

## Blatt: DTL

Michel Bünger

### DTL.01: Entscheidungsbäume mit CAL3 und ID3

#### CAL3

erster Durchlauf:

```
-----  
Alter(  
    ≥35(0:1)  
    <35: *  
)  
-----
```

```
-----  
Alter(  
    ≥35(0:1)  
    <35: (0:1)  
)  
-----
```

```
-----  
Alter(  
    ≥35(0:1, M:1)  
    <35: (0:1)  
)  
-----
```

```
-----  
Alter(  
    ≥35(0:1, M:2)  
    <35: (0:1)  
)  
-----
```

```

Alter(  

     $\geq 35(0:2, M:2) \rightarrow S_1 \text{ erreicht}, \frac{2}{4} < S_2;$   

    Differenzierung  

     $< 35: (0:1)$   

)
-----  

Alter(  

     $\geq 35($   

        hoch: (M:1 0:2)  

        niedrig:(M:1)  

)
-----  

<35: (0:2)  

)
-----  

Alter(  

     $\geq 35($   

        hoch: (M:1 0:2)  

        niedrig:(M:1)  

)
-----  

<35: (0:2, M:1)  

)
-----  

Erster Durchlauf abgeschlossen, Blätter nicht klar definiert.  

zweiter Durchlauf:  

-----  

Alter(  

     $\geq 35($   

        hoch: (M:1 0:3)  $\rightarrow S_1 \text{ erreicht}, \frac{3}{4} > S_2;$  Dominanz  

        von 0  

        niedrig:(M:1)  

)
-----  

<35: (0:2, M:1)  

)
-----  


```

```

Alter(
    ≥35(
        hoch: 0
        niedrig:(M:1)
    )

    <35: (0:3, M:1) →  $S_1$  erreicht,  $\frac{3}{4} > S_2$ ; Dominanz von
        0
    )

```

-----

```

Alter(
    ≥35(
        hoch: 0 ← Klasse M wird als Fehler
            klassifiziert
        niedrig:(M:1)
    )

    <35: 0
)

```

-----

```

Alter(
    ≥35(
        hoch: 0
        niedrig:(M:2)
    )

    <35: 0
)

```

-----

...

Für  $\text{Alter}(\geq 35(\text{hoch}:0))$  und  $\text{Alter}(<35: 0)$  kann keine Differenzierung mehr geschehen.  
Für  $\text{Alter}(\geq 35(\text{niedrig}: (M: 2)))$ , wird im vierten Durchlauf M: 4 gelten und somit  $S_1$  erreicht,  $\frac{4}{4} > S_2$ ; Dominanz von M.

...

finaler Ergebnisbaum:

```
-----  
    Alter(  
        ≥35(  
            hoch: 0  
            niedrig:M  
        )  
  
        <35: 0  
    )  
-----
```

## ID3

Berechnen der Entropie der Trainingsmenge:

$$H(S) = \frac{4}{7} \log_2(\frac{4}{7}) - \frac{3}{7} \log_2(\frac{3}{7})$$

$$H(S) = 0.985 \text{ Bit}$$

Berechnung der Informationsgewinne aller Attribute:

**Attribut: Alter**

- $\geq 35$ : O:2 M:2.  $H(\geq 35) = -\frac{2}{4} \log_2(\frac{2}{4}) - \frac{2}{4} \log_2(\frac{2}{4}) = 1$
- $< 35$ : O:2 M:1.  $H(< 35) = -\frac{2}{3} \log_2(\frac{2}{3}) - \frac{1}{3} \log_2(\frac{1}{3}) \approx 0.918$

Gewichtete Entropie von Alter:

$$H(S, \text{Alter}) = \frac{4}{7} H(\geq 35) + \frac{3}{7} H(< 35) \approx 0.965$$

Informationsgewinn für Alter:

$$Gain(S, \text{Alter}) = H(S) - H(S, \text{Alter}) = 0.02$$

### **Attribut: Einkommen**

- *hoch*: O:3 M:1.  $H(\text{hoch}) = -\frac{3}{4}\log_2(\frac{3}{4}) - \frac{1}{4}\log_2(\frac{1}{4}) \approx 0.811$
- *niedrig*: O:1 M:2.  $H(\text{niedrig}) = -\frac{1}{3}\log_2(\frac{1}{3}) - \frac{2}{3}\log_2(\frac{2}{3}) \approx 0.918$

Gewichtete Entropie von Einkommen:

$$H(S, \text{Einkommen}) = \frac{4}{7}H(\text{hoch}) + \frac{3}{7}H(\text{niedrig}) \approx 0.857$$

Informationsgewinn für Alter:

$$Gain(S, \text{Einkommen}) = H(S) - H(S, \text{Einkommen}) \approx 0.128$$

### **Attribut: Bildung**

- *Abitur*: O:1 M:2.  $H(\text{Abitur}) = -\frac{1}{3}\log_2(\frac{1}{3}) - \frac{2}{3}\log_2(\frac{2}{3}) \approx 0.918$
- *Bachelor*: O:1 M:1.  $H(\text{Bachelor}) = -\frac{1}{2}\log_2(\frac{1}{2}) - \frac{1}{2}\log_2(\frac{1}{2}) = 1$
- *Master*: O:2 M:0. Rein!  $H(\text{Master}) = 0$

Gewichtete Entropie von Bildung:

$$H(S, \text{Bildung}) = \frac{3}{7}H(\text{Abitur}) + \frac{2}{7}H(\text{Bachelor}) + \frac{2}{7}H(\text{Master}) \approx 0.679$$

Informationsgewinn für Alter:

$$Gain(S, \text{Bildung}) = H(S) - H(S, \text{Bildung}) \approx 0.306$$

Der Größte Informationsgewinn wird von Bildung versprochen, weshalb dies das erste Attribut sein wird, das behandelt wird.

### **Untermenge: Master**

Für die Untermenge 'Master' gehören alle Instanzen zu der Klasse O, was bedeutet, dass 'Master' rein ist und mit dem Blattknoten 'O' abschließt.

#### **0.0.1 Untermenge: Abitur**

Die Berechnungen der neuen Entropien zeigt hier, dass als nächstes Attribut, Einkommen den höchsten Informationsgewinn besitzt.

Zudem ist dann zu sehen, dass für  $\text{Abitur}(\text{Einkommen}(\text{hoch(O:1 M:0)}))$  und  $\text{Abitur}(\text{Einkommen}(\text{niedrig(O:0 M:2)}))$  gilt, d.h. beide sind rein und Blattknoten von jeweils 'O' und M'

## Untermenge: Bachelor

Die Berechnungen der neuen Entropien zeigen hier, dass als nächstes Attribut, Alter den höchsten Informationsgewinn besitzt.

Zudem ist dann zu sehen, dass für  $\text{Bachelor}(\text{Alter}(\geq 35(\text{O}:0 \text{ M}:1)))$  und  $\text{Bachelor}(\text{Alter}(\geq 35(\text{O}:1 \text{ M}:0)))$  gilt, d.h. beide sind rein und Blattknoten von jeweils 'M' und 'O'

fertiger Ergebnisbaum:

Bildung (

```

        Abitur(
            Einkommen(
                hoch: 0
                niedrig: M
            )
        )
        Bachelor(
            Alter(
                ≥35: M
                <35: O
            )
        )
    )
}

Master: O
)
```

## DTL.02: Pruning

Der gegebene Baum:

```
- - - - -  
x3(  
  x2(  
    x1(C, A),  
    x1(B, A)  
  ),  
  x1(  
    x2(C, B),  
    A  
  )  
)  
- - - - -
```

Lässt sich Kürzen, indem man die Allgemeine Transformationsregel:

$$x_1(x_2(a, b), x_2(c, d)) \leftrightarrow x_2(x_1(a, c), x_1(b, d))$$

Auf den ersten  $x_2(\dots)$  in  $x_3(x_2(\dots)\dots)$  anwendet.

Daraus entsteht:

```
- - - - -  
x3(  
  x1(  
    x2(C, B),  
    x2(A, A) → A  
  ),  
  x1(  
    x2(C, B),  
    A  
  )  
)  
- - - - -
```

Nun sind beide  $x_1(\dots)$  in  $x_3(x_1(\dots), x_1(\dots))$  gleich, weshalb sich der Baum letztlich ein letztes mal kürzen lässt:

$x_3$  (

$$x_1 ( \begin{array}{c} x_2 (C, B), \\ x_2 (A, A) \rightarrow A \end{array} )$$

)

Also ist der gekürzte Baum:  $x_3(x_1(x_2(C, B), A))$