

Inteligencia Artificial y Aprendizaje Automático
Maestría en Inteligencia Artificial Aplicada – Tecnológico de Monterrey
Actividad Semana 6: Árbol de Decisiones y Bosque Aleatorio

Riesgo Crediticio

Maestría en Inteligencia Artificial Aplicada
Prof. Luis Eduardo Falcón Morales

Tecnológico de Monterrey

Nombre(s): _____ Matrícula(s): _____

Esta Tarea se deberá resolver en equipos, de acuerdo a como fueron integrados en semanas pasadas.

El asignar un crédito sabemos que conlleva un riesgo para el prestamista en caso de que el deudor no pague al final la cantidad asignada, o inclusive el equivocarnos al negarle el préstamo a alguien que sí era confiable. Durante décadas se ha tratado de resolver dicho problema desde muchas áreas del conocimiento y en particular las técnicas de Aprendizaje Automático (Machine Learning) han brindado y siguen proporcionando nuevas formas de enfrentar este problema.

No existen muchas bases de datos abiertas bien documentadas sobre este problema, sin embargo los datos del archivo **SouthGermanCredit.asc** del South_German_Credit_Data_Set de la página de la UCI es una sobre la cual se hace mucha investigación en torno a minimizar el riesgo en la asignación de créditos. En esta tarea se trabajará con dichos datos y los puedes encontrar dentro del archivo zip de la siguiente liga: <https://archive.ics.uci.edu/ml/datasets/South+German+Credit>

En la página de la liga anterior también se encuentra el archivo **codetable.txt** en el archivo zip. Ahí puedes encontrar más información detallada sobre el significado y tipo de cada variable.

Al menos en las siguientes ligas puedes encontrar mayor información de dichos datos:

<https://www.kaggle.com/competitions/south-german-credit-prediction/overview>

<https://www.semanticscholar.org/paper/South-German-Credit-Data-Classification-Using-to-Religia-Pranoto/219c9968cfa2cbb802376cf88035fe5e664b4418>

<https://ieeexplore.ieee.org/document/9239944>

Así, con base a datos históricos, el objetivo es predecir si una persona es clasificada como confiable a la asignación de un crédito, o no lo es. La variable de salida se llama “kredit”.

1. Carga los datos y sustituye los nombres de las columnas del alemán al inglés de acuerdo a como se indica en la página de la UCI.

NOTA: Si lo deseas, puedes traducirlos y ponerlas en español.

2. Realiza una partición de los datos en el conjunto de entrenamiento del 85% y el de prueba de 15%. Los modelos se estarán entrenando con el método de validación cruzada, así que no es

necesario en este paso generar el conjunto de validación. Define como la variable X a todas las variables de entrada y a la variable Y como la variable de salida.

3. Como una primera aproximación (baseline) realizarás las siguientes transformaciones mínimas para generar los primeros modelos. En la misma página de la UCI se indica el tipo de variable de cada uno de los factores. Definen cuatro tipo de variables: categórica (categorical), ordinal (discretized quantitative), numérica (quantitative) y binaria (binary). Con base a dicha información realiza un Pipeline que incluya al menos las siguientes transformaciones:
 - a. Imputación a todas las variables de entrada, diferenciando entre el tipo de cada variable (decide y justifica que tipo de imputación realizas en cada caso).
 - b. Realiza un análisis de las variables numéricas (quantitative) de entrada y aplica una transformación que escale a todas ellas en un rango equiparable.
 - c. Aplica la transformación One-Hot encoding a las variables de entrada de tipo categórico y binaria. En particular, justifica por qué una variable binaria requeriría que se le aplique la transformación one-hot encoding. Por el momento dejar las variables ordinales sin transformar.
4. Llevarás un entrenamiento usando validación cruzada entre los siguientes tres modelos de aprendizaje automático: Regresión Logística, Árbol de Decisión y Bosque Aleatorio. Deberás llevar a cabo el entrenamiento de los tres de manera conjunta usando un ciclo FOR. Recuerda aplicar las transformaciones que definiste en tu Pipeline. El entrenamiento debe ser con las siguientes características:
 - a. Usa los parámetros predeterminados de cada modelo.
 - b. En cada iteración deben calcularse todas las siguientes métricas: accuracy, precision, recall, f1-score y Gmean. Todas estas métricas deben ser funciones que tú mismo debes definir (Es decir, no usar las funciones de dichas métricas que te proporciona scikit-learn. Sin embargo, sí puedes usar la información regresada por el método confusion_matrix() de scikit-learn para definir las métricas).
 - c. Usar validación cruzada estratificada con 5 particiones y con 3 repeticiones.
 - d. Imprimir el valor de todas estas métricas, tanto para los datos de entrenamiento, como para los de validación. Así como los diagramas de caja y bigotes de los tres modelos con la métrica “recall”. ¿Alguno de los modelos está subentrenado o sobreentrenado? Justifica tu respuesta.
 - e. En particular obtengamos algunas de las llamadas curvas de aprendizaje para algunos de estos casos. En cada gráfico debes incluir tus comentarios sobre el modelo generado:
 - i. Obtener las curvas de aprendizaje (learning_curve) en la cual se va incrementando el tamaño de la muestra para el modelo de regresión Logística con sus hiperparámetros predeterminados. Utilizar al menos 20 puntos en la partición de los conjuntos de entrenamiento y la métrica “f1-score”, como evaluación del desempeño de dicha función “learning_curve()”.
 - ii. Obtener las curvas de validación (validation_curve) en la cual se va incrementando la complejidad del hiperparámetro “max_depth” para el modelo de árbol de decisión con sus hiperparámetros predeterminados. Utilizar valores de máxima

profundidad desde 1 hasta 20 y con la métrica “f1-score” para la evaluación del desempeño del modelo.

- iii. Obtener las curvas de aprendizaje (`learning_curve`) en la cual se va incrementando el tamaño de la muestra para el modelo de regresión bosque aleatorio (`random forest`) con sus hiperparámetros predeterminados. Utilizar al menos 20 puntos en la partición de los conjuntos de entrenamiento y la métrica “recall”, como evaluación del desempeño del modelo.
5. Finalmente veamos la manera de mejorar los valores de los hiperparámetros de cada modelo, así como el problema del sobreentrenamiento de algunos de ellos. Para ello deberás usar el método `GridSearchCV()` de `scikit-learn`. Recuerda que este método hace una búsqueda de los mejores hiperparámetros de un modelo mediante el llamado formato de malla y aplicando validación cruzada. En cada caso puedes incrementar el máximo de iteraciones, “`max_iter`” para que tengas la convergencia adecuada para todas las combinaciones en cada modelo. Recuerda también aplicar las transformaciones que definiste en tu Pipeline. Para fines de este ejercicio se ha seleccionado para cada modelo una métrica diferente, que permita irte familiarizando con ellas. Puedes consultar su documentación de `GridSearchCV` en la siguiente liga:

https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.GridSearchCV.html

- a. Para el modelo de regresión logística realizar el entrenamiento buscando sus mejores hiperparámetros con `GridSearchCV()`. Los hiperparámetros que debes incluir en su búsqueda deben ser al menos los siguientes: `C`, `solver`, `class_weight` y `penalty`. En este caso deberás usar la métrica (scoring) “f1-score”. Imprime la mejor combinación de parámetros obtenidos, así como el valor del mejor desempeño (score) obtenido con la métrica f1. ¿Cuál es la utilidad de la métrica “f1-score”? Incluye tus conclusiones.

NOTA: Toma en cuenta que no todas las combinaciones de “solver” y “penalty” son posibles, para que lo tomes en cuenta al momento de realizar la búsqueda. Revisa la documentación.

- b. Con los mejores valores de los hiperparámetros encontrados con la métrica “f1-score” para el modelo de regresión logística, obtener las curvas de aprendizaje (`learning curve`), incrementando el tamaño del conjunto de entrenamiento al menos 20 veces. Si lo crees adecuado, puedes hacer los ajustes que consideres adecuados para mejorar el resultado y evitar el sobreentrenamiento o el subentrenamiento.
- c. Para el modelo de árbol de decisión (`decision tree`) realizar el entrenamiento buscando sus mejores hiperparámetros con `GridSearchCV()`. Los hiperparámetros que debes incluir en su búsqueda deben ser al menos los siguientes: `ccp_alpha`, `criterion`, `max_depth`, `min_samples_split` y `class_weight`. En este caso deberás usar la métrica (scoring) “precision”. Imprime la mejor combinación de parámetros obtenidos, así como el valor del mejor desempeño (score) obtenido con la métrica “precision”. ¿Cuál es la utilidad de la métrica “precision”? Incluye tus conclusiones.
- d. Con los mejores valores de los hiperparámetros encontrados con la métrica “precision” para el modelo de árbol de decisión, obtener las curvas de aprendizaje (`learning curve`), incrementando el tamaño del conjunto de entrenamiento al menos 20 veces. Si lo crees adecuado, puedes hacer los ajustes que consideres adecuados para mejorar el resultado y evitar el sobreentrenamiento o el subentrenamiento.

- e. Para el modelo de bosque aleatorio (random forest) realizar el entrenamiento buscando sus mejores hiperparámetros con GridSearchCV(). Los hiperparámetros que debes incluir en su búsqueda deben ser al menos los siguientes: ccp_alpha, criterion, max_depth, min_samples_split y class_weight. En este caso deberás usar la métrica (scoring) “recall”. Imprime la mejor combinación de parámetros obtenidos, así como el valor del mejor desempeño (score) obtenido con la métrica “recall”. ¿Cuál es la utilidad de la métrica “recall”? Incluye tus conclusiones.

NOTA: Toma en cuenta que el método de random forest puede tardar varios minutos en llevar a cabo

- f. Con los mejores valores de los hiperparámetros encontrados con la métrica “recall” para el modelo de bosque aleatorio, obtener las curvas de validación (validation curve), incrementando la complejidad del modelo a través del hiperparámetro “max_depth” con al menos 10 valores. Si lo crees adecuado, puedes hacer los ajustes que consideres adecuados para mejorar el resultado y evitar el sobreentrenamiento o el subentrenamiento.
6. Para cada uno de estos tres modelos, con las métricas que se consideraron en cada caso y usando el conjunto de Prueba que no has utilizado hasta ahora, obtener los modelos finales como se te indica a continuación. Deberás usar además como conjunto de entrenamiento el llamado modelo de entrenamiento “aumentado” que consiste en las datos que estuviste utilizando para entrenamiento y validación:
- a. Obtener el modelo de regresión logística con los mejores parámetros que hayas encontrado con la métrica f1-score utilizada. Imprimir el valor de dicha métrica e incluye tus conclusiones finales para este caso. Incluir un gráfico del árbol de decisión final obtenido.
 - b. Obtener el modelo de árbol de decisiones con los mejores parámetros que hayas encontrado con la métrica “precision” utilizada. Imprimir el valor de dicha métrica e incluye tus conclusiones finales para este caso.
 - c. Obtener el modelo de bosque aleatorio con los mejores parámetros que hayas encontrado con la métrica “recall” utilizada. Imprimir el valor de dicha métrica e incluye tus conclusiones finales para este caso.