

Árboles de Decisión

García Prado, Sergio

2 de marzo de 2017

Resumen

En este documento se analiza el comportamiento de los algoritmos de generación de árboles de decisión ID3 y J48 desde el punto de vista de la discretización de atributos (en el caso de ID3 de forma manual). Para ello se ha utilizado el conjunto de datos Thoracic Surgery Data Data Set[3]. La herramienta utilizada para el aprendizaje automático ha sido WEKA[2].

1. ¿POR QUÉ NO SE PUEDE APLICAR DIRECTAMENTE ID3?

El conjunto de datos utilizado[3] está formado por **470** instancias, las cuales describen resultados acerca de la esperanza de vida tras operaciones de cáncer pulmonar. El conjunto de datos presenta atributos con las siguientes características:

- 10 atributos de tipo categórico binario.
- 1 atributo de tipo categórico con 3 valores distintos.
- 1 atributo de tipo categórico con 4 valores distintos.
- 1 atributo de tipo categórico con 7 valores distintos.
- 1 atributo de tipo numérico real en el rango [1,44,6,3].
- 1 atributo de tipo numérico real en el rango [0,96,86,3].
- 1 atributo de tipo numérico entero en el rango [21,87].
- 1 atributo de tipo categórico binario, que representa el valor de la clase.

La razón por la cual no se puede aplicar el algoritmo *ID3* sobre este conjunto de datos es la existencia de atributos de naturaleza numérica, para los cuales dicho algoritmo no está diseñado.

2. REALICE LAS MODIFICACIONES PREVIAS EN EL FICHERO DE DATOS PARA QUE PUEDA LLEVAR A CABO LO ANTERIOR Y PROPORCIONE LOS RESULTADOS APLICANDO EL MÉTODO DE RETENCIÓN O HOLD-OUT PARA LA FORMACIÓN DEL EXPERIMENTO

Para poder aplicar el algoritmo *ID3* existen distintas alternativas, como la discretización de los atributos numéricos o la eliminación de los mismos. Nótese que la eliminación no es una buena alternativa debido a que de esta manera se deja de utilizar información valiosa para la tarea de clasificación, por tanto, en este caso se ha decidido realizar una discretización.

El método que se ha utilizado para dicha tarea ha sido el más básico posible, realizando una partición de 10 intervalos sobre el rango de valores que toma cada uno de los atributos numéricos.

3. VOLVIENDO AL FICHERO ORIGINAL, PASE EL ALGORITMO J48 TAMBIÉN APLICANDO EL MÉTODO DE RETENCIÓN O HOLD-OUT. ANALICE EL ÁRBOL OBTENIDO SIN PODA. ¿SE PODRÍA PRESCINDIR DE ALGÚN ATRIBUTO? SI ES ASÍ, HÁGALO Y COMPARE LOS RESULTADOS DE NUEVO

4. HABRÁ NOTADO QUE CUANDO SE USA ALGÚN ATRIBUTO NUMÉRICO, IMPLÍCITAMENTE SE APLICA UNA DISCRETIZACIÓN AL PLANTEAR LAS DIFERENTES RAMAS DEL ÁRBOL A PARTIR DE ÉL. ¿POR QUÉ ES MÁS EFICIENTE ESTA TÉCNICA QUE LA APLICADA EN 2?

5. PLANTEE, ENTONCES, UNA DISCRETIZACIÓN BASADA EN EL PUNTO ANTERIOR, AUNQUE NO RESULTE SER BINARIA, SINO EN TANTOS TRAMOS COMO INDUZCAN LOS VALORES USADOS AL FORMAR LAS RAMAS DEL ÁRBOL CON J48 SIN PODA (HOLD-OUT). INTRODUZCA ESTE FICHERO DE NUEVO AL ALGORITMO ID3. COMPARE LOS RESULTADOS CON EL J48 DE 3

REFERENCIAS

- [1] CALONGE CANO, T., AND ALONSO GONZÁLEZ, C. J. Técnicas de Aprendizaje Automático, 2016/17.
- [2] THE UNIVERSITY OF WAIKATO. Weka. <http://www.cs.waikato.ac.nz/ml/weka/>.
- [3] UCI MACHINE LEARNING REPOSITORY. Thoracic Surgery Data Data Set. <http://archive.ics.uci.edu/ml/datasets/Thoracic+Surgery+Data>.