0.99):
8     lower_bound = df[column].quantile(
9         lower_quantile)
10    upper_bound = df[column].quantile(
11        upper_quantile)
12    df[column] = np.where(df[column] < lower_bound
13        , lower_bound, df[column])
14    df[column] = np.where(df[column] > upper_bound
15        , upper_bound, df[column])
16    return df
17
18 for column in change_column:
19     spo = cap_floor(spo, column)

```

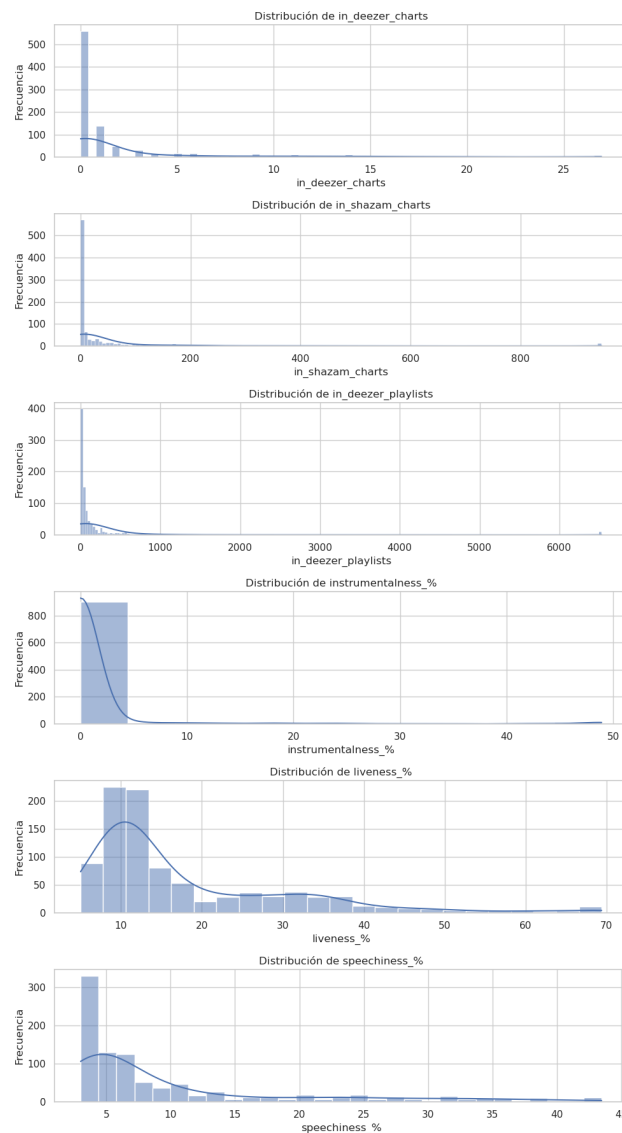
Listing 5: Tratamiento de valores atípicos

```

1 sns.set(style="whitegrid")
2
3 num_columns = len(change_column)
4 fig, axes = plt.subplots(num_columns, 1, figsize
5     =(10, 3 * num_columns))
6
7 for i, column in enumerate(change_column):
8     ax = axes[i] if num_columns > 1 else axes
9     if pd.api.types.is_numeric_dtype(spo[column]):
10         sns.histplot(spo[column], kde=True, ax=ax)
11         ax.set_xlabel(column)
12         ax.set_ylabel('Frecuencia')
13     else:
14         sns.countplot(y=spo[column], ax=ax)
15         ax.set_xlabel('Frecuencia')
16         ax.set_ylabel(column)
17         ax.set_title(f'Distribución de {column}')
18
19 plt.tight_layout()
20 plt.show()

```

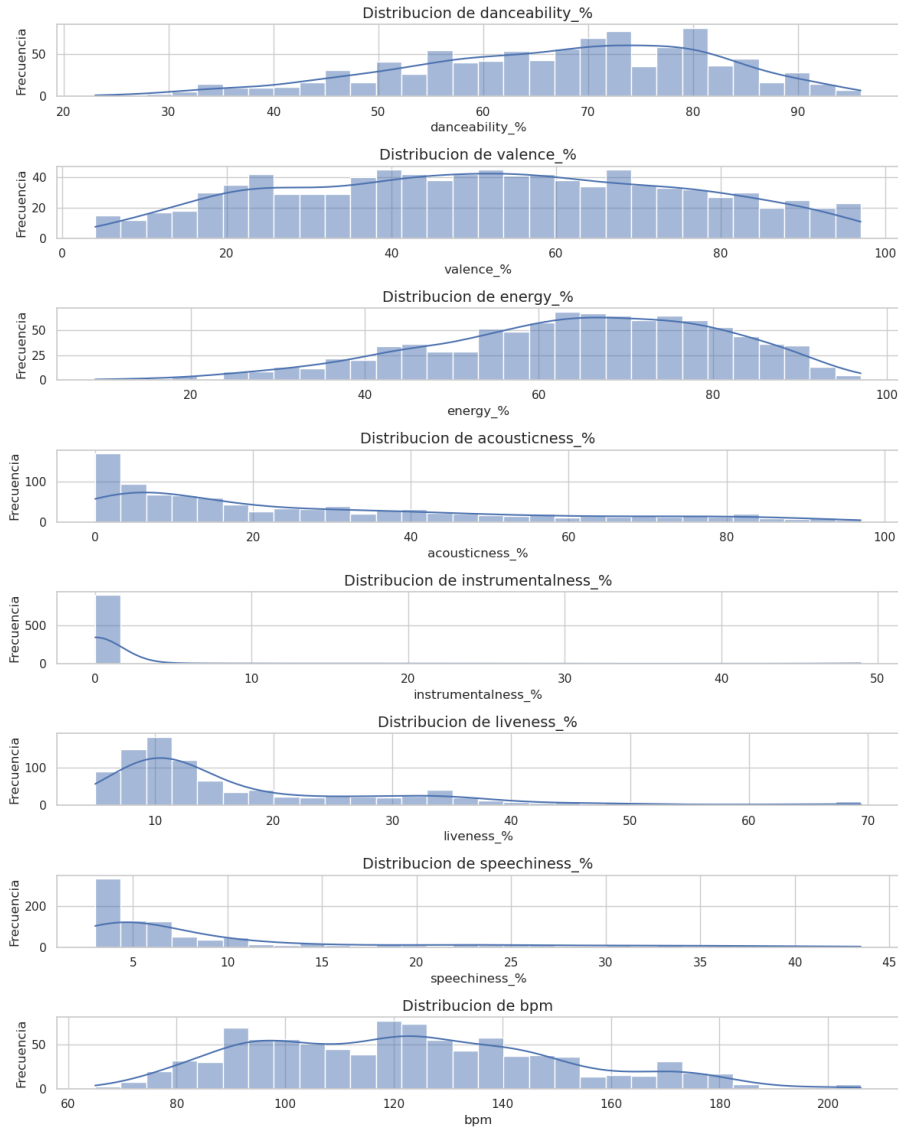
Listing 6: Código para las visualizaciones



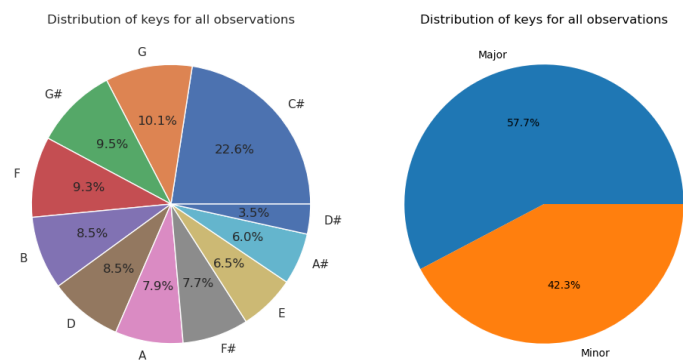
Al comparar las gráficas, podemos notar un cambio significativo en varias de ellas, pero también observamos que algunas se mantienen similares. Es evidente que el uso de la técnica de capping y flooring resulta muy útil para tratar valores atípicos, lo que nos permite continuar con un análisis más detallado de los datos.

4. Análisis de Datos (Data Analysis)

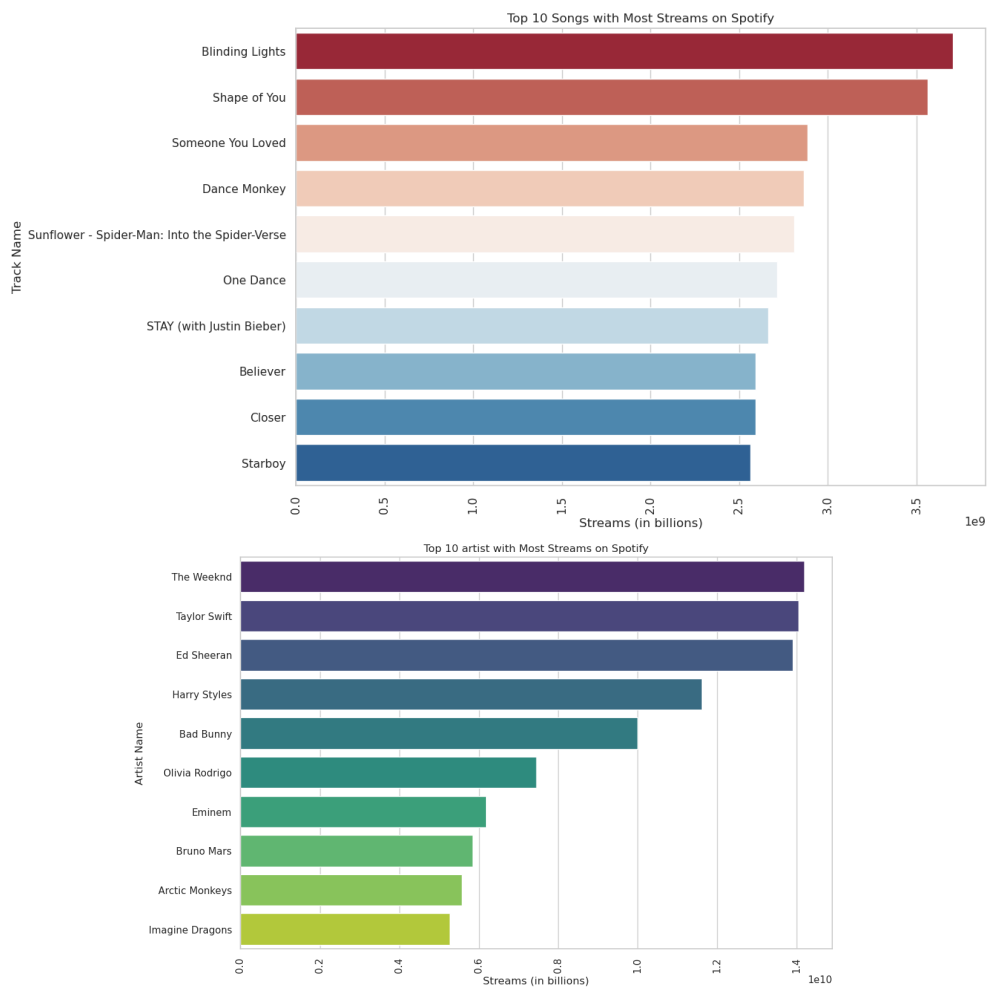
En esta sección, analizamos las distribuciones de las variables que representan los atributos de las canciones, como el porcentaje de energía, danzabilidad, entre otros. El objetivo es entender el comportamiento de estas variables y sus tipos de distribución, considerando que los datos ya han sido tratados.



Además de las variables continuas, también se analizaron las variables discretas, particularmente 'key' y 'mode'. Se utilizaron gráficos de pastel para visualizar los porcentajes de cada valor único, proporcionando una visión clara de la distribución de estas variables.



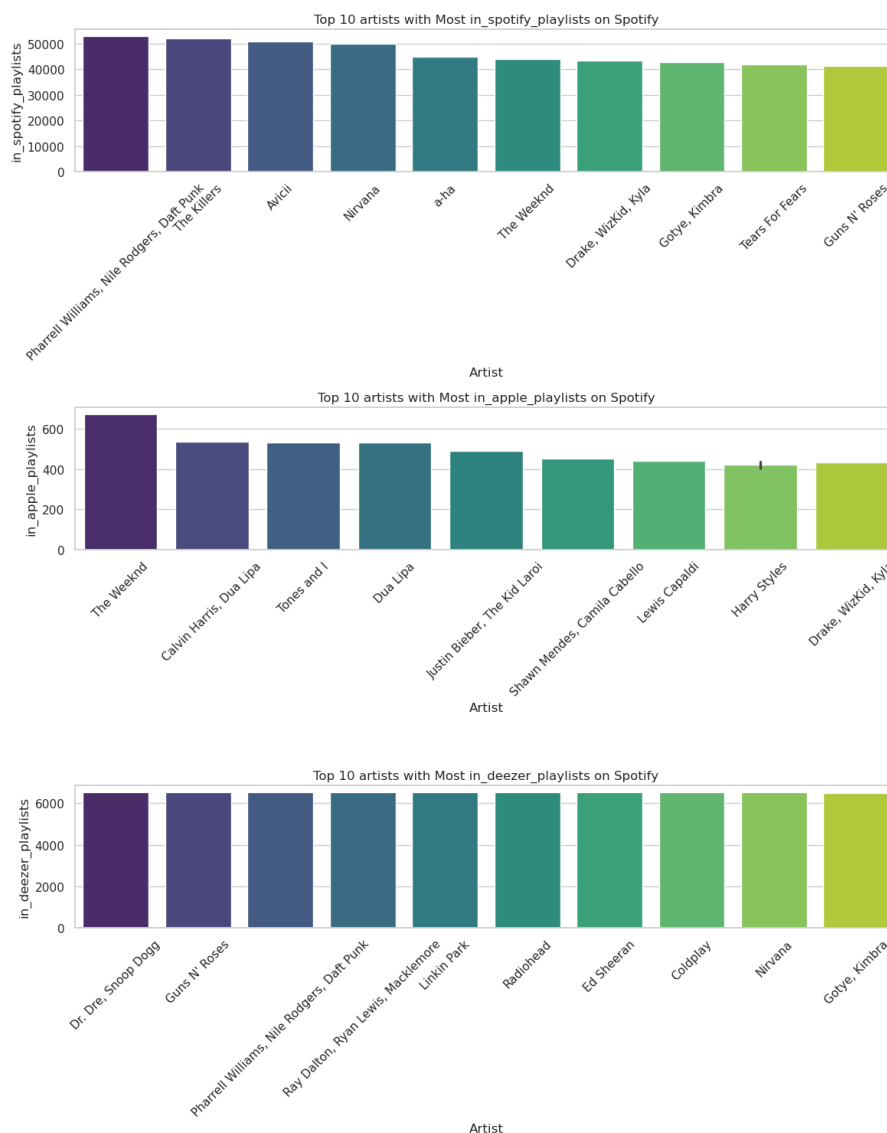
Conociendo las frecuencias y porcentajes de las variables, procedemos a analizar las variables relacionadas con streams, playlists, charts, artistas y canciones. Esto nos permitirá identificar a los mejores artistas, las mejores canciones en diferentes plataformas, entre otros aspectos destacados.



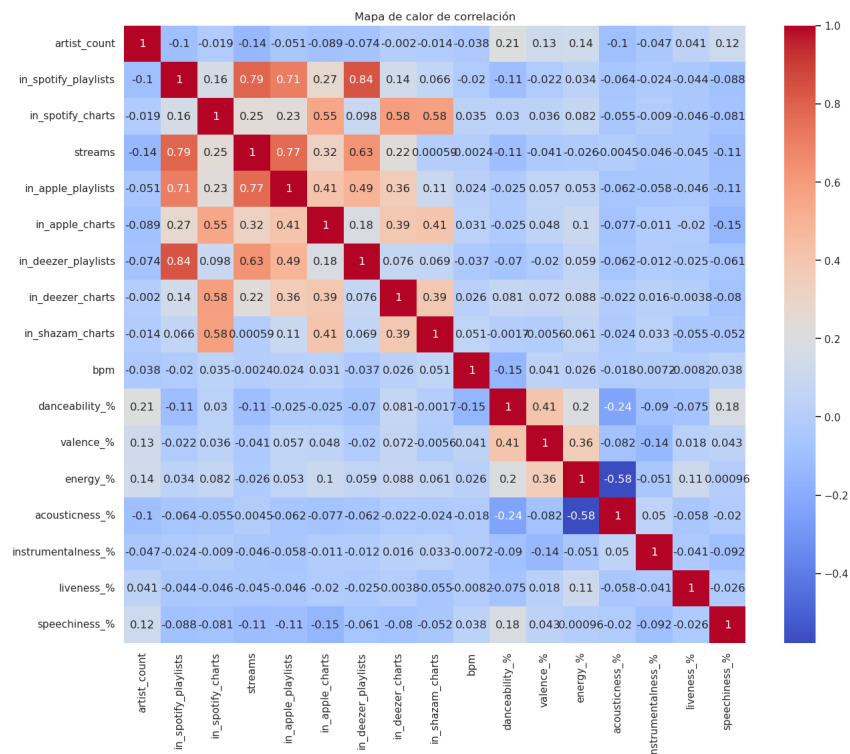


Algunos artistas muestran consistencia en su popularidad a través de múltiples plataformas. Por ejemplo, The Weeknd aparece frecuentemente en las listas de reproducción tanto de Spotify como de Apple Music.

Existen diferencias notables entre las plataformas. Artistas que son populares en Spotify pueden no serlo en Deezer y viceversa. Esto puede deberse a las diferencias en la base de usuarios y algoritmos de recomendación de cada plataforma.



Finalmente, examinamos las correlaciones entre las columnas utilizando mapas de calor. El mapa de calor de correlaciones proporciona una representación visual de las relaciones entre las variables del conjunto de datos. En este mapa, cada celda representa el coeficiente de correlación entre dos variables, donde un valor de 1 indica una correlación positiva perfecta, -1 indica una correlación negativa perfecta, y 0 indica que no hay correlación lineal entre las variables.



1. Correlación Alta:

Las variables relacionadas con playlists y charts en las diferentes plataformas tienen correlaciones significativamente altas. Por ejemplo, `in_spotify_playlists` tiene una alta correlación con `in_spotify_charts` (0.79) y `in_apple_playlists` (0.71). Esto sugiere que las canciones que son populares en una plataforma tienden a ser populares en otras también.

2. Correlación Moderada:

`streams` tiene una correlación moderada con `in_spotify_charts` (0.63) y `in_apple_playlists` (0.77), indicando que el número de reproducciones está relacionado con la aparición en listas de éxitos. Las características de la canción como `energy_ %` y `acousticness_ %` tienen una correlación inversa moderada (-0.58), sugiriendo que las canciones con alta energía tienden a tener baja acústica.

3. Correlación Baja:

Variables como `bpm`, `danceability_ %`, y `speechiness_ %` muestran baja correlación con la mayoría de las otras variables, indicando que estos atributos individuales de las canciones no están fuertemente relacionados con su popularidad o presencia en playlists y charts.

5. Reducción de Dimensiones (Dimensionality Reduction)

El uso de múltiples variables en un conjunto de datos ocasiona alta dimensionalidad, lo cual implica varios problemas en el momento de la modelación, un fenómeno conocido como la "maldición de la dimensionalidad". Para evitar estos problemas, en este proyecto se decidió realizar una reducción de dimensiones.

Primero, se normalizaron los datos, pero no todos los datos, sino los importantes para crear los clusters. Nos enfocamos en datos que representen si una canción tuvo éxito o no. El siguiente código muestra este proceso y las variables utilizadas:

```
1 from sklearn.preprocessing import StandardScaler
2
3 # Variables relevantes
4 rank = ['in_spotify_charts', 'in_apple_charts', '
        in_deezer_charts', 'in_shazam_charts', '
        in_spotify_playlists', 'in_apple_playlists', '
        in_deezer_playlists', 'streams']
5 X1 = spo[rank]
6
7 # Normalizar datos
8 scaler = StandardScaler()
9 X1_scaled = scaler.fit_transform(X1)
```

Listing 7: Normalizacion

La normalización es un paso crucial porque ajusta los valores de las variables a una escala común sin distorsionar las diferencias en los rangos de valores. Esto es especialmente importante para técnicas de reducción de dimensiones.

Para la reducción de dimensiones, utilizamos la técnica de PCA (Análisis de Componentes Principales). El PCA es una técnica que ayuda a mantener los componentes principales de cada variable, reduciéndolos a un menor número de variables que aún retienen la mayor parte de la variabilidad presente en los datos originales.

```
1 from sklearn.decomposition import PCA
2
3 # Aplicar PCA
4 pca = PCA(n_components=3)
5 X1_pca = pca.fit_transform(X1_scaled)
6
7 principal_components = pd.DataFrame(X1_pca,
```



```

8 columns=['componente_1',
          , 'componente_2', '
          componente_3'])

```

Listing 8: Reduccion de dimensiones

El PCA es muy útil ya que permite reducir los datos a componentes principales, eliminando los datos redundantes o menos significativos. En este caso, reducimos las dimensiones a tres componentes principales, lo cual facilita el análisis posterior y mejora la eficiencia de los modelos de clasificación.

6. Clustering

El clustering es una técnica de aprendizaje no supervisado que se utiliza para agrupar objetos similares en clusters. En el contexto de nuestro análisis de canciones en Spotify, queremos identificar grupos de canciones que comparten características similares, lo cual puede ser útil para entender patrones de popularidad y éxito.

Para determinar el número óptimo de clusters, utilizamos el método del Silhouette Score. Este método evalúa la coherencia dentro de los clusters (qué tan similares son los puntos dentro de un mismo cluster) y la separación entre los clusters (qué tan diferentes son los puntos de diferentes clusters).

El siguiente gráfico muestra el Silhouette Score para diferentes números de clusters:

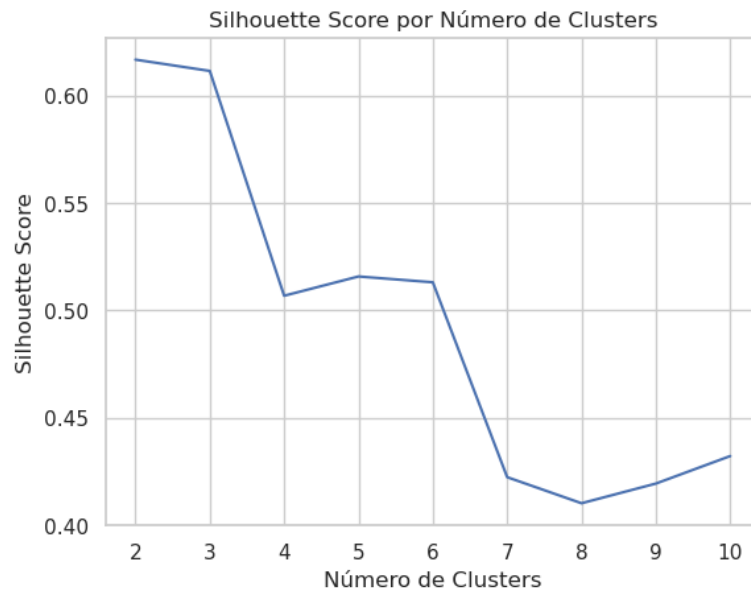
```

1 from sklearn.cluster import KMeans
2 from sklearn.metrics import silhouette_score
3
4 silhouette_scores = []
5 range_n_clusters = list(range(2, 11))
6
7 for n_clusters in range_n_clusters:
8     kmeans = KMeans(n_clusters=n_clusters, random_state=4,
9                     n_init=10, max_iter=1000)
10    kmeans.fit(principal_components[['componente_1', '
11        componente_2', 'componente_3']])
12    cluster_labels = kmeans.labels_
13    silhouette_avg = silhouette_score(principal_components
14        [['componente_1', 'componente_2', 'componente_3']],
15        cluster_labels)
16    silhouette_scores.append(silhouette_avg)
17
18 plt.plot(range_n_clusters, silhouette_scores)
19 plt.xlabel('N mero de Clusters')
20 plt.ylabel('Silhouette Score')
21 plt.title('Silhouette Score por N mero de Clusters')

```

```
18 plt.show()
```

Listing 9: Números de clusters



El gráfico sugiere que el número óptimo de clusters es 2, ya que muestra el mayor Silhouette Score en ese punto.

Aplicamos el algoritmo K-Means con 2 clusters y visualizamos los resultados en un gráfico tridimensional:

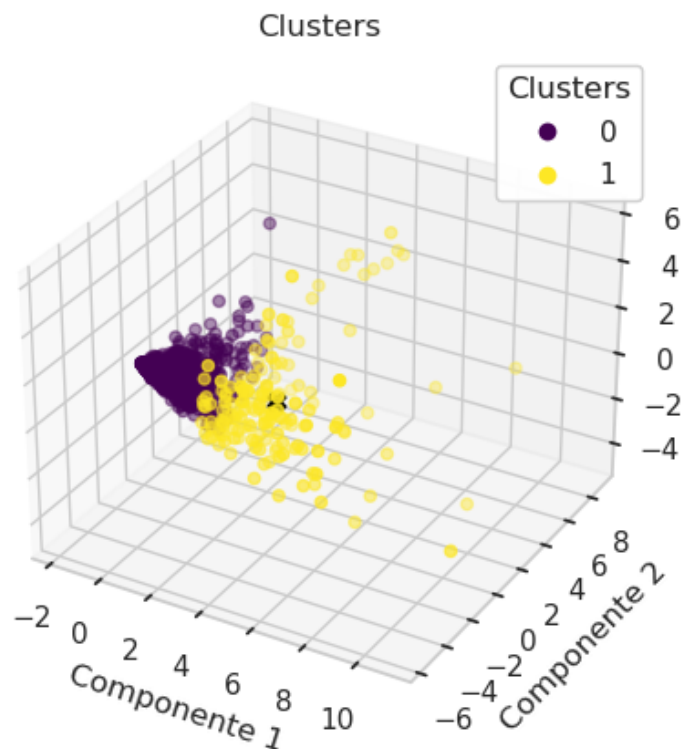
```
1 from sklearn.cluster import KMeans
2
3 kmeans = KMeans(n_clusters=2, random_state=4, n_init=10,
4                 max_iter=1000)
5 kmeans.fit(principal_components)
6
7 clusters = kmeans.labels_
8 centroids = kmeans.cluster_centers_
9
10 principal_components['cluster'] = clusters
11
12 fig = plt.figure()
13 ax = fig.add_subplot(111, projection='3d')
14
15 scatter = ax.scatter(principal_components['componente_1'],
16                     principal_components['componente_2'],
17                     principal_components['componente_3'],
18                     c=clusters, cmap='viridis')
19
20 ax.scatter(centroids[:, 0], centroids[:, 1], centroids[:, 2],
```

```

20         s=50, color='black', marker='X', label='
           Centroids')
21
22     ax.set_xlabel('Componente 1')
23     ax.set_ylabel('Componente 2')
24     ax.set_zlabel('Componente 3')
25     ax.set_title('Clusters')
26
27     # A adir una leyenda
28     legend1 = ax.legend(*scatter.legend_elements(), title="
           Clusters")
29     ax.add_artist(legend1)
30
31     plt.show()

```

Listing 10: Clusters



El Silhouette Score del modelo K-Means con 2 clusters es 0.6166, lo que indica una buena cohesión y separación entre los clusters:

```

1     from sklearn.metrics import silhouette_score
2
3     silhouette_avg = silhouette_score(principal_components[['
           componente_1', 'componente_2', 'componente_3']],
           clusters)
4     print(f'Silhouette Score: {silhouette_avg}')

```

Listing 11: Evaluacion de clusters

En resumen, mediante el clustering, hemos podido identificar grupos de canciones con características similares. Este análisis proporciona una base sólida para explorar más a fondo los patrones de éxito de las canciones en Spotify y otras plataformas.

7. Modelación de Clasificación (Classification Modeling)

En esta sección, se describen los pasos seguidos para el modelado de clasificación y se presentan brevemente los modelos utilizados.

Preparación de Datos Primero, los datos fueron estandarizados y divididos en conjuntos de entrenamiento y prueba.

```
1 from sklearn.preprocessing import StandardScaler
2 from sklearn.model_selection import train_test_split
3
4 scaler = StandardScaler()
5 X = scaler.fit_transform(spo[per])
6 y = spo['clusters']
7 X_train, X_test, y_train, y_test = train_test_split(X, y,
    test_size=0.3, random_state=42)
```

Listing 12: Preparacion de datos

Búsqueda de Mejores Parámetros Para optimizar los modelos, se utilizaron las técnicas de validación cruzada 'GridSearchCV' y 'RandomizedSearchCV'.

```
1 from sklearn.model_selection import GridSearchCV,
2   RandomizedSearchCV
3
4 def entrenar(param, modelo, X, y):
5     grid = RandomizedSearchCV(param_distributions=param,
6                               n_jobs=-1,
7                               n_iter=10,
8                               cv=4,
9                               estimator=modelo,
10                              error_score='raise')
11
12     grid.fit(X, y)
13     return grid, grid.best_estimator_, grid.best_score_,
14           grid.best_params_
15
16 def entrenar_g(param, modelo, X, y):
17     grid = GridSearchCV(param_grid=param,
```

```

15         n_jobs=-1,
16         cv=4,
17         estimator=modelo,
18         error_score='raise')
19     grid.fit(X, y)
20     return grid, grid.best_estimator_, grid.best_score_,
        grid.best_params_

```

Listing 13: Validacion cruzada

Parámetros de los Modelos Se definieron los siguientes parámetros para cada modelo:

```

1  # SVM (M quinas de vectores de soporte)
2  svm_param = {'C': [0.1, 1, 10, 100], 'gamma': [0.1, 0.01,
3              0.001, 0.0001], 'kernel': ['linear', 'rbf']}
4
5  # Random Forest
6  rf_param = dict(n_estimators=list(range(1, 100, 25)),
7                  criterion=['gini', 'entropy'],
8                  max_depth=[x for x in list(range(2, 5))] +
9                      [None],
10                 min_samples_split=[x for x in list(range(2,
11                     4))],
12                 min_samples_leaf=[x for x in list(range(2,
13                     4))],
14                 max_features=[None] + [i * .05 for i in
15                     list(range(2, 4))],
16                 max_leaf_nodes=list(range(2, 10)) + [None],
17                 min_impurity_decrease=[x * .10 for x in
18                     list(range(2, 4))],
19                 oob_score=[True, False],
20                 warm_start=[True, False],
21                 class_weight=[None, 'balanced'],
22                 max_samples=[None])
23
24  # Regresión Logística
25  rl_param = {'C': [0.1, 1, 10, 100]}
26
27  # Red Neuronal
28  rn_param = dict(activation = ['identity', 'logistic', 'tanh',
29                              'relu'],
30                  solver = ['lbfgs', 'sgd', 'adam'],
31                  alpha = np.arange(0.0001, 0.001, 0.0001),
32                  learning_rate = ['constant', 'invscaling',
33                              'adaptive'])
34
35  # Ada Boost
36  ab_param = dict(n_estimators = range(2,10),
37                  learning_rate = np.arange(0.1,1,0.1),
38                  algorithm = ['SAMME.R'])
39
40  # Análisis Discriminante
41  ad_param = dict(solver = ['svd', 'lsqr', 'eigen'])

```

Listing 14: Parametros

Balanceo de Clases Se utilizó la técnica 'SMOTETomek' para balancear las clases en los datos de entrenamiento.

```
1 from imblearn.combine import SMOTETomek
2
3 os_us = SMOTETomek()
4 X_train_res, y_train_res = os_us.fit_resample(X_train,
        y_train)
```

Listing 15: Balance de clases

Descripción de los Modelos de Clasificación A continuación, se describen brevemente los modelos de clasificación utilizados en el análisis.

1. Support Vector Machine (SVM)

- Descripción: Encuentra un hiperplano óptimo que separa las clases de datos.
- Ventajas: Eficaz en espacios de alta dimensión y con un margen de separación claro.
- Desventajas: No adecuado para conjuntos de datos muy grandes.
- Mejores Parámetros:

```
1 {'C': 0.1, 'gamma': 0.1, 'kernel': 'rbf'}
```

- Precisión: 0.5771375464684014

2. Random Forest

- Descripción: Utiliza un conjunto de árboles de decisión para mejorar la precisión.
- Ventajas: Maneja bien grandes conjuntos de datos y es robusto ante el overfitting.
- Desventajas: Puede ser lento y consumir mucha memoria.
- Mejores Parámetros:

```
1 {'warm_start': True,
2   'oob_score': True,
3   'n_estimators': 1,
4   'min_samples_split': 2,
5   'min_samples_leaf': 3,
6   'min_impurity_decrease': 0.2,
7   'max_samples': None,
8   'max_leaf_nodes': 5,
9   'max_features': None,
10  'max_depth': 2,
11  'criterion': 'entropy',
12  'class_weight': 'balanced'}
```

- Precisión: 0.5018587360594795

3. Regresión Logística

- Descripción: Modela la probabilidad de una clase específica utilizando una función logística.
- Ventajas: Simple de implementar y coeficientes fácilmente interpretables.
- Desventajas: No funciona bien con relaciones no lineales.
- Mejores Parámetros:

```
1 {'C': 1}
```

- Precisión: 0.5585501858736059

4. Red Neuronal (Neural Network)

- Descripción: Modela relaciones complejas utilizando capas de neuronas.
- Ventajas: Capaz de modelar relaciones no lineales complejas.
- Desventajas: Requiere grandes cantidades de datos y es computacionalmente intensiva.
- Mejores Parámetros:

```
1 {'solver': 'lbfgs',
2   'learning_rate': 'constant',
3   'alpha': 0.0006000000000000001,
4   'activation': 'relu'}
```

- Precisión: 0.8596654275092936

5. AdaBoost

- Descripción: Combina varios modelos débiles para crear un clasificador fuerte.
- Ventajas: Mejora la precisión de clasificadores débiles.
- Desventajas: Sensible al ruido y a los datos atípicos.
- Mejores Parámetros:

```
1 {'n_estimators': 6, 'learning_rate': 0.6,
   'algorithm': 'SAMME.R'}
```

- Precisión: 0.6078066914498141

6. Análisis Discriminante

- Descripción: Clasifica utilizando combinaciones lineales de características.

- Ventajas: Rápido y eficiente para conjuntos de datos pequeños.
- Desventajas: No adecuado si no se cumplen las suposiciones de normalidad y homogeneidad de varianza.
- Mejores Parámetros:

```
1 {'solver': 'svd'}
```

- Precisión: 0.5548327137546468

El análisis realizado demuestra que la Red Neuronal es el modelo más efectivo para este conjunto de datos, logrando una precisión de 0.8596654275092936. Los resultados indican que la Red Neuronal no solo superó a los demás modelos, sino que también mostró una notable diferencia en precisión en comparación con AdaBoost, que fue el segundo mejor modelo con una precisión de 0.6078066914498141. Esta diferencia sugiere que la Red Neuronal es mucho más eficaz para capturar las complejidades inherentes en los datos.

8. Evaluación del Modelo (Model Evaluation)

Para determinar la confiabilidad de los modelos, se utilizarán las siguientes herramientas de evaluación:

1. **Métricas de Evaluación:** Para obtener una comprensión numérica del rendimiento de cada modelo.
2. **Matriz de Confusión:** Para visualizar las predicciones correctas e incorrectas.
3. **Curva de Aprendizaje:** Para analizar el rendimiento del modelo en función de la cantidad de datos de entrenamiento.

Se han definido funciones para facilitar la implementación de estas evaluaciones:

```
1 def plot_learning_curve(model, X, y):
2     train_sizes, train_scores, test_scores = learning_curve(
3         model, X, y, cv=5, n_jobs=-1, train_sizes=np.linspace
4         (0.1, 1.0, 10))
5     train_mean = np.mean(train_scores, axis=1)
6     train_std = np.std(train_scores, axis=1)
7     test_mean = np.mean(test_scores, axis=1)
8     test_std = np.std(test_scores, axis=1)
9
10    plt.plot(train_sizes, train_mean, 'o-', color='r', label='
11    Training score')
```



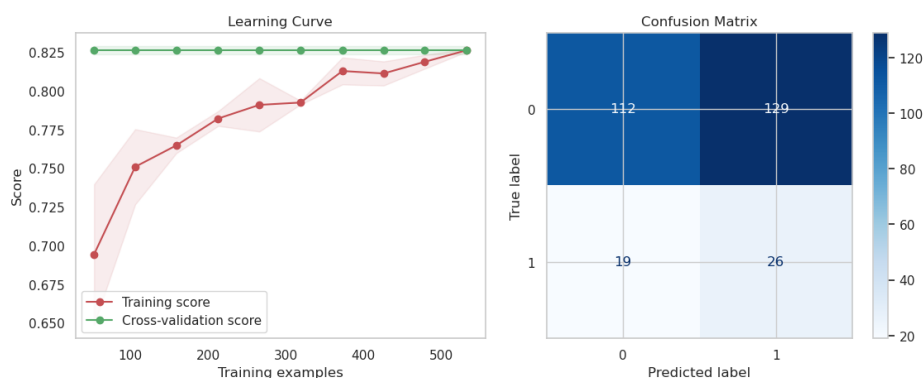
```

9     plt.plot(train_sizes, test_mean, 'o-', color='g', label='
      Cross-validation score')
10    plt.fill_between(train_sizes, train_mean - train_std,
      train_mean + train_std, alpha=0.1, color='r')
11    plt.fill_between(train_sizes, test_mean - test_std,
      test_mean + test_std, alpha=0.1, color='g')
12
13    plt.title('Learning Curve')
14    plt.xlabel('Training examples')
15    plt.ylabel('Score')
16    plt.legend(loc='best')
17    plt.grid()
18    plt.show()
19
20    def plot_confusion_matrix(model, X_test, y_test):
21        y_pred = model.predict(X_test)
22        cm = confusion_matrix(y_test, y_pred)
23        disp = ConfusionMatrixDisplay(confusion_matrix=cm)
24        disp.plot(cmap='Blues')
25        plt.title('Confusion Matrix')
26        plt.show()
27
28    def print_classification_report(model, X_test, y_test):
29        y_pred = model.predict(X_test)
30        report = classification_report(y_test, y_pred)
31        print("Classification Report:\n", report)

```

Listing 16: Evaluacion de modelos

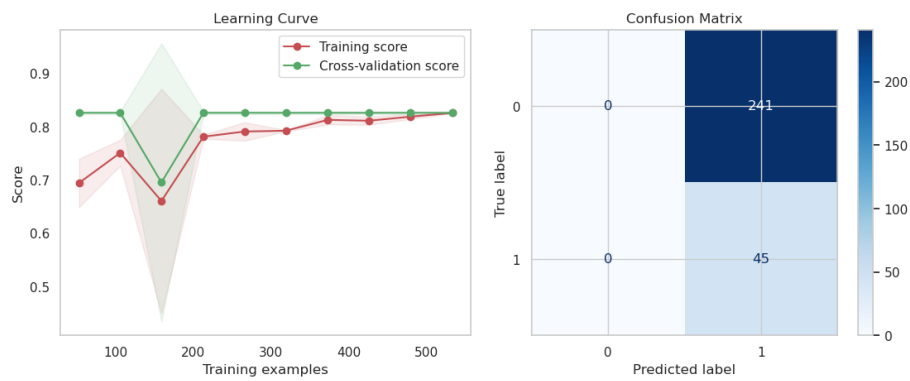
■ Support Vector Machine (SVM)



Classification Report:					
	precision	recall	f1-score	support	
0	0.85	0.46	0.60	241	
1	0.17	0.58	0.26	45	
accuracy				0.48	286

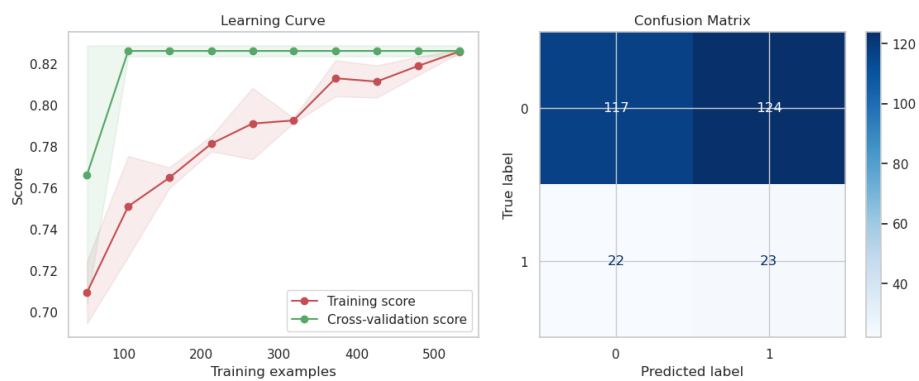
8	macro avg	0.51	0.52	0.43	286
9	weighted avg	0.75	0.48	0.55	286

■ Random Forest



1	Classification Report:				
2		precision	recall	f1-score	support
3					
4	0	0.00	0.00	0.00	241
5	1	0.16	1.00	0.27	45
6					
7	accuracy			0.16	286
8	macro avg	0.08	0.50	0.14	286
9	weighted avg	0.02	0.16	0.04	286

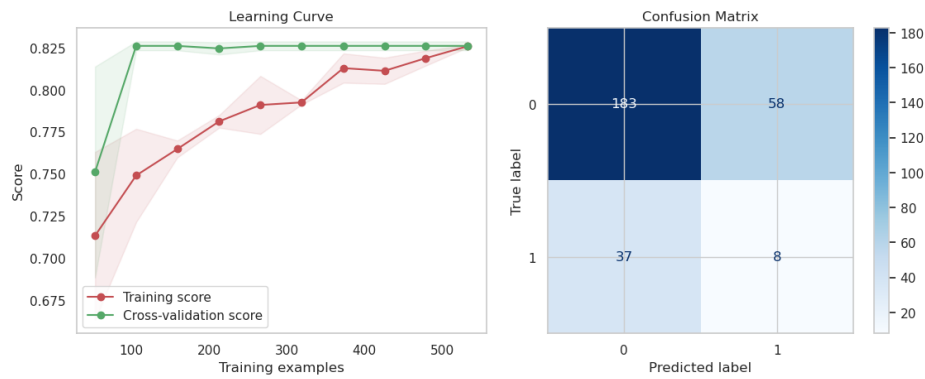
■ Regresión Logística



1	Classification Report:				
2		precision	recall	f1-score	support
3					
4	0	0.84	0.49	0.62	241
5	1	0.16	0.51	0.24	45

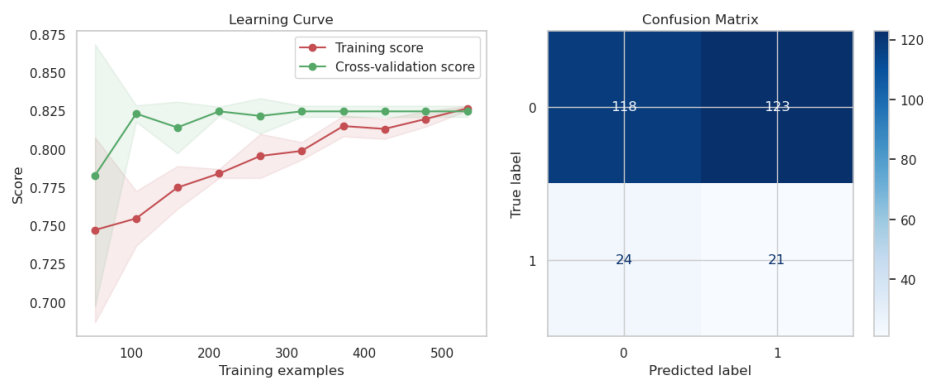
6					
7	accuracy			0.49	286
8	macro avg	0.50	0.50	0.43	286
9	weighted avg	0.73	0.49	0.56	286

■ Red Neuronal (Neural Network)



1	Classification Report:				
2		precision	recall	f1-score	support
3					
4	0	0.83	0.76	0.79	241
5	1	0.12	0.18	0.14	45
6					
7	accuracy			0.67	286
8	macro avg	0.48	0.47	0.47	286
9	weighted avg	0.72	0.67	0.69	286

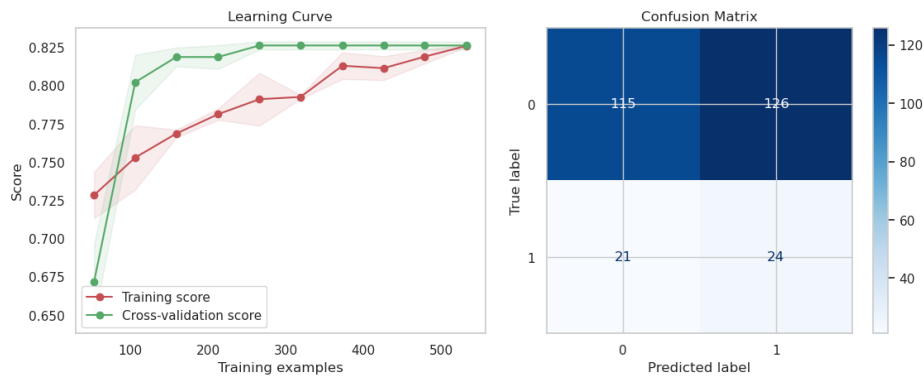
■ AdaBoost



1		Classification Report:			
2		precision	recall	f1-score	support
3					

4	0	0.83	0.49	0.62	241
5	1	0.15	0.47	0.22	45
6					
7	accuracy			0.49	286
8	macro avg	0.49	0.48	0.42	286
9	weighted avg	0.72	0.49	0.55	286

■ Análisis Discriminante



1	Classification Report:				
2		precision	recall	f1-score	support
3					
4	0	0.85	0.48	0.61	241
5	1	0.16	0.53	0.25	45
6					
7	accuracy			0.49	286
8	macro avg	0.50	0.51	0.43	286
9	weighted avg	0.74	0.49	0.55	286

Teniendo en cuenta los resultados obtenidos, se observa que el modelo de Red Neuronal es el que presenta la mayor precisión, alcanzando una precisión de 0.67. Además, es importante destacar que, aunque el modelo de AdaBoost tuvo un rendimiento aceptable con una precisión de 0.49, otros modelos como el SVM, Random Forest, Regresión Logística y Análisis Discriminante no superaron esta precisión.

9. Conclusiones

El tratamiento, análisis y modelación de datos es un proceso exhaustivo y complejo que requiere una serie de habilidades especializadas. Primero,

en el tratamiento de datos, es fundamental asegurar que los datos estén en condiciones óptimas. Los datos crudos a menudo no son útiles sin un proceso adecuado de limpieza y preparación, lo cual exige experiencia y el uso de diversas herramientas.

En el análisis de datos, es crucial tener una comprensión clara de los objetivos y necesidades del cliente o del proyecto. No se trata solo de generar gráficos, sino de plantear preguntas significativas y encontrar respuestas concretas utilizando los datos disponibles. Esta fase requiere una combinación de creatividad y análisis crítico para interpretar correctamente la información.

La modelación de datos es quizás la etapa más laboriosa y técnica. Implica no solo la creación de modelos predictivos, sino también la búsqueda de formas de maximizar el valor de los datos, ya sea mediante la creación de clusters, la reducción de dimensiones, o la selección de características adecuadas. Cada modelo tiene sus propias ventajas y desventajas:

- Red Neuronal: Alta precisión y capacidad para modelar relaciones no lineales complejas, aunque requiere muchos datos y es computacionalmente intensiva.
- Support Vector Machine (SVM): Eficaz en espacios de alta dimensión, pero no adecuado para conjuntos de datos muy grandes.
- Random Forest: Robusto ante el overfitting y maneja bien grandes conjuntos de datos, pero puede ser lento y consumir mucha memoria.
- Regresión Logística: Simple de implementar y fácil de interpretar, aunque no maneja bien relaciones no lineales.
- AdaBoost: Mejora la precisión de clasificadores débiles, aunque es sensible al ruido y a los datos atípicos.
- Análisis Discriminante: Rápido y eficiente para conjuntos de datos pequeños, pero no adecuado si no se cumplen ciertas suposiciones estadísticas.

En resumen, el trabajo de un científico de datos no se limita a la creación de modelos predictivos, sino que también incluye la preparación, análisis e interpretación de datos de manera efectiva. Este proyecto ha demostrado la importancia de cada una de estas etapas y ha subrayado la necesidad de un enfoque meticuloso para obtener resultados precisos y útiles.