

Proyecto final IA

David Santiago Orozco López

Inteligencia Artificial
Presentado A: Francisco Carlos Calderón
Bocanegra

Pontificia Universidad
Javeriana
2022



1. Tabla de contenido

1.	Tabla de contenido.....	2
1.	INTRODUCCIÓN	2
1.1	Planteamiento del problema.....	2
1.2	Definición del problema.....	2
1.2.1	Problema general	2
1.2.2	Especificación del problema.....	2
1.1	Objetivos	2
1.4.1	Objetivo general.....	2
2.	MARCO TEÓRICO	2
2.1	Matching Learning	2
2.2	¿Cómo funciona el Matching Learning?	2
2.3	Métodos de aprendizaje automático	2
2.4	Algoritmos comunes de aprendizaje automático.....	3
3.	ESTADO DEL ARTE	3
3.1	Casos de uso de aprendizaje automático en el mundo real.....	3
4.	PROPUESTA DE LA SOLUCION.....	4
4.1	Propuesta de la solución	4
4.2	Procedimiento de la solución.....	4
4.1	Resultados de la solución	5
5.	ANÁLISIS DE LOS RESULTADOS.	5
6.	CONCLUSIONES.....	5
7.	REFERENCIAS	5

1. INTRODUCCIÓN

1.1 Planteamiento del problema

La inteligencia artificial (IA), está cada vez más presente en nuestras vidas. Una buena definición de la misma sería la combinación de algoritmos, que intentan simular algunas acciones de los humanos o mejor aún, ir más allá de la inteligencia humana.

En este sentido, el siguiente informe va a presentar la aplicación de la inteligencia artificial por medio del uso del Matching Learning el cual busca desarrollar un clasificador de ciertas características de manera automática, a un conjunto de datos que se caracteriza por denominar ciertos parámetros a diferentes artistas musicales con una canción en específico, ejemplo: (Shakira-Waka Waka- Bailable 100%).

1.2 Definición del problema

1.2.1 Problema general

Como clasificar un dataset a partir de las características que tengan canciones musicales.

1.2.2 Especificación del problema

Más específicamente se usarán 3 clasificadores, los cuales son: regresión logística, SVN, KNN. Ya con estos clasificadores se encontrará el que mejor los clasifique por medio de parámetros como MCC y f1_score.

1.1 Objetivos

1.4.1 Objetivo general

1. Encontrar el mejor clasificador de los 3, para parametrizar ciertas canciones en diferentes categorías.

2. MARCO TEÓRICO

2.1 Matching Learning

El aprendizaje automático es una rama de la inteligencia artificial (IA) y la informática que se centra en el uso de datos y algoritmos para imitar la forma en que los humanos aprenden, mejorando gradualmente su precisión. [1]

2.2 ¿Cómo funciona el Matching Learning?

1. Un proceso de decisión: en general, los algoritmos de aprendizaje automático se utilizan para hacer una predicción o clasificación. En función de algunos datos de entrada, que se pueden etiquetar o no, su algoritmo producirá una estimación sobre un patrón en los datos.
2. Una función de error: una función de error evalúa la predicción del modelo. Si hay ejemplos conocidos, una función de error puede hacer una comparación para evaluar la precisión del modelo.
3. Un proceso de optimización del modelo: si el modelo puede ajustarse mejor a los puntos de datos en el conjunto de entrenamiento, entonces los pesos se ajustan para reducir la discrepancia entre el ejemplo conocido y la estimación del modelo. El algoritmo repetirá este proceso de “evaluar y optimizar”, actualizando los pesos de forma autónoma hasta alcanzar un umbral de precisión.

2.3 Métodos de aprendizaje automático

- Aprendizaje automático supervisado.
- Aprendizaje automático no supervisado.
- Aprendizaje semisupervisado.

2.4 Algoritmos comunes de aprendizaje automático

- **Neural networks:** las redes neuronales simulan la forma en que funciona el cerebro humano, con una gran cantidad de nodos de procesamiento vinculados. Las redes neuronales son buenas para reconocer patrones y juegan un papel importante en aplicaciones que incluyen traducción de lenguaje natural, reconocimiento de imágenes, reconocimiento de voz y creación de imágenes.
- **Linear regression:** este algoritmo se utiliza para predecir valores numéricos, basándose en una relación lineal entre diferentes valores. Por ejemplo, la técnica podría usarse para predecir los precios de la vivienda en función de los datos históricos del área.
- **Logistic regression:** este algoritmo de aprendizaje supervisado hace predicciones para variables de respuesta categóricas, como respuestas de "sí/no" a las preguntas. Se puede utilizar para aplicaciones como la clasificación de spam y el control de calidad en una línea de producción.
- **Clustering:** mediante el aprendizaje no supervisado, los algoritmos de agrupamiento pueden identificar patrones en los datos para poder agruparlos. Las computadoras pueden ayudar a los científicos de datos al identificar las diferencias entre los elementos de datos que los humanos han pasado por alto.
- **Decision trees:** los árboles de decisión se pueden usar tanto para predecir valores numéricos (regresión) como para clasificar datos en categorías. Los árboles de decisión utilizan una secuencia de ramificación de decisiones vinculadas que se pueden representar con un diagrama de árbol. Una de las ventajas de los árboles de decisión es que son

fáciles de validar y auditar, a diferencia de la caja negra de la red neuronal.

- **Random forests:** en un bosque aleatorio, el algoritmo de aprendizaje automático predice un valor o categoría al combinar los resultados de varios árboles de decisión.

3. ESTADO DEL ARTE

3.1 Casos de uso de aprendizaje automático en el mundo real.

Reconocimiento de voz: también se conoce como reconocimiento de voz automático (ASR), reconocimiento de voz por computadora o conversión de voz a texto, y es una capacidad que utiliza el procesamiento de lenguaje natural (NLP) para traducir el habla humana a un formato escrito. Muchos dispositivos móviles incorporan reconocimiento de voz en sus sistemas para realizar búsquedas por voz, por ejemplo, Siri, o mejorar la accesibilidad para enviar mensajes de texto.

Servicio al cliente: Los chatbots en línea están reemplazando a los agentes humanos a lo largo del viaje del cliente, cambiando la forma en que pensamos sobre la participación del cliente en los sitios web y las plataformas de redes sociales. Los chatbots responden preguntas frecuentes (FAQ) sobre temas como el envío, o brindan asesoramiento personalizado, venta cruzada de productos o sugerencias de tallas para los usuarios. Los ejemplos incluyen [agentes virtuales](#) en sitios de comercio electrónico; bots de mensajería, utilizando Slack y Facebook Messenger; y tareas que suelen realizar los asistentes virtuales y los asistentes de voz.

Visión por computadora: esta tecnología de inteligencia artificial permite que las computadoras obtengan información significativa de imágenes digitales, videos y otras entradas visuales, y luego tomen la acción apropiada. Impulsada por redes

artistas musicales con sus respectivas canciones.

4.2 Procedimiento de la solución

Declaración de bibliotecas para usar cada uno de los clasificadores:

```

from sklearn.model_selection import train_test_split|r\n",
"import pandas as pd|r\n",
"import numpy as np|r\n",
"import matplotlib.pyplot as plt|r\n",
"import csv|r\n",
"from sklearn import metrics|r\n",
"from sklearn.decomposition import PCA|r\n",
"|r\n",
"from sklearn.neighbors import KNeighborsClassifier|r\n",
"from sklearn.svm import SVC|r\n",
"from sklearn.linear_model import LogisticRegression|r\n",
"|r\n",
"from sklearn.preprocessing import minmax_scale|r\n",
"from sklearn.preprocessing import StandardScaler|r\n",
"|r\n",
"from sklearn.metrics import classification_report|r\n",
"from sklearn.metrics import confusion_matrix|r\n",
"from sklearn.metrics import roc_auc_score +***AUC-ROC***|r\n",
"from sklearn.metrics import MatthewsCorrcoef +***MCC***|r\n",
"from sklearn.metrics import f1_score +***F1***|r\n",
"from sklearn.metrics import plot_roc_curve"

```

Lectura del dataset y su respectiva limpieza con PCA la cual se utiliza para describir un conjunto de datos en términos de nuevas variables no correlacionadas.

```

"dataset = pd.read_csv('train.csv') #Se extraen los datos \r\n",
"\r\n",
"#Procesamiento de los datos\r\n",
"dataset = dataset.dropna() #Elimina las filas con valores NaN, que según Kaggle es de un porcentaje de 2% \r\n",
"dataset = dataset.drop(['ArtistName', axis=1]) #Elimina la columna del nombre del artista\r\n",
"dataset = dataset.drop(['Track Name', axis=1]) #Elimina la columna del nombre de la canción\r\n",
"\r\n",
"x = dataset.drop(['Class'], axis = 1).values\r\n",
"y = dataset['Class'].values\r\n",
"\r\n",
"scaler = StandardScaler()\r\n",
"scaler.fit(x)\r\n",
"x = scaler.fit_transform(x)\r\n",
"\r\n",
"#PCA\r\n",
"pca = PCA(n_components=10, svd_solver='full')\r\n",
"pca.fit(x)\r\n",
"print(pca.explained_variance_ratio_)\r\n",
"\r\n",
"x_train, x_test, y_train, y_test = train_test_split(x, y, random_state = 0, test_size = 0.2)\r\n",
"\r\n",
"x_train = pca.fit_transform(x_train)\r\n",
"y_train = y_train\r\n",
"x_test = pca.fit_transform(x_test)\r\n",
"y_test = y_test\r\n",

```

Primer método de clasificación regresión logística.

[illegible]

Diseñar un software en Python que permita la clasificación de manera automática de un dataset que es basado en la calificación de características, por medio de escalas (numéricas y de porcentaje) a diferentes

Segundo método de clasificación SVM.

```
<code># SVM
# Importar librerías
from sklearn.svm import SVC
from sklearn.metrics import accuracy_score, f1_score
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler

# Cargar datos
X = data[['x1', 'x2', 'x3', 'x4', 'x5', 'x6', 'x7', 'x8', 'x9', 'x10', 'x11', 'x12', 'x13', 'x14', 'x15', 'x16', 'x17', 'x18', 'x19', 'x20', 'x21', 'x22', 'x23', 'x24', 'x25', 'x26', 'x27', 'x28', 'x29', 'x30', 'x31', 'x32', 'x33', 'x34', 'x35', 'x36', 'x37', 'x38', 'x39', 'x40', 'x41', 'x42', 'x43', 'x44', 'x45', 'x46', 'x47', 'x48', 'x49', 'x50', 'x51', 'x52', 'x53', 'x54', 'x55', 'x56', 'x57', 'x58', 'x59', 'x60', 'x61', 'x62', 'x63', 'x64', 'x65', 'x66', 'x67', 'x68', 'x69', 'x70', 'x71', 'x72', 'x73', 'x74', 'x75', 'x76', 'x77', 'x78', 'x79', 'x80', 'x81', 'x82', 'x83', 'x84', 'x85', 'x86', 'x87', 'x88', 'x89', 'x90', 'x91', 'x92', 'x93', 'x94', 'x95', 'x96', 'x97', 'x98', 'x99', 'x100']]
y = data['target']

# Dividir datos en entrenamiento y prueba
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Escalar datos
scaler = StandardScaler()
X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test)

# Entrenar modelo SVM
svm = SVC(kernel='rbf')
svm.fit(X_train, y_train)

# Predecir
y_pred = svm.predict(X_test)

# Calcular métricas
accuracy = accuracy_score(y_test, y_pred)
f1 = f1_score(y_test, y_pred)

# Imprimir resultados
print("MCC Logistic Regression: ", mcc)
print("F1 score Logistic Regression: ", f1)</code>
```

Tercer método de clasificación KNN.

```
<code># KNN
# Importar librerías
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import accuracy_score, f1_score
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler

# Cargar datos
X = data[['x1', 'x2', 'x3', 'x4', 'x5', 'x6', 'x7', 'x8', 'x9', 'x10', 'x11', 'x12', 'x13', 'x14', 'x15', 'x16', 'x17', 'x18', 'x19', 'x20', 'x21', 'x22', 'x23', 'x24', 'x25', 'x26', 'x27', 'x28', 'x29', 'x30', 'x31', 'x32', 'x33', 'x34', 'x35', 'x36', 'x37', 'x38', 'x39', 'x40', 'x41', 'x42', 'x43', 'x44', 'x45', 'x46', 'x47', 'x48', 'x49', 'x50', 'x51', 'x52', 'x53', 'x54', 'x55', 'x56', 'x57', 'x58', 'x59', 'x60', 'x61', 'x62', 'x63', 'x64', 'x65', 'x66', 'x67', 'x68', 'x69', 'x70', 'x71', 'x72', 'x73', 'x74', 'x75', 'x76', 'x77', 'x78', 'x79', 'x80', 'x81', 'x82', 'x83', 'x84', 'x85', 'x86', 'x87', 'x88', 'x89', 'x90', 'x91', 'x92', 'x93', 'x94', 'x95', 'x96', 'x97', 'x98', 'x99', 'x100']]
y = data['target']

# Dividir datos en entrenamiento y prueba
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Escalar datos
scaler = StandardScaler()
X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test)

# Entrenar modelo KNN
knn = KNeighborsClassifier(n_neighbors=5)
knn.fit(X_train, y_train)

# Predecir
y_pred = knn.predict(X_test)

# Calcular métricas
accuracy = accuracy_score(y_test, y_pred)
f1 = f1_score(y_test, y_pred)

# Imprimir resultados
print("MCC Logistic Regression: ", mcc)
print("F1 score Logistic Regression: ", f1)</code>
```

4.1 Resultados de la solución

Parámetros de evaluación de la regresión logística.

```
<code>{
  "text": [
    "MCC Logistic Regression: 0.2700127837497189\n",
    "F1 score Logistic Regression: 0.4028776978417266\n"
  ]
}</code>
```

Parámetros de evaluación de la SVM

```
<code>{
  "text": [
    "MCC Logistic Regression: 0.2505451525138137\n",
    "F1 score Logistic Regression: 0.39272111722386804\n"
  ]
}</code>
```

Parámetros de evaluación de la KNN.

```
<code>{
  "text": [
    "MCC Logistic Regression: 0.2827439179658535\n",
    "F1 score Logistic Regression: 0.4168429961912823\n"
  ]
}</code>
```

5. ANÁLISIS DE LOS RESULTADOS.

A partir de los resultados mostrados anteriormente se puede analizar que los datos están demasiado dispersos entre si ya que muestra una correlación en cada una de ellas bastante baja, esto se debe principalmente a que las variables que se están, son subjetivas y se concentran en todo los rangos posibles, es decir que al ser tantas variables medidas, cada decil hablando en términos estadísticos contiene por lo menos 1800 datos, por lo cual la correlación de cada uno de los datos va a ser bastante grande. También si se habla del

valor F vemos que es más grande en todos los casos con respecto a la correlación esto quiere decir que la precisión y la recuperación del modelo es más mayor, con respecto a la correlación de los datos.

6. CONCLUSIONES.

Con este proyecto se puede llegar a la conclusión de que el mejor clasificador que hay para este caso es KNN teniendo un 0.29 de correlación de datos y 0.42 de F1, sin embargo, estos números son bastante bajos, lo que quiere decir que es un modelo bastante malo y que al ser tan subjetivo para cada valoración no se puede clasificar de mejor manera.

7. REFERENCIAS

[1]

<https://www.ibm.com/cloud/learn/machine-learning>