
Data Mining und Maschinelles Lernen

Wintersemester 2015/2016

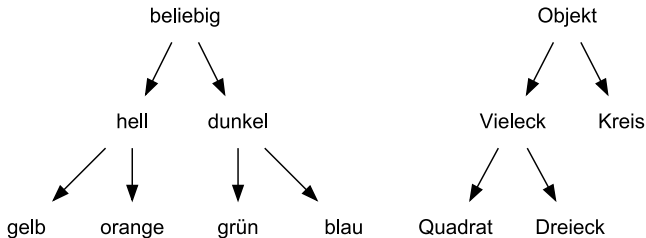
Lösungsvorschlag für das 3. Übungsblatt



TECHNISCHE
UNIVERSITÄT
DARMSTADT

Aufgabe 1: Version Space, Generalisierung und Spezialisierung (1)

Gegeben sei folgende Hierarchie von Begriffen:



Beobachtet werden Objekte, die durch Begriffspaare charakterisiert werden, die man an der untersten Ebene dieser Taxonomien finden kann (also z.B. „blaues Dreieck“). Konzepte können auch höherliegende Begriffe verwenden (also z.B. „dunkles Vieleck“). Überlegen Sie sich eine Generalisierungsvorschrift, die diese Taxonomien verwendet.

Aufgabe 1: Version Space, Generalisierung und Spezialisierung (2)



a) Wie sieht die minimale Generalisierung der Objekte „blauer Kreis“ und „grünes Dreieck“ aus?

Lösung: Gesucht wurde hier die gemeinsame Generalisierung zweier Objekte. Geht man von „grün“ und „blau“ den kürzesten Weg nach oben, landet man bei „dunkel“. Geht man von „Kreis“ und „Dreieck“ den minimalen Schritt nach oben, landet man bei „Objekt“.

⇒ „*dunkles Objekt*“

b) Wie sehen minimale Spezialisierungen des Konzepts „helles Objekt“ aus, sodaß das Beispiel „oranger Kreis“ nicht mehr abgedeckt wird?

Lösung: Entweder bei „hell“ eine Ebene weiter runter (ohne bei „orange“ zu landen) oder bei „Objekt“ (ohne bei „Kreis“ zu landen):

„gelbes Objekt“ und „helles Vieleck“. „Gelbes Vieleck“ ist keine zulässige Lösung, da diese Spezialisierung nicht minimal ist.

Aufgabe 1: Version Space, Generalisierung und Spezialisierung (3)

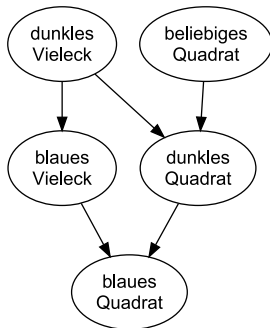
c) Gegeben seien folgende S und G-Sets:

G: { dunkles Vieleck, beliebiges Quadrat }

S: { blaues Quadrat }

Skizzieren Sie den Version Space, der durch diese Mengen definiert wird.

Lösung:



Aufgabe 1: Version Space, Generalisierung und Spezialisierung (4)



d) Wie würden Sie mit Hilfe des oben gegebenen Version Spaces die folgenden Beispiele klassifizieren (mit Begründung):

<u>Objekt</u>	<u>Klasse</u>
blaues Quadrat	
blauer Kreis	
blaues Dreieck	

Lösung: Das erste Beispiel ist im S-set enthalten und wird daher als positiv klassifiziert. Das zweite Beispiel ist nicht im G-set enthalten und wird daher als negativ klassifiziert. Beispiel 3 ist ein „dunkles Vieleck“, aber kein „beliebiges Quadrat“. Im S-set ist es nicht enthalten. Daher würde es nicht klassifizierbar sein und ein „?“ erhalten. Wenn Beispiele teilweise im S-Set und im G-Set sind (d.h., von diesen Hypothesen abgedeckt werden), können diese nie eindeutig klassifiziert werden. (Vergleiche *Lernen einzelner Regeln*, Folien *How to Classify these Examples?*)

Zusatzfrage: Wie würde man das Beispiel *blaues Vieleck* klassifizieren?

Aufgabe 1: Version Space, Generalisierung und Spezialisierung (5)

e) Gegeben seien wiederum die S- und G-sets aus c). Wie verändern sich die Sets nach Eintreffen des Beispiels (es wird der Candidate Elimination Algorithmus betrachtet):

gelbes Dreieck +

Lösung: Betrachtet wird der Candidate Elimination Algorithmus für das Auftreten eines positiven Beispiels.

G-Set: Hier wird jede Hypothese im G-Set entfernt, die das Beispiel nicht abdeckt. Da „gelbes Dreieck“ $\notin G$ ist, werden beide Hypothesen entfernt. Dann gilt $G = \{\}$.

S-Set: Da zum S-set nur die Hypothesen hinzugefügt werden, für die es eine Hypothese im G-set gibt, die genereller ist, werden keine Hypothesen hinzugefügt (da das G-set leer ist). Außerdem wurde das aktuelle Beispiel nicht abgedeckt. Daher gilt $S = \{\}$.

Wie interpretieren Sie dieses Ergebnis?

Lösung: Mit der Hinzunahme des Beispiels „gelbes Dreieck“ + ist das Konzept mit dem vorliegenden Version Space nicht mehr lernbar. Wenn man sich die Beispiele betrachten würde, die zur Entstehung des Version Spaces geführt haben, und das aktuelle Beispiel hinzunehmen würde, so wäre diese Beispielmenge mit dem Candidate Elimination Algorithmus nicht mehr lernbar.

Aufgabe 2: Version Space und Candidate Elimination Algorithmus (1)



Überlegen Sie sich eine geeignete Sprache, um den Candidate Elimination Algorithmus um die Behandlung von numerischen Daten zu erweitern.

a) Wie sieht eine passende Generalisierungsspezialisierungsvorschrift aus?

Lösung: Wir verwenden als Hypothesensprache Intervalle. Bei Generalisierungen sind diese abgeschlossen und bei Spezialisierungen entsprechend offen. Demnach sind die Hypothesen des G -Sets offen beziehungsweise die des S -Sets abgeschlossen.

Spezialisierungen erfolgen durch Einschränkung beziehungsweise Aufteilung von Intervallen. Generalisierungen entstehen durch Erweiterung von Intervallen.

Aufgabe 2: Version Space und Candidate Elimination Algorithmus (2)

b) Berechnen Sie den Version Space für folgende Beispiele:

Nr.	A1	A2	Klasse
1	0.5	1.5	-
2	1.1	1.2	+
3	1.8	1.0	+
4	1.5	2.1	-
5	2.1	1.2	-

Lösung: Wir verwenden den wie oben beschrieben modifizierten Candidate Elimination Algorithmus.

1. Wir beginnen mit:

- ▶ $G_0 = \{ \langle ?, ? \rangle \}$
- ▶ $S_0 = \{ \langle \emptyset, \emptyset \rangle \}$

2. Als nächsten erhalten wir folgendes Beispiel (0.5, 1.5, -):

- ▶ $S_1 = S_0$, da keine Hypothese in S_0 das Beispiel abdeckt
- ▶ $\{ \langle ?, ? \rangle \}$ deckt das Beispiel ab, deshalb muß diese Hypothese entfernt und minimal spezialisiert werden. Alle möglichen minimalen Spezialisierungen werden hinzugefügt, da G_0 die generellste Hypothese beinhaltet.
 $\Rightarrow G_1 = \{ \langle (-\infty, 0.5), ? \rangle, \langle (0.5, \infty), ? \rangle, \langle ?, (-\infty, 1.5) \rangle, \langle ?, (1.5, \infty) \rangle \}$

Aufgabe 2: Version Space und Candidate Elimination Algorithmus (3)

3. Nächstes Beispiel (1.1, 1.2, +):

- ▶ $S_1 = \{ \langle \emptyset, \emptyset \rangle \}$, $G_1 = \{ \langle (-\infty, 0.5), ? \rangle, \langle (0.5, \infty), ? \rangle, \langle ?, (-\infty, 1.5) \rangle, \langle ?, (1.5, \infty) \rangle \}$
- ▶ *Wir entfernen alle Hypothesen aus G_1 , die das Beispiel nicht abdecken:*
 $\langle (-\infty, 0.5), ? \rangle, \langle ?, (1.5, \infty) \rangle \Rightarrow G_2 = \{ \langle (0.5, \infty), ? \rangle, \langle ?, (-\infty, 1.5) \rangle \}$
- ▶ $\langle \emptyset, \emptyset \rangle$ aus S_1 deckt das Beispiel nicht ab und muß generalisiert werden:
 $\Rightarrow S_2 = \{ \langle [1.1, 1.1], [1.2, 1.2] \rangle \}$ (alle Hypothesen in G_2 sind genereller)

4. Nächstes Beispiel (1.8, 1.0, +):

- ▶ *Alle Hypothesen aus G_2 decken das Beispiel ab:* $\Rightarrow G_3 = G_2$
- ▶ $\langle [1.1, 1.1], [1.2, 1.2] \rangle$ aus S_2 deckt das Beispiel nicht ab und muß generalisiert werden: $\Rightarrow S_3 = \{ \langle [1.1, 1.8], [1.0, 1.2] \rangle \}$

Aufgabe 2: Version Space und Candidate Elimination Algorithmus (4)

5. Nächstes Beispiel (1.5, 2.1, -):

▶ $S_3 = \{ \langle [1.1, 1.8], [1.0, 1.2] \rangle \}$, $G_3 = \{ \langle (0.5, \infty), ? \rangle, \langle ?, (-\infty, 1.5) \rangle \}$

▶ Keine Hypothese in S_3 deckt das Beispiel ab:

⇒ $S_4 = S_3$

▶ Wir spezialisieren alle Hypothesen in G_3 , die das Beispiel abdecken:

$\langle (0.5, \infty), ? \rangle$

Mögliche Spezialisierungen werden hinzugefügt, falls eine Hypothese in S_3 spezifischer ist:

- ▶ $\langle (0.5, 1.5), ? \rangle$, keine Hypothese (in S) ist spezifischer
- ▶ $\langle (1.5, \infty), ? \rangle$, keine Hypothese (in S) ist spezifischer
- ▶ $\langle (0.5, \infty), (-\infty, 2.1) \rangle$, $\langle [1.1, 1.8], [1.0, 1.2] \rangle$ ist spezifischer
- ▶ $\langle (0.5, \infty), (2.1, \infty) \rangle$, keine Hypothese (in S) ist spezifischer

Alle anderen Hypothesen in G_3 (also $\langle ?, (-\infty, 1.5) \rangle$) bleiben unverändert.

⇒ $G_4 = \{ \langle ?, (-\infty, 1.5) \rangle, \langle (0.5, \infty), (-\infty, 2.1) \rangle \}$

Aufgabe 2: Version Space und Candidate Elimination Algorithmus (5)



6. Nächstes Beispiel (2.1, 1.2, -):

▶ $S_4 = \{ \langle [1.1, 1.8], [1.0, 1.2] \rangle \}$, $G_4 = \{ \langle ?, (-\infty, 1.5) \rangle, \langle (0.5, \infty), (-\infty, 2.1) \rangle \}$

▶ Keine Hypothese in S_4 deckt das Beispiel ab:

$\Rightarrow S_5 = S_4$

▶ Wir spezialisieren alle Hypothesen in G_4 , die das Beispiel abdecken:

$\langle (0.5, \infty), (-\infty, 2.1) \rangle, \langle ?, (-\infty, 1.5) \rangle$

Betrachten wir die möglichen Spezialisierungen der ersten Hypothese:

- ▶ $\langle (0.5, 2.1), (-\infty, 2.1) \rangle, \langle [1.1, 1.8], [1.0, 1.2] \rangle$ ist spezifischer
- ▶ $\langle (2.1, \infty), (-\infty, 2.1) \rangle$, keine Hypothese ist spezifischer
- ▶ $\langle (0.5, \infty), (-\infty, 1.2) \rangle$, keine Hypothese ist spezifischer
- ▶ $\langle (0.5, \infty), (1.2, 2.1) \rangle$, keine Hypothese ist spezifischer

Die möglichen Spezialisierungen der zweiten Hypothese lauten wie folgt:

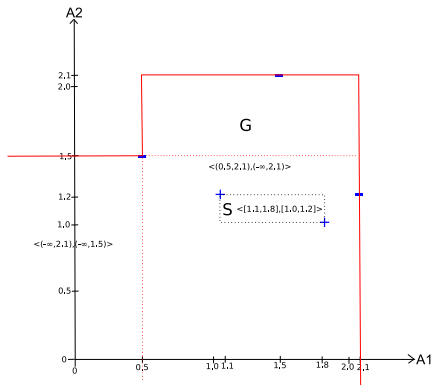
- ▶ $\langle (-\infty, 2.1), (-\infty, 1.5) \rangle, \langle [1.1, 1.8], [1.0, 1.2] \rangle$ ist spezifischer
- ▶ $\langle (2.1, \infty), (-\infty, 1.5) \rangle$, keine Hypothese ist spezifischer
- ▶ $\langle ?, (-\infty, 1.2) \rangle$, keine Hypothese ist spezifischer
- ▶ $\langle ?, (1.2, 1.5) \rangle$, keine Hypothese ist spezifischer

$\Rightarrow G_5 = \{ \langle (0.5, 2.1), (-\infty, 2.1) \rangle, \langle (-\infty, 2.1), (-\infty, 1.5) \rangle \}$

Aufgabe 2: Version Space und Candidate Elimination Algorithmus (6)

c) Skizzieren Sie das S-Set, das G-Set, und den Version Space im \mathbb{R}^2 .

Lösung:



Aufgabe 2: Version Space und Candidate Elimination Algorithmus (7)



d) Vertauschen Sie die Rolle der positiven und negativen Beispiele. Was passiert dann?

Lösung:

1. *Wir beginnen mit:*

▶ $G_0 = \{ \langle ?, ? \rangle \}$

▶ $S_0 = \{ \langle \emptyset, \emptyset \rangle \}$

2. *Als nächsten erhalten wir folgendes Beispiel (0.5, 1.5, +):*

▶ *Alle Hypothesen aus G_0 decken das Beispiel ab:*

⇒ $G_1 = G_0$

▶ *$\langle \emptyset, \emptyset \rangle$ aus S_0 deckt das Beispiel nicht ab und muß generalisiert werden:*

⇒ $S_1 = \{ \langle [0.5, 0.5], [1.5, 1.5] \rangle \}$

Aufgabe 2: Version Space und Candidate Elimination Algorithmus (8)



3. Nächstes Beispiel (1.1, 1.2, -):

- ▶ Keine Hypothese in S_1 deckt das Beispiel ab:
 $\Rightarrow S_2 = S_1$
 - ▶ $\{<?, ? >\}$ aus G_1 deckt das Beispiel ab und muß spezialisiert werden:
 - ▶ $<(-\infty, 1.1), ? >$, $<[0.5, 0.5], [1.5, 1.5] >$ ist spezifischer
 - ▶ $<(1.1, \infty), ? >$, keine Hypothese ist spezifischer
 - ▶ $<?, (-\infty, 1.2) >$, keine Hypothese ist spezifischer
 - ▶ $<?, (1.2, \infty) >$, $<[0.5, 0.5], [1.5, 1.5] >$ ist spezifischer
- $$\Rightarrow G_2 = \{<(-\infty, 1.1), ? >, <?, (1.2, \infty) >\}$$

4. Nächstes Beispiel (1.8, 1.0, -):

- ▶ Keine Hypothese in S_2 deckt das Beispiel ab:
 $\Rightarrow S_3 = S_2$
- ▶ Keine Hypothese in G_2 deckt das Beispiel ab:
 $\Rightarrow G_3 = G_2$

Aufgabe 2: Version Space und Candidate Elimination Algorithmus (9)

5. Nächstes Beispiel (1.5, 2.1, +):

- ▶ Wir entfernen alle Hypothesen aus G_3 , die das Beispiel nicht abdecken:
 $\langle (-\infty, 1.1), ? \rangle$
 $\Rightarrow G_4 = \{ \langle ?, (1.2, \infty) \rangle \}$
- ▶ $\langle [0.5, 0.5], [1.5, 1.5] \rangle$ aus S_3 deckt das Beispiel nicht ab und muß generalisiert werden:
 $\Rightarrow S_3 = \{ \langle [0.5, 1.5], [1.5, 2.1] \rangle \}$

6. Nächstes Beispiel (2.1, 1.2, +):

- ▶ Wir entfernen alle Hypothesen aus G_4 , die das Beispiel nicht abdecken:
 $\langle ?, (1.2, \infty) \rangle \Rightarrow G_5 = \emptyset$
- ▶ $\langle [0.5, 1.5], [1.5, 2.1] \rangle$ aus S_4 deckt das Beispiel nicht ab und muß generalisiert werden. Es können aber keine gültigen Generalisierungen gefunden werden, da $G_5 = \emptyset$ (Vergleiche "Lernen einzelner Regeln", Folie "Candidate Elimination Algorithm", Textstelle "some hypothesis...")
 $\Rightarrow S_5 = \emptyset$

Da sowohl das G-set und das S-set leer sind, folgt, dass diese Beispielmenge nicht mit dem Algorithmus und der gegebenen Hypothesensprache lernbar ist.

Aufgabe 3: Candidate Elimination Algorithmus (1)

Gegeben sei ein Datensatz mit drei Attributen:

Haarfarbe: *blond, braun, schwarz*
Größe: *klein, groß*
Augenfarbe: *grün, blau*

Der Hypothesenraum besteht aus Disjunktionen (Oder-Verknüpfungen) von maximal einem Wert pro Attribut, einer speziellsten Theorie *false*, die keine Beispiele abdeckt, und einer allgemeinsten Theorie *true*, die alle Beispiele abdeckt.

Zum Beispiel deckt die Hypothese $blond \vee blau$ alle Personen ab, die entweder blond oder blauäugig sind (in der Datenmenge aus Aufgabe 3b) sind das z.B. die Beispiele 1, 3, 4).

Beachte: Hypothesen wie $blond \vee braun$, die mehrere Werte desselben Attributs verwenden, sind nicht im Hypothesenraum.

Aufgabe 3: Candidate Elimination Algorithmus (2)



TECHNISCHE
UNIVERSITÄT
DARMSTADT

a) Geben Sie in dieser Hypothesensprache alle minimalen Generalisierungen und Spezialisierungen der Hypothese $blond \vee blau$ an.

Lösung: Durch Hinzufügen einer weiteren disjunktiven Terms kann man nur mehr Beispiele abdecken, daher generalisiert diese Operation. Wegstreichen eines Teilterms spezialisiert dagegen.

Minimale Generalisierungen = $\{ blond \vee blau \vee klein, blond \vee blau \vee groß \}$

Minimale Spezialisierungen = $\{ blond, blau \}$

Aufgabe 3: Candidate Elimination Algorithmus (3)



TECHNISCHE
UNIVERSITÄT
DARMSTADT

b) Folgende Beispiele treffen in dieser Reihenfolge ein:

1	<i>braun</i>	<i>groß</i>	<i>blau</i>	+
2	<i>braun</i>	<i>klein</i>	<i>grün</i>	–
3	<i>schwarz</i>	<i>klein</i>	<i>blau</i>	–
4	<i>blond</i>	<i>klein</i>	<i>grün</i>	+

Das erste Beispiel kodiert also eine Person, die braune Haare und blaue Augen hat und groß ist. Führen Sie auf diesen Beispielen den Candidate-Elimination Algorithmus zur Berechnung des Version Spaces durch und geben Sie nach jedem Schritt das *S*-Set und das *G*-Set an.

Aufgabe 3: Candidate Elimination Algorithmus (4)



Lösung:

$$G_0 = \{ \text{true} \}$$

$$S_0 = \{ \text{false} \}$$

$$G_1 = \{ \text{true} \}$$

$$S_1 = \{ \text{braun, groß, blau} \} \text{ (das sind drei Hypothesen!)}$$

$$G_2 = \{ \text{blond} \vee \text{groß} \vee \text{blau, schwarz} \vee \text{groß} \vee \text{blau} \}$$

$$S_2 = \{ \text{groß, blau} \} \text{ (die Hypothese } \text{braun} \text{ würde nun inkonsistent sein).}$$

$$G_3 = \{ \text{blond} \vee \text{groß} \} \text{ (schwarz} \vee \text{groß} \vee \text{blau} \text{ wird inkonsistent und zu } \text{groß} \text{ spezialisiert, was spezieller ist als } \text{blond} \vee \text{groß} \text{)}$$

$$S_3 = \{ \text{groß} \} \text{ (blau wird nun inkonsistent)}$$

$$G_4 = \{ \text{blond} \vee \text{groß} \}$$

$$S_4 = \{ \text{blond} \vee \text{groß} \} \text{ (auf } \text{groß} \vee \text{grün} \text{ können wir nicht generalisieren, da diese Theorie keine Spezialisierung eines Elements in } G \text{ ist).}$$

Der Version Space konvergiert also zu einer einzigen Lösung.