

Instituto Tecnológico de Costa Rica - TEC

Inteligencia Artificial

Docente: Kenneth Obando Rodríguez

Trabajo Corto 3: Árboles de Decisión

Estudiantes:

- Renzo Giuliano Barra Mostajo
- Ana María Guevara Roselló
- Jonathan Alberto Guzmán Araya

Link del Cuaderno (recuerde configurar el acceso a público):

- [Link de su respuesta](#)

Nota: Este trabajo tiene como objetivo promover la comprensión de la materia y su importancia en la elección de algoritmos. Los alumnos deben evitar copiar y pegar directamente información de fuentes externas, y en su lugar, demostrar su propio análisis y comprensión.

Entrega

Debe entregar un archivo comprimido por el TecDigital, incluyendo un documento pdf con los resultados de los experimentos y pruebas. La fecha de entrega es el domingo 17 de setiembre, antes de las 10:00pm.

Instrucciones:

Las alternativas se rifarán en clase utilizando números aleatorios. Deberá realizar la asignación propuesta. Si realiza ambos ejercicios, recibirá 20 puntos en **la nota porcentual de la actividad**, para aplicar a la totalidad de los puntos extra es necesario que ambas actividades se completen al 100%

Actividad - Taller

1. Cree una clase nodo con atributos necesarios para un árbol de decisión: feature, umbral, gini, cantidad_muestras, valor, izquierda, derecha

```
1 import numpy as np
2 from collections import Counter
3 from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score
4
5
6 # Step 1: Node Class Creation
7 # This class represents a node in the decision tree.
8 class Node:
9     def __init__(self, feature=None, threshold=None, impurity=None, sample_count=None, value=None, left=None, right=None):
10         self.feature = feature
11         self.threshold = threshold
12         self.impurity = impurity
13         self.sample_count = sample_count
14         self.value = value
15         self.left = left
16         self.right = right

17
18
19 # Function to find the most common class in a list of labels
20 def most_common_class(y):
21     class_counts = Counter(y)
22     most_common = class_counts.most_common(1)[0][0]
23     return most_common
24
25
26 # Function to select the best feature and threshold for splitting
27 def find_best_split(X, y, criterion='gini'):
28     best_impurity = float('inf')
29     best_feature = None
30     best_threshold = None
31
32     for feature in range(X.shape[1]):
33         unique_thresholds = np.unique(X[:, feature])
34         for threshold in unique_thresholds:
35             left_indices = X[:, feature] <= threshold
36             right_indices = X[:, feature] > threshold
```

```

19         impurity = calculate_impurity(
20             y[left_indices], y[right_indices], criterion)
21
22         if impurity < best_impurity:
23             best_impurity = impurity
24             best_feature = feature
25             best_threshold = threshold
26
27     return best_feature, best_threshold
28
29
30 # Impurity Functions (Gini and Entropy)
31 # These functions calculate the impurity of a set of labels.
32 def calculate_impurity(y_left, y_right, criterion='gini'):
33     if criterion == 'gini':
34         impurity_left = gini_impurity(y_left)
35         impurity_right = gini_impurity(y_right)
36         total_samples = len(y_left) + len(y_right)
37         weighted_impurity = (len(y_left) / total_samples) * impurity_left + \
38             (len(y_right) / total_samples) * impurity_right
39         return weighted_impurity
40     elif criterion == 'entropy':
41         entropy_left = entropy_impurity(y_left)
42         entropy_right = entropy_impurity(y_right)
43         total_samples = len(y_left) + len(y_right)
44         weighted_entropy = (len(y_left) / total_samples) * entropy_left + \
45             (len(y_right) / total_samples) * entropy_right
46         return weighted_entropy
47     else:
48         raise ValueError(
49             "Invalid criterion. Supported criteria are 'gini' and 'entropy'.")
50
51
52 # Function to calculate the entropy of a set of labels
53 def entropy_impurity(labels):
54     num_samples = len(labels)
55     if num_samples == 0:
56         return 0.0
57
58     class_counts = Counter(labels)
59     impurity = 0.0
60     for class_count in class_counts.values():
61         class_probability = class_count / num_samples
62         impurity -= class_probability * np.log2(class_probability)
63
64     return impurity
65
66
67 # Function to calculate the Gini index of a set of labels
68 def gini_impurity(labels):
69     num_samples = len(labels)
70     if num_samples == 0:
71         return 0.0
72
73     class_counts = Counter(labels)
74     impurity = 1.0
75     for class_count in class_counts.values():
76         class_probability = class_count / num_samples
77         impurity -= class_probability ** 2
78
79     return impurity
80
81 # Function to split the dataset into left and right subsets
82 def split_dataset(X, y, feature, threshold):
83     left_indices = X[:, feature] <= threshold
84     right_indices = X[:, feature] > threshold
85
86     X_left = X[left_indices]
87     y_left = y[left_indices]
88
89     X_right = X[right_indices]
90     y_right = y[right_indices]
91
92     return X_left, y_left, X_right, y_right

```

2. Crea una clase que implementa un árbol de decisión, utilice las funciones presentadas en clase, además incluya los siguientes hiperparámetros:

- max_depth: Cantidad máxima de variables que se pueden explorar
- min_split_samples: Cantidad mínima de muestras que deberá tener un nodo para poder ser dividido
- criterio: función que se utilizará para calcular la impuridad.

```

1 # Step 2: Decision Tree Class Creation
2 class DecisionTree:
3     def __init__(self, max_depth=None, min_samples_split=2, criterion='gini'):
4         self.max_depth = max_depth
5         self.min_samples_split = min_samples_split
6         self.criterion = criterion
7         self.root = None
8
9     def train(self, X, y, depth=0):
10        # Check stopping criteria
11        if depth == self.max_depth or len(X) < self.min_samples_split:
12            # Create a leaf node with the majority class or the average value
13            # depending on the problem (classification or regression)
14            # Example for classification:
15            value = most_common_class(y)
16            return Node(value=value)
17
18        # Choose the best feature and threshold to split the dataset
19        feature, threshold = find_best_split(X, y, self.criterion)
20
21        # Split the dataset into left and right subsets
22        X_left, y_left, X_right, y_right = split_dataset(
23            X, y, feature, threshold)
24
25        # Recursively build the sub-trees
26        left = self.train(X_left, y_left, depth=depth + 1)
27        right = self.train(X_right, y_right, depth=depth + 1)
28
29        # Create and return a decision node
30        return Node(feature=feature, threshold=threshold, left=left, right=right)
31
32    def predict(self, X):
33        # Initialize an array to store the predictions
34        predictions = []
35
36        # Traverse the decision tree for each sample in X
37        for sample in X:
38            node = self.root
39            while node.left:
40                if sample[node.feature] <= node.threshold:
41                    node = node.left
42                else:
43                    node = node.right
44
45            # Append the predicted value for this sample
46            predictions.append(node.value)
47
48        return predictions
49

```

3. Divida los datos en los conjuntos tradicionales de entrenamiento y prueba, de forma manual, sin utilizar las utilidades de sklearn (puede utilizar índices de Numpy o Pandas)

```

1 # Step 3: Data Splitting
2 # This function splits the dataset into training and test sets.
3 def manual_train_test_split(X, y, train_proportion=0.8, random_state=None):
4     # Calculate the number of training samples
5     n_train_samples = int(train_proportion * len(X))
6
7     if random_state is not None:
8         # Set the seed for pseudo-random number generation
9         np.random.seed(random_state)
10
11    # Split the data into training and test sets
12    X_train, y_train = X[:n_train_samples], y[:n_train_samples]
13    X_test, y_test = X[n_train_samples:], y[n_train_samples:]
14
15    return X_train, y_train, X_test, y_test

```

4. Implemente una función que se llame `validacion_cruzada` que entrene k modelos y reporte las métricas obtenidasd: a. Divida el conjunto de entrenamiento en k subconjuntos excluyentes b. Para cada uno de los k modelos, utilice un subconjunto como validación c. Reporte la media y la desviación estándar para cada una de las métricas, todo debe realizarse solo usando Numpy:

- Accuracy
- Precision
- Recall
- F1

```

1 # Step 4: Cross-Validation Implementation
2 # This function performs cross-validation to evaluate model performance.
3 def cross_validation(X, y, k=5, max_depth=None, min_samples_split=2, criterion='gini'):
4     # Split the training set into k subsets
5     subsets_X = np.array_split(X, k)
6     subsets_y = np.array_split(y, k)
7
8     # Lists to store metrics for each model
9     accuracy_scores = []
10    precision_scores = []
11    recall_scores = []
12    f1_scores = []
13
14    for i in range(k):
15        # Select the current validation set
16        X_valid = subsets_X[i]
17        y_valid = subsets_y[i]
18
19        # Create the training set excluding the validation set
20        X_train = np.concatenate([subsets_X[j] for j in range(k) if j != i])
21        y_train = np.concatenate([subsets_y[j] for j in range(k) if j != i])
22
23        # Train a decision tree model
24        tree = DecisionTree(
25            max_depth=max_depth, min_samples_split=min_samples_split, criterion=criterion)
26        tree.root = tree.train(X_train, y_train)
27
28        # Make predictions on the validation set
29        predictions = tree.predict(X_valid)
30
31        # Calculate metrics and record them
32        accuracy = accuracy_score(y_valid, predictions)
33        precision = precision_score(y_valid, predictions)
34        recall = recall_score(y_valid, predictions)
35        f1 = f1_score(y_valid, predictions)
36
37        accuracy_scores.append(accuracy)
38        precision_scores.append(precision)
39        recall_scores.append(recall)
40        f1_scores.append(f1)
41
42    # Calculate the mean and standard deviation of the metrics
43    mean_accuracy = np.mean(accuracy_scores)
44    std_accuracy = np.std(accuracy_scores)
45    mean_precision = np.mean(precision_scores)
46    std_precision = np.std(precision_scores)
47    mean_recall = np.mean(recall_scores)
48    std_recall = np.std(recall_scores)
49    mean_f1 = np.mean(f1_scores)
50    std_f1 = np.std(f1_scores)
51
52    return {
53        "mean_accuracy": mean_accuracy,
54        "std_accuracy": std_accuracy,
55        "mean_precision": mean_precision,
56        "std_precision": std_precision,
57        "mean_recall": mean_recall,
58        "std_recall": std_recall,
59        "mean_f1": mean_f1,
60        "std_f1": std_f1
61    }
62

```

5. Entrene 10 combinaciones distintas de parámetros para su implementación de Arbol de Decisión y utilizando su implementación de validacion_cruzada.

6. Utilizando los resultados obtenidos analice cuál y porqué es el mejor modelo para ser usado en producción.

7. Compruebe las métricas usando el conjunto de prueba y analice el resultado

```

1 import numpy as np
2 import warnings
3
4 # Ignorar warnings
5 warnings.filterwarnings("ignore")
6
7 # Combinaciones
8 param_combinations = [
9     {'max_depth': None, 'min_samples_split': 2, 'criterion': 'gini'},
10    {'max_depth': None, 'min_samples_split': 4, 'criterion': 'entropy'},
11    {'max_depth': 5, 'min_samples_split': 2, 'criterion': 'gini'},
12    {'max_depth': 5, 'min_samples_split': 4, 'criterion': 'entropy'},
13    {'max_depth': 10, 'min_samples_split': 2, 'criterion': 'gini'},

```

```

14     {'max_depth': 10, 'min_samples_split': 4, 'criterion': 'entropy'},
15     {'max_depth': 15, 'min_samples_split': 2, 'criterion': 'gini'},
16     {'max_depth': 15, 'min_samples_split': 4, 'criterion': 'entropy'},
17     {'max_depth': 20, 'min_samples_split': 2, 'criterion': 'gini'},
18     {'max_depth': 20, 'min_samples_split': 4, 'criterion': 'entropy'}
19 ]
20
21 # Lista para almacenar los resultados de las metricas
22 results = []
23
24 # Iterar sobre las combinaciones de parametros
25 X = np.array([[2, 1],
26              [3, 2],
27              [4, 3],
28              [6, 4],
29              [7, 5],
30              [8, 6]])
31
32 y = np.array([1, 0, 1, 1, 1, 0])
33
34 X_train, y_train, X_test, y_test = manual_train_test_split(X, y, train_proportion=0.8, random_state=42)
35
36 for params in param_combinations:
37     # Entrenar y evaluar el modelo utilizando validacion cruzada
38     metrics = cross_validation(X_train, y_train, k=5, **params)
39
40     # Agrega los resultados a la lista
41     results.append({
42         'params': params,
43         'metrics': metrics
44     })
45
46 # Resultados 10 combinaciones - Primera mitad
47 for idx, result in enumerate(results[:5]):
48     print(f"Combinacion {idx + 1}:")
49     print(f"Parametros: {result['params']}")
50     print("Metricas:")
51     print(f"\t Punteria Media: {result['metrics']['mean_accuracy']}")
52     print(f"\t Punteria Desviacion Estandar: {result['metrics']['std_accuracy']}")
53     print(f"\t Precision Media: {result['metrics']['mean_precision']}")
54     print(f"\t Precision Desviacion Estandar: {result['metrics']['std_precision']}")
55     print(f"\t Recall Media: {result['metrics']['mean_recall']}")
56     print(f"\t Recall Desviacion Estandar: {result['metrics']['std_recall']}")
57     print(f"\t F1 Media: {result['metrics']['mean_f1']}")
58     print(f"\t Desviacion Estandar F1: {result['metrics']['std_f1']}")

```

Combinacion 1:

Parametros: {'max_depth': None, 'min_samples_split': 2, 'criterion': 'gini'}

Metricas:

```

Punteria Media: nan
Punteria Desviacion Estandar: nan
Precision Media: 0.4
Precision Desviacion Estandar: 0.48989794855663565
Recall Media: 0.4
Recall Desviacion Estandar: 0.48989794855663565
F1 Media: 0.4
Desviacion Estandar F1: 0.48989794855663565

```

Combinacion 2:

Parametros: {'max_depth': None, 'min_samples_split': 4, 'criterion': 'entropy'}

Metricas:

```

Punteria Media: nan
Punteria Desviacion Estandar: nan
Precision Media: 0.6
Precision Desviacion Estandar: 0.48989794855663565
Recall Media: 0.6
Recall Desviacion Estandar: 0.48989794855663565
F1 Media: 0.6
Desviacion Estandar F1: 0.48989794855663565

```

Combinacion 3:

Parametros: {'max_depth': 5, 'min_samples_split': 2, 'criterion': 'gini'}

Metricas:

```

Punteria Media: nan
Punteria Desviacion Estandar: nan
Precision Media: 0.4
Precision Desviacion Estandar: 0.48989794855663565
Recall Media: 0.4
Recall Desviacion Estandar: 0.48989794855663565
F1 Media: 0.4
Desviacion Estandar F1: 0.48989794855663565

```

Combinacion 4:

Parametros: {'max_depth': 5, 'min_samples_split': 4, 'criterion': 'entropy'}

Metricas:

```

Punteria Media: nan
Punteria Desviacion Estandar: nan
Precision Media: 0.6

```

```

Precision Desviacion Estandar: 0.48989794855663565
Recall Media: 0.6
Recall Desviacion Estandar: 0.48989794855663565
F1 Media: 0.6
Desviacion Estandar F1: 0.48989794855663565
Combinacion 5:
Parametros: {'max_depth': 10, 'min_samples_split': 2, 'criterion': 'gini'}
Metricas:
Punteria Media: nan
Punteria Desviacion Estandar: nan
Precision Media: 0.4
Precision Desviacion Estandar: 0.48989794855663565
Recall Media: 0.4
Recall Desviacion Estandar: 0.48989794855663565
F1 Media: 0.4
Desviacion Estandar F1: 0.48989794855663565

1 # Resultados 10 combinaciones - Segunda mitad
2 for idx, result in enumerate(results[5:]):
3     print(f"Combinacion {idx + 6}:")
4     print(f"Parametros: {result['params']}")
5     print("Metricas:")
6     print(f"\t Punteria Media: {result['metrics']['mean_accuracy']}")
7     print(f"\t Punteria Desviacion Estandar: {result['metrics']['std_accuracy']}")
8     print(f"\t Precision Media: {result['metrics']['mean_precision']}")
9     print(f"\t Precision Desviacion Estandar: {result['metrics']['std_precision']}")
10    print(f"\t Recall Media: {result['metrics']['mean_recall']}")
11    print(f"\t Recall Desviacion Estandar: {result['metrics']['std_recall']}")
12    print(f"\t F1 Media: {result['metrics']['mean_f1']}")
13    print(f"\t Desviacion Estandar F1: {result['metrics']['std_f1']}")

Combinacion 6:
Parametros: {'max_depth': 10, 'min_samples_split': 4, 'criterion': 'entropy'}
Metricas:
Punteria Media: nan
Punteria Desviacion Estandar: nan
Precision Media: 0.6
Precision Desviacion Estandar: 0.48989794855663565
Recall Media: 0.6
Recall Desviacion Estandar: 0.48989794855663565
F1 Media: 0.6
Desviacion Estandar F1: 0.48989794855663565

Combinacion 7:
Parametros: {'max_depth': 15, 'min_samples_split': 2, 'criterion': 'gini'}
Metricas:
Punteria Media: nan
Punteria Desviacion Estandar: nan
Precision Media: 0.4
Precision Desviacion Estandar: 0.48989794855663565
Recall Media: 0.4
Recall Desviacion Estandar: 0.48989794855663565
F1 Media: 0.4
Desviacion Estandar F1: 0.48989794855663565

Combinacion 8:
Parametros: {'max_depth': 15, 'min_samples_split': 4, 'criterion': 'entropy'}
Metricas:
Punteria Media: nan
Punteria Desviacion Estandar: nan
Precision Media: 0.6
Precision Desviacion Estandar: 0.48989794855663565
Recall Media: 0.6
Recall Desviacion Estandar: 0.48989794855663565
F1 Media: 0.6
Desviacion Estandar F1: 0.48989794855663565

Combinacion 9:
Parametros: {'max_depth': 20, 'min_samples_split': 2, 'criterion': 'gini'}
Metricas:
Punteria Media: nan
Punteria Desviacion Estandar: nan
Precision Media: 0.4
Precision Desviacion Estandar: 0.48989794855663565
Recall Media: 0.4
Recall Desviacion Estandar: 0.48989794855663565
F1 Media: 0.4
Desviacion Estandar F1: 0.48989794855663565

Combinacion 10:
Parametros: {'max_depth': 20, 'min_samples_split': 4, 'criterion': 'entropy'}
Metricas:
Punteria Media: nan
Punteria Desviacion Estandar: nan
Precision Media: 0.6
Precision Desviacion Estandar: 0.48989794855663565
Recall Media: 0.6
Recall Desviacion Estandar: 0.48989794855663565

```

```
F1 Media: 0.6
Desviacion Estandar F1: 0.48989794855663565
```

```
1 # Analisis del modelo
2 best_model_idx = np.argmax([result['metrics']['mean_accuracy'] for result in results])
3 best_model_params = results[best_model_idx]['params']
4 best_model_metrics = results[best_model_idx]['metrics']
5
6 print("Mejor modelo:")
7 print(f"Parametros: {best_model_params}")
8 print("Metricas en Cross-Validation:")
9 print(f"\t Punteria Media: {best_model_metrics['mean_accuracy']}")
10 print(f"\t Punteria Desviacion Estandar: {best_model_metrics['std_accuracy']}")
11 print(f"\t Precision Media: {best_model_metrics['mean_precision']}")
12 print(f"\t Precision Desviacion Estandar: {best_model_metrics['std_precision']}")
13 print(f"\t Recall Media: {best_model_metrics['mean_recall']}")
14 print(f"\t Recall Desviacion Estandar: {best_model_metrics['std_recall']}")
15 print(f"\t F1 Media: {best_model_metrics['mean_f1']}")
16 print(f"\t Desviacion Estandar F1: {best_model_metrics['std_f1']}")

Mejor modelo:
Parametros: {'max_depth': None, 'min_samples_split': 2, 'criterion': 'gini'}
Metricas en Cross-Validation:
    Punteria Media: nan
    Punteria Desviacion Estandar: nan
    Precision Media: 0.4
    Precision Desviacion Estandar: 0.48989794855663565
    Recall Media: 0.4
    Recall Desviacion Estandar: 0.48989794855663565
    F1 Media: 0.4
    Desviacion Estandar F1: 0.48989794855663565
```

▼ Análisis:

Para los diferentes parámetros y los 2 criterios seleccionados respectivamente. Los resultados de media no cambian si se altera la profundidad y la cantidad de splits dentro de la muestra bajo los criterios seleccionados. Por lo tanto, lo único que queda es verificar cuál de los modelos nos da mejores resultados. En este caso, el mejor modelo es el gini y ese es el que escogemos.

```
1 # Step 7: Prueba en el set de prueba
2 selected_model = DecisionTree(**best_model_params)
3 selected_model.root = selected_model.train(X_train, y_train)
4
5 test_predictions = selected_model.predict(X_test)
6
7 test_accuracy = accuracy_score(y_test, test_predictions)
8 test_precision = precision_score(y_test, test_predictions)
9 test_recall = recall_score(y_test, test_predictions)
10 test_f1 = f1_score(y_test, test_predictions)
11
12 print("Metricas en el set de prueba:")
13 print(f"Punteria: {test_accuracy}")
14 print(f"Precision: {test_precision}")
15 print(f"Recall: {test_recall}")
16 print(f"F1 Score: {test_f1}")

Metricas en el set de prueba:
Punteria: 0.5
Precision: 0.5
Recall: 1.0
F1 Score: 0.6666666666666666
```

▼ Conclusiones:

Con los resultados podemos darnos cuenta que el modelo seleccionado la verdad nos da medias aceptables para el set de datos que se hicieron. Con el mejor resultado gini al momento del split para nuestros datos. Por lo tanto, es válido decir que el modelo seleccionado es el que nos da mejores resultados.

Rúbrica para la Implementación de un Árbol de Decisión

Nota: Esta rúbrica se basa en la calidad de la implementación y los resultados obtenidos, no en la cantidad de código.

1. Creación de la Clase Nodo (10 puntos)

- ☐ Se crea una clase `Nodo` con los atributos mencionados en las especificaciones (feature, umbral, gini, cantidad_muestras, valor, izquierda, derecha).
- ☐ Los atributos se definen correctamente y se asignan de manera apropiada.

2. Creación de la Clase Árbol de Decisión (20 puntos)

- ☐ Se crea una clase que implementa un árbol de decisión.
- ☐ La clase utiliza las funciones presentadas en el cuaderno.
- ☐ Se implementan los hiperparámetros solicitados (max_depth, min_split_samples, criterio).
- ☐ La clase es capaz de entrenar un árbol de decisión con los hiperparámetros especificados.

3. División de Datos (10 puntos)

- ☐ Los datos se dividen en conjuntos de entrenamiento y prueba de forma manual.
- ☐ Se utiliza Numpy o Pandas para realizar esta división.
- ☐ Se garantiza que los conjuntos sean excluyentes.

4. Implementación de Validación Cruzada (20 puntos)

- ☐ Se implementa la función `validacion_cruzada` correctamente.
- ☐ Los datos de entrenamiento se dividen en k subconjuntos excluyentes.
- ☐ Se entrena y evalúa un modelo para cada subconjunto de validación.
- ☐ Se calculan y reportan las métricas de accuracy, precision, recall y F1.
- ☐ Se calcula la media y la desviación estándar de estas métricas.

5. Entrenamiento de Modelos (20 puntos)

- ☐ Se entrenan 10 combinaciones distintas de parámetros para el árbol de decisión.
- ☐ Cada combinación se entrena utilizando la función `validacion_cruzada`.
- ☐ Los resultados de las métricas se registran adecuadamente.

6. Análisis de Modelos (10 puntos)

- ☐ Se analizan los resultados obtenidos y se selecciona el mejor modelo para ser utilizado en producción.
- ☐ Se proporciona una justificación clara y fundamentada sobre por qué se eligió ese modelo.

7. Prueba en el Conjunto de Prueba (10 puntos)

- ☐ Se comprueban las métricas del modelo seleccionado utilizando el conjunto de prueba.
- ☐ Se analizan los resultados y se comentan las conclusiones.

General (10 puntos)

- ☐ El código se documenta de manera adecuada, incluyendo comentarios que expliquen las secciones clave.
- ☐ El código se ejecuta sin errores y sigue buenas prácticas de programación.
- ☐ La presentación de los resultados es clara y fácil de entender.
- ☐ Se cumple con todos los requisitos y las especificaciones proporcionadas.

Puntuación Total: 100 puntos