
Prohledávání prostoru hypotéz

Pro "klasickou úlohu strojového učení" potřebujeme:

- DATA
- hypotézy
- míru, nakolik daná hypotéza odpovídá datům

Data

- **DATA**: příklad na tabuli,
- jeden atribut je **cílový**, u nás EnjoySport, ostatní jsou **vstupní**
- cílem učení je najít hypotézu – funkci, která na základě vstupních parametrů správně určí cílový atribut
- **pozitivní příklady** jsou data s hodnotou cílového atributu Yes, **negativní příklady** jsou data s hodnotou cílového atributu No.

Prostor hypotéz

- **Hypotézy** formulujeme v určitém vyjadřovacím jazyce v našem případě konjunkce testů vstupních atributů, které charakterizují hodnotu cílového atributu Yes
 - hypotézy jsou formátu $\langle ?, Cold, High, ?, ?, ? \rangle$, kde
 - znak na pozici odpovídá podmínce na odpovídající vstupní (ne-cílový) atribut
 - znakem je buď konkrétní hodnota atributu, znak ? nekladoucí žádnou podmínku na daný atribut, znak \emptyset odpovídající nesplnitelné podmínce
- Pro binární atributy máme $4^{|\text{počet atributů}|}$ hypotéz, hypotézy obsahující \emptyset jsou ekvivalentní, tj. máme $3^{|\text{počet atributů}|} + 1$.
- Budeme prohledávat systematicky.

-
- Prostor hypotéz je **částečně uspořádaný inkluzí** $h_1 >_g h_2$ hypotéza h_1 je **obecnější** než h_2 (píšeme $h_1 >_g h_2$), pokud každý příklad splňující h_1 splňuje i h_2 . V tom případě se h_2 nazývá **specifičtější** než h_1 .
 - Např. $\langle ?, ?, \dots, ? \rangle$ je obecnější než $\langle \text{Sunny}, ?, \dots, \text{Same} \rangle$.
 - **Nejobecnější hypotéza** je $\langle ?, ?, \dots, ? \rangle$, tu splňují všechna data
 - **maximálně specifická hypotéza** je $\langle \emptyset, \emptyset, \dots, \emptyset \rangle$, kterou nesplňuje žádný záznam.
 - Prostor všech hypotéz tvoří svaz, viz. obrázek na tabuli.

Ohodnocovací funkce

- **ohodnocovací funkce** určuje, nakolik hypotéza odpovídá datům.
- Hledáme takovou hypotézu, kterou by splňovaly všechny pozitivní příklady a nesplňoval žádný negativní příklad.
- tj. aby byla implikace *hypoteza* \Rightarrow (*EnjoySport* = *Yes*) pro všechna data pravdivá
- (to lze, pokud máme data bez náhody a šumu).

Nalezení maximálně specifické hypotézy odpovídající datům

Algoritmus FIND-S

1. $h \leftarrow \langle \emptyset, \dots, \emptyset \rangle$ max. specifická hypotéza
2. pro každý pozitivní příklad x v datech
pro každou podmínku na atribut $A_i = a_i$ v h
Pokud příklad x nesplňuje $A_i = a_i$
nahraď podmínku nejbližší obecnější podmínkou,
kterou x splňuje
jinak nech h beze změny
3. vydej hypotézu h

Ale:

- Je hypotéza nalezená FIND–S jediná konzistentní s daty?
- Proč tedy volit ji, ne nějakou maximálně obecnou či něco mezi?
- V jiném prostoru hypotéz nemusí být ani maximálně specifická hypotéza jednoznačná.
- **Budeme hledat všechny hypotézy konzistentní s daty.**
- Pokud nejsou trénovací data konzistentní, máme problém. Řešení je jiný typ hypotéz a jiná ohodnocovací funkce.

Prostor verzí

Prostor verzí vzhledem k prostoru hypotéz H a trénovacích dat D je podmnožina hypotéz z H konzistentní s trénovacími daty D ,

$$VS_{H,D} = \{h \in H \mid \text{Consistent}(h, D)\}$$

- Tento prostor může být charakterizován obecnou a specifickou hranicí; každá hypotéza mezi těmito hranicemi spadá do prostoru verzí.
- **Obecná hranice G** vzhledem k H a D je množina maximálně obecných hypotéz z H konzistentních s daty, tj.

$$G = \{ g \in H \mid \text{Consistent}(g, D) \& (\neg \exists g^l \in H) [(g^l \succ_g g) \& \text{Consistent}(g^l, D)] \}$$

-
- **Specifická hranice S** vzhledem k H a D je množina maximálně specifických hypotéz z H konzistentních s daty, tj.

$$S = \{ s \in H \mid \text{Consistent}(g, D) \& \\ (\neg \exists s^l \in H) [(s >_g s^l) \& \text{Consistent}(s^l, G)] \}$$

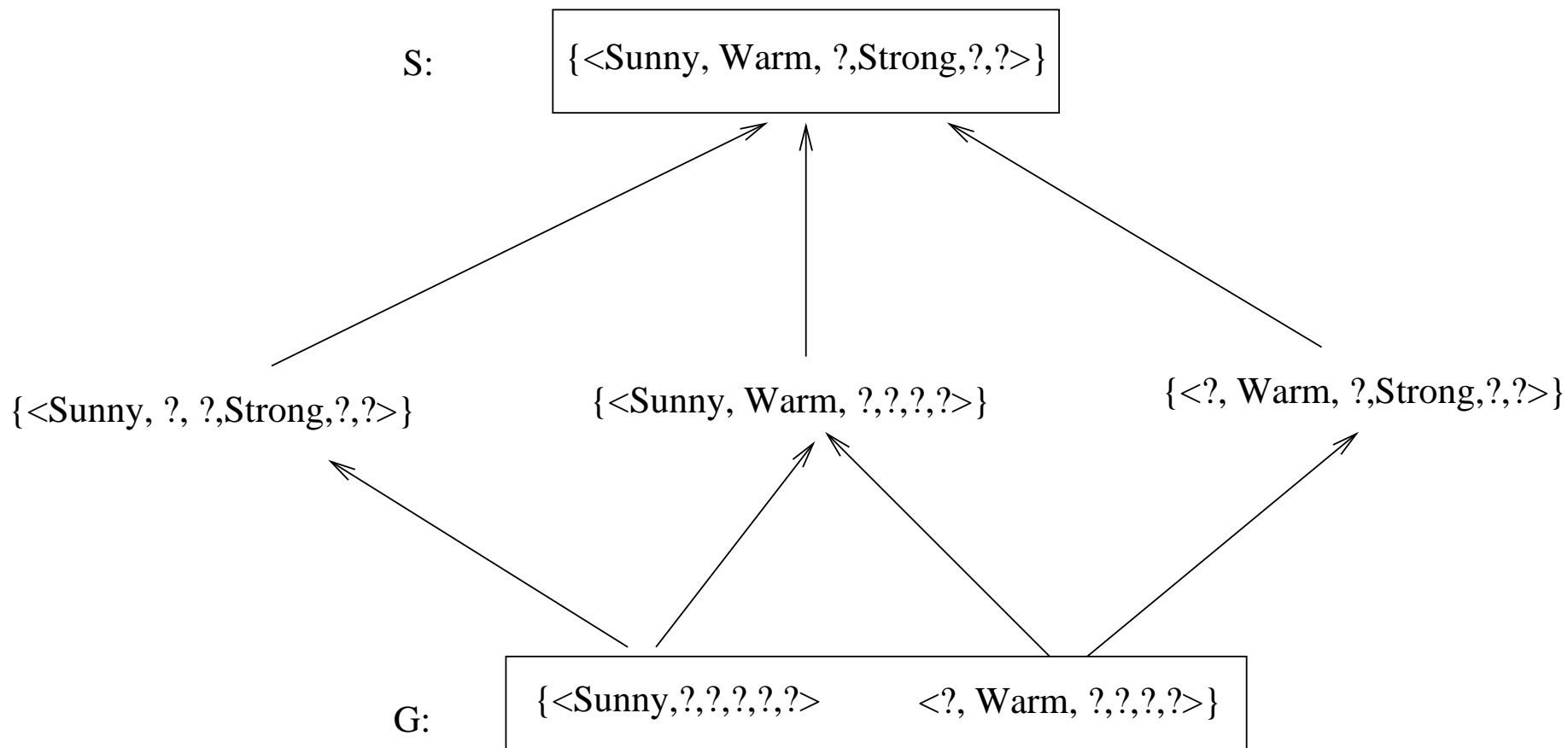


Figure 1: Prostor verzí s částečným uspořádáním inkluzí.

Algoritmus **Candidate–Elimination**

$G \leftarrow$ maximálně obecné hypotézy v H

$S \leftarrow$ maximálně specifické hypotézy v H

pokračuje

Pro každý trénovací příklad d , **do**

If d je pozitivní příklad

Odstraň z G všechny hypotézy nekonzistentní s d

For each $s \in S$, s nekonzistentní s d

Odstraň s z S

Přidej do S všechna h ; minimální zobecnění s taková, že

h je konzistentní s d a zároveň $\exists g \in G; g >_g h$

Odstraň z S hypotézy, které nejsou maximálně specifické v S

If d je negativní příklad

Odstraň z S všechny hypotézy nekonzistentní s d

For each $g \in G$, g nekonzistentní s d

Odstraň g z G

Přidej do G všechna h ; minimálně specifičtější než g taková, že

h je konzistentní s d a zároveň $\exists s \in S; h >_g s$

Odstraň z G hypotézy, které nejsou maximálně obecné v G