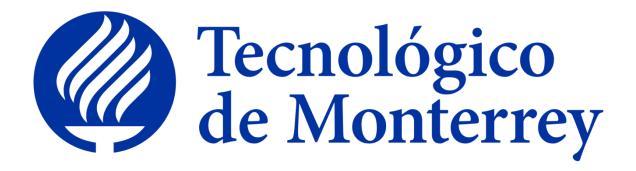
Instituto Tecnológico de Estudios Superiores de Monterrey



Inteligencia artificial avanzada para la ciencia de datos I (Gpo 101)

Profesores:

Jorge Adolfo Uresti

Actividad:Momento de Retroalimentación: Módulo 2 Implementación de una técnica de aprendizaje máquina sin el uso de un framework.

Eric Manuel Navarro Martínez A01746219

1 de Septiembre de 2024

Comenzamos con las siguientes funciones de apoyo para que el código pueda funcionar de forma adecuada:

normalization

 Se encarga de normalizar los datos basándonos en la Z score, tome en cuenta esta normalización debido a que decrementa el impacto de los outliers en los datos esto ya que mi dataset contiene outliers considerables debido a las calificaciones que puede que aunque bajas hayan tenido un gran desempeño en taquilla y viceversa. (Effects Of Normalization Techniques On Logistic Regression In Data Science, 2022)

```
def normalization(data):
    #This normalization utilizes Z-score
    mean = np.mean(data, axis=0)
    std = np.std(data, axis=0)
    return (data-mean)/std
```

sigmoid

 \circ

- Calcula la función sigmoidal, al igual que tiene unos stops que permite que no se haga un overflow al usar datos extremadamente grandes evita igualmente dividir entre 0
- r squared
 - Calcula la r cuadrada usando la suma de cuadrados al igual
- confusion matrix
 - Realiza los cálculos de acuracy, specifity, precision, recall, f1 score a través del cálculo de los true positives, true negatives, false positives y false negatives que puedan ser obtenidos de una predicción errónea o certera de una buena o mala película, al igual que tiene un rango de error.

```
confusion_matrix(y_true, y_pred, threshold=0.5):
FN = 0 #false negative
y_true = np.array(y_true)
y_pred = np.array(y_pred)
predicted = (v pred >= threshold).astype(int)
#Check if the prediction is corre
TP = np.sum((predicted == 1) & (y_true == 1))
TN = np.sum((predicted == 0) & (y_true == 0))
FP = np.sum((predicted == 1) & (y_true == 0))
FN = np.sum((predicted == 0) & (y_true == 1))
 \textbf{accuracy = (TP+TN)/(TP+TN+FP+FN) if (TP+TN+FP+FN) > 0 else 0} \\ \textbf{ #avoid dividing by zero and crashing the program } 
specifity = (TN)/(TN+FP) if (TN+FP) > 0 else 0 #avoid dividing by zero and crashing the program
precision = (TP)/(TP+FP) if (TP+FP) > 0 else 0 #avoid dividing by zero and crashing the program
F1 = 2*precision*recall/(precision+recall) if (precision+recall) > 0 else 0 #avoid dividing by zero and crashing the program
return {'True Positives: ': TP, 'True Negatives: ': TN,
          'False Positives:': FP, 'False Negatives: ': FN, 'Accuracy': accuracy, 'Specifity': specifity,
              'Precision': precision, 'F1 Score': F1,
    'Recall: ': recall}
```

linear

 Realizamos a la muestra y parámetros un producto punto, evaluando una función lineal por x elementos de las muestras y de los parámetros

```
def linear(params, sample):
return np.dot(params, sample) # Apply dot function to the parameters and the sample
for the parameters and the sample
```

- binary cross entropy
 - Permite que se calcule el error medio al momento de usar el logistical regression al igual que tien un stop a números demasiado pequeños para que no se realicen operaciones como lo es el log(0)

```
def binary_cross_entropy(y_true, y_pred):
    y_pred = np.clip(y_pred, 1e-15, 1 - 1e-15) # Prevent log(0) issues
    return -np.mean(y_true * np.log(y_pred) + (1 - y_true) * np.log(1 - y_pred))
```

- gradient descent
 - En la función de gradient descent hacemos diversos calculos los cuales nos permiten en la longitud de los parametros que enviamos y con un alfa el cual nos da pequeños saltos para estar mejorando de forma más precisa el mismo, estas predicciones toman la sigmoide nuevamente para encontrar el error predecido, y actualizarse continuamente.

```
def gradient_descent(params, samples, y, alfa):
    adjusted = np.copy(params) # Doesn't modify the original array
    for i in range (len(params)):
        c = 0
        for j in range(len(samples)):
            prediction = sigmoid(linear(params, samples[j]))
            error = prediction - y[j] #Compute the prediction error
        c += error*samples[j][i] #C gets updated based on the error and the sample
        adjusted[i] = params[i]-alfa*(1/len(samples))*c #Update the parameter
        return adjusted
```

- logistic regression train:
 - Entrenamos con el dataset de training para poder obtener los parámetros adecuados que sean capaces de crear un modelo que pueda predecir de la manera más precisa las películas y su calidad, para hacer esto tomamos los valores de nuestro dataset de training, el cual vamos a normalizar para poder crear mejores predicciones en el y train el cual es la columna que queremos predecir.
 - Esta función entonces usa el gradient descent en conjunto con la evaluación lineal a la cual le realizaremos su sigmoide para que sea logística, por último tomamos en cuenta el error el cual permite que la función esté activa siempre y cuando el error sea mayor a 0.000001.
 - Aquí es donde se ven reflejadas las funciones ya que estará mejorando continuamente, por lo que guardamos el error actual y en qué época nos encontramos hasta que lleguemos al error menor; en este momento se considera el entrenamiento terminado y los parámetros obtenidos

```
def logistic_regression_train(x_train, y_train, alfa):
         # Normalizing the data
         normalized_x_train = normalization(x_train)
          # Initial parameters
          params = np.zeros(x_train.shape[1]) #Set the parameters to an array of zeros
          previous error = float('inf') #Set the previous error to infinity
          while True:
             # Update the parameters using gradient descent
 85
              params = gradient_descent(params, normalized_x_train, y_train, alfa)
             y_pred_prob = [sigmoid(linear(params, sample)) for sample in normalized_x_train]
              #Calculate the error
              current_error = binary_cross_entropy(y_train, y_pred_prob)
              epochs += 1
              print(f'Epoch: {epochs}')
              print('Parameters:', params)
              print('Current error:', current_error)
              print('Previous error:', previous_error)
              print('-----
              if abs(previous_error - current_error) < 0.000001:</pre>
                  print("Finished training")
                  break
              previous_error = current_error
106
          return params
```

logistic_regression_valid

 A diferencia del logistic regression de training esta vez si queremos calcular la matriz de errores ya que podemos hacer la comparativa una vez que recibimos los params de el logistic regression de training, este despliega estos datos para una vez que se termine la comparación se pueda hacer un análisis del mismo

probabilities

 En esta función imprimiremos todas las probabilidades de nuestro dataset de que sean una buena o mala pelicula, esto usando los parametros obtenidos de nuestro modelo entrenado

```
#Compare the model with the test set AKA validation dataset

def probabilities(x_test, params):

# Normaliza los datos de prueba
normalized_x_test = normalization(x_test)

# Calcula las probabilidades para el conjunto de prueba
y_pred_prob = [sigmoid(linear(params, sample)) for sample in normalized_x_test]

# Imprime todas las probabilidades junto con su interpretación
for i, prob in enumerate(y_pred_prob):

print(f"Title {testing_data['Title'].iloc[i+1]}: Probability = {prob:.4f}")

if prob >= 0.5:

print("Expected a good movie with this probability: ", prob)
else:

print("Expect a bad movie with this probability: ", prob)
print('------')

return y_pred_prob
```

Tomaremos el siguiente dataset de IMDB, el cual lo limpie previamente para limitar la cantidad de columnas y datos que recibe para predecir, manteniendo solamente, el título el cual no lo usaremos en ningún momento de la predicción, rating, votes y revenue en millions, el dato a predecir es la metascore de la misma.

Están divididos en validation, testing y training para que pueda analizarlo de forma eficaz usando pandas, con training midiendo 251 datos, validation 30 y 30 para test. Al momento de analizar con respecto a solamente training toma 3785 epochs en contrar un parametro adecuado al tener un error menor al 0.000001

Posteriormente guardamos estos parametros para realizar una corrida con el dataset de validation el cual nos despliega los siguientes valores de matriz de confusión.

True Positives: 17, True Negatives: 11, False Positives: 0. False Negatives: 3,

Accuracy: 0.9032258064516129,

Specifity: 1.0, Precision: 1.0,

F1 Score: 0.9189189189189

Recall: : 0.85

Por último para comparar las probabilidades que lanza el modelo usando el test, podemos ver las siguientes probabilidades para cada caso, con varias peliculas las cuales el modelo considerará malas o buenas.

```
Title Stardust:

Expected a good movie with this probability: 0.8037309536424588

Title American Hustle:

Expected a good movie with this probability: 0.6068540956275106

Title Jennifer's Body:

Expect a bad movie with this probability: 0.04769989847721063

Title Midnight in Paris:

Expected a good movie with this probability: 0.7656455969949044
```

En conclusión podemos decir que el modelo realiza su cometido ya que tiene probabilidades vistas para películas que no conozca haciendo un análisis, que si es una buena o mala película.

```
Sample 1: Probability = 0.2867
Expect a bad movie with this probability: 0.28668046935816976
Sample 2: Probability = 0.6694
Expected a good movie with this probability: 0.6694057962504452
Sample 3: Probability = 0.3996
Expect a bad movie with this probability: 0.3996397562533787
Sample 4: Probability = 0.3938
Expect a bad movie with this probability: 0.39380684995698306
Sample 5: Probability = 0.6663
Expected a good movie with this probability: 0.6662897755501807
Sample 6: Probability = 0.3453
Expect a bad movie with this probability: 0.3453457568546494
_____
Sample 7: Probability = 0.6813
Expected a good movie with this probability: 0.6813459947091449
Sample 8: Probability = 0.2152
Expect a bad movie with this probability: 0.21516767839700324
Sample 9: Probability = 0.2247
Expect a bad movie with this probability: 0.22471293804442397
Sample 10: Probability = 0.7926
Expected a good movie with this probability: 0.7926396873731616
_____
Sample 11: Probability = 0.6221
Expected a good movie with this probability: 0.6220710804509563
_____
Sample 12: Probability = 0.0886
Expect a bad movie with this probability: 0.08859725967713639
```

Sample 13: Probability = 0.6172 Expected a good movie with this probability: 0.6171973879352592 Sample 14: Probability = 0.6796 Expected a good movie with this probability: 0.679576435317383 _____ Sample 15: Probability = 0.8328 Expected a good movie with this probability: 0.8328496036933907 _____ Sample 16: Probability = 0.7473 Expected a good movie with this probability: 0.7472800729560392 Sample 17: Probability = 0.8604 Expected a good movie with this probability: 0.8603736446466942 Sample 18: Probability = 0.8282 Expected a good movie with this probability: 0.828195829069278 Sample 19: Probability = 0.3952 Expect a bad movie with this probability: 0.3951972534731864 Sample 20: Probability = 0.3972 Expect a bad movie with this probability: 0.3972297160582582 Sample 21: Probability = 0.5806 Expected a good movie with this probability: 0.5805729453495858 Sample 22: Probability = 0.3163 Expect a bad movie with this probability: 0.3163095651914222 Sample 23: Probability = 0.5736 Expected a good movie with this probability: 0.5736440103839169 Sample 24: Probability = 0.2010 Expect a bad movie with this probability: 0.20096736628277476

Sample 25: Probability = 0.1523

Expect a bad movie with this probability: 0.1523350839045311

Sample 26: Probability = 0.8037

Expected a good movie with this probability: 0.8037309536424588

Sample 27: Probability = 0.6069

Expected a good movie with this probability: 0.6068540956275106

Sample 28: Probability = 0.0477

Expect a bad movie with this probability: 0.04769989847721063

Sample 29: Probability = 0.7656

Expected a good movie with this probability: 0.7656455969949044