

Machine Learning - Lab 4 AdaBoost

RUBEN MARTINEZ GONZALEZ

Marzo 2024

1 Introducción

En este informe, se presenta el trabajo realizado en el laboratorio 4 de Machine Learning. El objetivo principal es implementar el algoritmo AdaBoost con clasificación binaria para demostrar cómo se puede mejorar la precisión de los clasificadores débiles al combinarlos en un modelo fuerte. Para entrenar el modelo que representa el algoritmo, se utilizarán los datos del archivo `dataCircle.txt`. Como complemento, se visualizarán las líneas de decisión de cada clasificador para una mejor comprensión de su funcionamiento.

La solución se implementó en Python utilizando la biblioteca NumPy para el procesamiento de matrices y Matplotlib para la visualización de gráficos.

La implementación está desplegada en un notebook de Google Colab, el cual se puede acceder a través del siguiente enlace: [Colab](#)

2 Implementación

2.1 Procesamiento de datos

Para el procesamiento de datos, se cargan los datos del archivo de texto `dataCircle.txt` y se extraen las primeras dos columnas como las coordenadas (x, y). Luego, se separan los ejemplos positivos y negativos siendo los primeros 40 ejemplos positivos y los últimos a partir del 41 ejemplos negativos. Finalmente, se grafican los ejemplos positivos en verde y los ejemplos negativos en rojo.

```
1 # Load the data from the text file
2 data = np.loadtxt('dataCircle.txt')
3
4 # Extract the first two columns as the x and y coordinates
5 x = data[:, 0]
6 y = data[:, 1]
7
8 # Separate the positive and negative examples
9 positive_examples = data[:40, :]
10 negative_examples = data[40:, :]
11
12 # Plot the positive examples in green
13 plt.scatter(positive_examples[:, 0], positive_examples[:, 1], c='g', label='
    Positive examples')
14
15 # Plot the negative examples in red
16 plt.scatter(negative_examples[:, 0], negative_examples[:, 1], c='r', label='
    Negative examples')
17
18 # Add labels and title
19 plt.xlabel('X')
20 plt.ylabel('Y')
21 plt.title('Data Plot')
22
```

```

23 # Show the plot
24 plt.legend()
25 plt.show()
26

```

Listing 1: Data processing

Al ejecutar este bloque de código, se obtiene la siguiente gráfica:

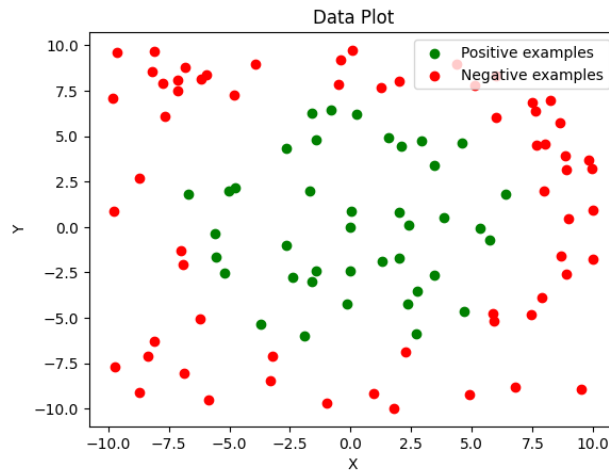


Figure 1: Gráfica de los datos

2.2 Implementación de AdaBoost

Para la implementación del algoritmo, se diseñó la clase `AdaBoost` con 3 métodos principales: `train`, `predict` y `weak_classifier`, y un método auxiliar `error`. Con estos métodos, es posible entrenar un modelo para clasificar nuevos datos basados en el conocimiento adquirido de los datos de entrenamiento. La implementación de la clase `AdaBoost` se describe a continuación:

- Clase `AdaBoost`:

Esta clase se inicializa recibiendo por parámetro el número de clasificadores a entrenar, y además, se inicializan las listas `classifiers` y `alphas` que almacenarán los clasificadores débiles y los valores de α respectivamente. La clase también define los métodos `train`, `predict`, `weak_classifier` y `compute_error` descritos posteriormente.

```

1 class AdaBoost:
2     def __init__(self, n_estimators):
3         self.n_estimators = n_estimators #cantidad de clasificadores
4         self.alphas = []
5         self.classifiers = []
6
7     def train(self, X, y):
8
9     def predict(self, X):
10
11     def weak_classifier(self, X, y, D):
12
13     def compute_error(self, y, y_pred):
14

```

Listing 2: AdaBoost class

- Método train

El metodo train recibe como parámetros los datos de entrenamiento X y las etiquetas y. Inicializa los pesos de las muestras con 1/total para que todas las muestras igual peso y la suma quede normalizada. Luego, se itera sobre el número de clasificadores débiles a entrenar y para cada clasificador:

- Se encuentra el clasificador débil que minimiza el error invoca el método `weak_classifier` descrito posteriormente.
- Se calcula el error del clasificador actual invocando el método `compute_error` descrito posteriormente.
- Se calcula el valor de α del clasificador actual mediante la fórmula:

$$\alpha = 0.5 * \log((1 - error)/error).$$
 Este valor de α representará la contribución del clasificador débil al modelo final.
- Se agrega el clasificador actual a la lista de clasificadores y el valor de α a la lista de alphas.
- Se obtienen las predicciones del clasificador actual accediendo a la propiedad `predictions` del clasificador.
- Se actualizan los pesos de las muestras para que las mal clasificadas tengan más peso en las siguientes iteraciones. Para actualizar los pesos se almacena en la variable `weights` el resultado de multiplicar los pesos por $e^{-\alpha*y*predictions}$ y se normalizan dividiendo por la suma de los pesos.

```

1 def train(self, X, y):
2     # Inicializaci n de pesos
3     n_samples, n_features = X.shape
4
5     # Se inicializan los pesos de las muestras con 1/total
6     weights = np.ones(n_samples) / n_samples
7
8     for _ in range(self.n_estimators):
9
10        # encuentra el clasificador d bil que minimiza el error
11        classifier = self._weak_classifier(X, y, weights)
12
13        # calcula el error del clasificador actual
14        error = self._compute_error(classifier, X, y, weights)
15
16        # calcula alpha del clasificador actual (contribuci n al modelo final
17        )
18        alpha = 0.5 * np.log((1 - error) / error)
19
20        # agrega el clasificador actual a la lista de clasificadores
21        self.classifiers.append(classifier)
22
23        # agrega el alpha del clasificador actual a la lista de alphas.
24        self.alphas.append(alpha)
25
26        # Se obtienen las predicciones del clasificador actual.
27        predictions = classifier['predictions']
28
29        #Se actualizan los pesos de las muestras
30        #las muestras mal clasificadas tengan m s peso en las siguientes
31        iteraciones
32        weights *= np.exp(-alpha * y * predictions)
33
34        #Se normalizan los pesos
35        weights /= np.sum(weights)

```

Listing 3: train method

- Método predict

El método predict recibe como parámetro las muestras X y retorna las predicciones finales del modelo. Para predecir, se inicializa un vector de ceros que contendrá las predicciones finales, y luego se itera sobre cada clasificador y su correspondiente α . Para cada clasificador, se suman las predicciones acumuladas multiplicadas por el α del clasificador. Finalmente, si el valor es positivo, se clasifica como 1, si es negativo como -1.

```

1 def predict(self, X):
2     n_samples = X.shape[0]
3
4     #inicializa un vector de ceros que contendr las predicciones finales
5     predictions = np.zeros(n_samples)
6
7     #Se itera sobre cada clasificador y su correspondiente alpha
8     for alpha, classifier in zip(self.alphas, self.classifiers):
9         #Para cada clasificador: predicciones acumuladas = predicciones * peso
10        predictions += alpha * classifier['predictions']
11    #Si el valor es positivo, se clasifica como 1, si es negativo como -1
12    return np.sign(predictions)
13

```

Listing 4: predict method

- Método weak_classifier

El método weak_classifier recibe como parámetros las muestras X, las etiquetas y y los pesos de las muestras. Su objetivo es encontrar el clasificador débil que minimiza el error ponderado. Para ello, realiza la siguiente secuencia:

- Extrae el número de muestras y características del conjunto de datos.
- Inicializa el mejor error con infinito y el mejor clasificador con None.
- Itera sobre cada característica del conjunto de datos y para cada característica:
- Itera sobre cada valor de la característica actual como posible umbral y para cada umbral:
 - * Inicializa un vector de predicciones en 1 del tamaño de las muestras.
 - * Actualiza las predicciones a -1 donde el valor es menor que el umbral.
 - * Calcula el error ponderado sumando los pesos de las muestras mal clasificadas.
 - * Si el error actual es menor que el mejor error, se actualiza el mejor error y el mejor clasificador.
- Retorna el mejor clasificador encontrado.

```

1 def _weak_classifier(self, X, y, weights):
2     n_samples, n_features = X.shape
3     best_error = np.inf
4     best_classifier = None
5     #Itera sobre cada caracter stica del conjunto de datos
6     for feature in range(n_features):
7         #Itera sobre cada valor de la caracter stica actual como posible
8         umbral.
9         for threshold in np.unique(X[:, feature]):
10
11             #inicializa vector de predicciones en 1.
12             predictions = np.ones(n_samples)
13
14             # actualiza las predicciones a -1 donde el valor es menor que el
15             umbral
16             predictions[X[:, feature] < threshold] = -1
17
18             #error ponderado sumando los pesos de las muestras mal
19             clasificadas

```

```

17         error = np.sum(weights[y != predictions])
18
19         #Si el error actual es menor, se actualiza el mejor error
20         if error < best_error:
21             best_error = error
22             best_classifier = {'feature': feature, 'threshold': threshold,
23                               'predictions': predictions.copy()}
24         return best_classifier

```

Listing 5: weak_classifier method

- Método error

El método `compute_error` recibe como parámetros el clasificador actual, las muestras `X`, las etiquetas `y`, y los pesos de las muestras. Su objetivo es calcular el error ponderado del clasificador actual. Para ello, simplemente suma los pesos de las muestras mal clasificadas.

```

1 def _compute_error(self, classifier, X, y, weights):
2     return np.sum(weights[y != classifier['predictions']])
3

```

Listing 6: compute_error method

3 Entrenamiento

Para entrenar el modelo usando la clase `AdaBoost`, se cargan los datos del archivo `dataCircle.txt` y se separan las características y las etiquetas. Luego, se inicializa un objeto `AdaBoost` con 10 clasificadores débiles y se entrena el modelo invocando el método `train`. Finalmente, se realizan predicciones sobre los datos de entrenamiento y se calcula la precisión del modelo. Para calcular la precisión, simplemente se compara las predicciones con las etiquetas y se calcula el promedio de las coincidencias. Al finalizar la ejecución, se imprime en pantalla la precisión del modelo.

```

1 # Load the data
2 data = np.loadtxt('dataCircle.txt')
3
4 # Separate features and labels
5 x = data[:, :2]
6 y = data[:, 2]
7
8 # Initialize AdaBoost with 10 weak classifiers
9 adaboost = AdaBoost(n_estimators=10)
10
11 # Train AdaBoost
12 adaboost.train(x, y)
13
14 # Make predictions
15 predictions = adaboost.predict(x)
16
17 # Calculate accuracy
18 accuracy = np.mean(predictions == y)
19 print("Accuracy:", accuracy)
20

```

Listing 7: Training the model

4 Resultados

El ejecutar el entrenamiento del modelo, se obtiene una precisión del 0.686 lo cual está cercano a un 70% de precisión. Para poder visualizar mejor el modelo, se grafican los datos de entrenamiento y las líneas de decisión de cada clasificador débil. Para ello, se grafican los datos de entrenamiento positivos en verde y los negativos en rojo. Luego, se itera sobre los clasificadores débiles y para cada clasificador:

- Se obtiene la característica y el umbral en el clasificador actual.
- Si la característica es 0, se grafica una línea vertical en el umbral, si no, se grafica una línea horizontal. (en ambos casos de color azul)

Finalmente, se muestra la gráfica con los datos de entrenamiento y las líneas de decisión de los clasificadores débiles.

```
1 # Plot the positive and negative examples
2 plt.scatter(positive_examples[:, 0], positive_examples[:, 1], c='g', label='
    Positive examples')
3 plt.scatter(negative_examples[:, 0], negative_examples[:, 1], c='r', label='
    Negative examples')
4
5 #itera sobre los clasificadores d biles
6 for classifier in adaboost.classifiers:
7 # Obtiene la caracter stica y el umbral en el clasificador actual.
8 feature = classifier['feature']
9 threshold = classifier['threshold']
10
11 # l neas de los clasificadores de color azul
12 if feature == 0:
13 plt.plot([threshold, threshold], [min(x[:, 1]), max(x[:, 1])], c='b', linestyle='
    --')
14 else:
15 plt.plot([min(x[:, 0]), max(x[:, 0])], [threshold, threshold], c='b', linestyle='
    --')
16
17 # Add labels and title
18 plt.xlabel('X')
19 plt.ylabel('Y')
20 plt.title('Data Plot with Weak Classifiers')
21
22 # Show the plot
23 plt.legend()
24 plt.show()
25
```

Listing 8: Plotting the decision boundaries

Al ejecutar este bloque de código, se obtiene la siguiente gráfica:

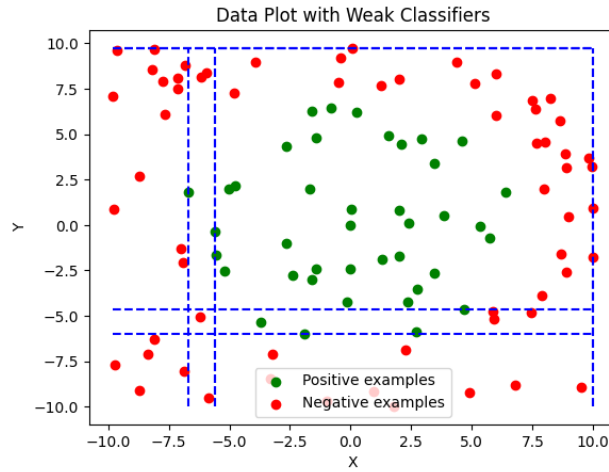


Figure 2: Gráfica de los datos con los clasificadores débiles

5 Conclusiones

En este laboratorio, se implementó el algoritmo AdaBoost para clasificación binaria. Se entrenó el modelo con 10 clasificadores para los datos del archivo `dataCircle.txt` y se obtuvo una precisión del 68.6%. Esto demuestra que el algoritmo AdaBoost es capaz de mejorar la precisión de los clasificadores débiles al combinarlos en un modelo fuerte. Además, se visualizaron las líneas de decisión de cada clasificador débil, lo que permite entender cómo se combinan para formar el modelo final.