

Esercitazione #3- Data Mining

D'Andrea Giuseppe

Domanda 1

a)

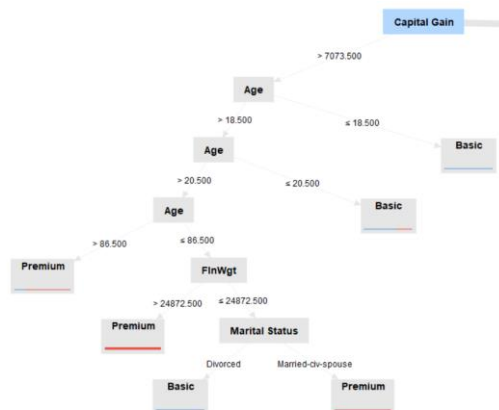
L'albero di decisione con l'algoritmo Decision Tree usando l'intero dataset (Utenti.xls) per il training, settando il "minimal gain" a 0.01 e mantenendo la configurazione di default per gli altri parametri, ha come attributo più selettivo "Capital Gain". I valori possibili di "Capital Gain" sono stati discretizzati ottenendo due range: " ≤ 7073.500 " e " > 7073.500 ".

b)

L'altezza dell'albero di decisione dipende dalle dimensioni e dalle caratteristiche del dataset utilizzato e rappresenta il cammino massimo, espresso come numero di archi attraversati, tra la radice ed una foglia. Per il dataset preso in esame (Utenti.xls), con "minimal gain" a 0.01 e mantenendo la configurazione di default per gli altri parametri, l'altezza dell'albero di decisione con l'algoritmo Decision Tree risulta essere pari a 9.

c)

In generale, un partizionamento puro identifica idealmente delle partizioni pure. Una foglia è pura se ha tutti gli elementi etichettati come appartenenti alla stessa classe. Nell'albero di decisione preso in esame, ci sono dei partizionamenti puri (come quello mostrato in figura) oltre a molti partizionamenti quasi puri.



Domanda 2

La profondità dell'albero di decisione varia in base alle dimensioni e alle caratteristiche del dataset. Il parametro "maximal depth" può essere usato per limitare l'altezza massima dell'albero di decisione. Come visto in precedenza, settando il "minimal gain" a 0.01 e mantenendo la configurazione di default per gli altri parametri (maximal depth = 10), l'altezza dell'albero di decisione risulta essere pari a 9. Impostando un valore di "maximal depth" minore di 10 (in RapidMiner il parametro tiene conto dei nodi e non degli archi attraversati nel cammino massimo, e dunque è necessario incrementare il valore di un'unità rispetto al valore di altezza dell'albero) la ricorsione dell'algoritmo viene interrotta ad un livello superiore e si genera un albero potenzialmente meno dettagliato.

L'algoritmo Decision Tree con la configurazione di default utilizza il "Gain Ratio" come criterio per calcolare il guadagno, che viene calcolato in un nodo prima di partizionarlo. Il nodo viene diviso se il suo guadagno è maggiore del guadagno minimo, configurabile in RapidMiner con il parametro "minimal gain". Un valore più alto di "minimal gain" si traduce in meno partizionamenti e quindi in un albero più piccolo. Un valore troppo alto impedirà completamente la divisione e verrà generato un albero con un singolo nodo. Viceversa, mantenendo costante "maximal depth" e decrementando il "minimal gain" si nota come l'albero di decisione diventi sempre più fitto e dunque più dettagliato.

Domanda 3

a)

Utilizzando "Native Country" come attributo di classe, l'albero di decisione con l'algoritmo Decision Tree usando l'intero dataset (Utenti.xls) per il training, settando il "minimal gain" a 0.01 e mantenendo la configurazione di default per gli altri parametri, ha come attributo più selettivo "Education-Num". I valori possibili di "Education-Num" sono stati discretizzati ottenendo due range: " ≤ 3.500 " e " > 3.500 ".

b)

L'altezza dell'albero di decisione dipende dalle dimensioni e dalle caratteristiche del dataset utilizzato e rappresenta il cammino massimo, espresso come numero di archi attraversati, tra la radice ed una foglia. Utilizzando "Native Country" come attributo di classe, per il dataset preso in esame (Utenti.xls), con "minimal gain" a 0.01 e mantenendo la configurazione di default per gli altri parametri, l'altezza dell'albero di decisione con l'algoritmo Decision Tree risulta essere pari a 9.

c)

In generale, un partizionamento puro identifica idealmente delle partizioni pure. Una foglia è pura se ha tutti gli elementi etichettati come appartenenti alla stessa classe. Nell'albero di decisione preso in esame, essendo l'attributo di classe non binario, non ci sono dei partizionamenti puri ma solo alcuni partizionamenti quasi puri.

Domanda 4

L'accuratezza misura il numero di predizione corrette rispetto al numero totale di predizioni fatte. È un indice di valutazione che funziona bene se la distribuzione delle classi è abbastanza bilanciata.

L'accuratezza aumenta al diminuire del parametro "maximal depth" anche se il livello di dettaglio dell'albero di decisione sta diminuendo sempre più. Questo fenomeno, per cui il modello diventa troppo semplice per classificare nuovi record in modo corretto, prende il nome di "underfitting".

Si può anche verificare il fenomeno opposto, quello dell'"overfitting", quando l'albero di decisione è molto dettagliato e ramificato. Così facendo il modello generato risulta troppo incentrato sul dataset di training per poter classificare in modo corretto nuovi record. Questo fenomeno si verifica quando si aumenta o si mantiene costante il valore di "maximal depth" e si diminuisce progressivamente il valore di "minimal gain".

[illegible]