

# A Covering-based Algorithm for Classification: PRISM

Instructor: Dr. Lisa Fan

Speaker: Xiaofei Deng

Department of Computer Science

University of Regina

Regina, Saskatchewan, Canada S4S 0A2

E-mail: [deng200x@cs.uregina.ca](mailto:deng200x@cs.uregina.ca)

CS831: Knowledge Discovery in Databases

## Outline

- 1 **Background knowledge: ID3**
- 2 **Problem statement**
  - The problems of ID3
  - What causes this problem in ID3? (the inherent weakness)
- 3 **The PRISM algorithm**
  - An Information theoretic approach: PRISM
  - The basic steps of PRISM
  - An example for basic steps
  - Results of the example
  - Difference between ID3 and PRISM

## The basic idea of ID3.

- 1 Greedy Algorithm.
  - Select the attribute that contributes the maximum Information Gain.
- 2 Inductive bias: prefers small trees over large trees.
  - A short tree but might be a wide tree.
- 3 Its efficiency.
  - Been proved in theory by Quinlan.
  - Works well in chess endgames.

## Disadvantages of the representation of rules.

- 1 Difficult to manipulate for expert systems.

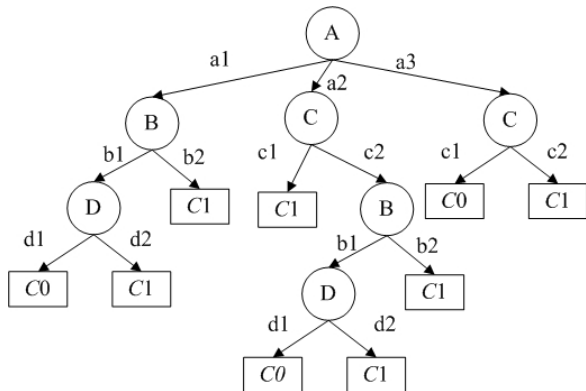
### Extract rules about a single classification

- Need to examine the whole tree.
- Partial solution: converting Decision Trees(DT) into a set of rules.
- Problems: There're rules can't easily be represented by DT.

### Example: extract rules about $C_0$ from a DT

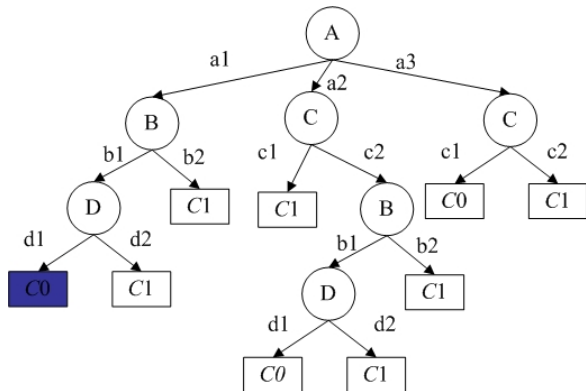
- *Rule1* :  $b_1 \wedge d_1 \rightarrow C_0$ , *Rule2* :  $a_3 \wedge c_1 \rightarrow C_0$ .
- Assume only two rules about  $C_0$ .
- Assume no attributes common to both Rules.

## Cont. (Extracting rules about $C_0$ )



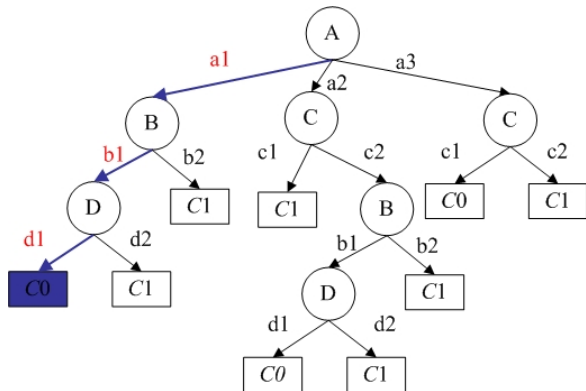
**Figure:** Extracting rules about  $C_0$  from decision tree

## Cont. (Extracting rules about $C_0$ )



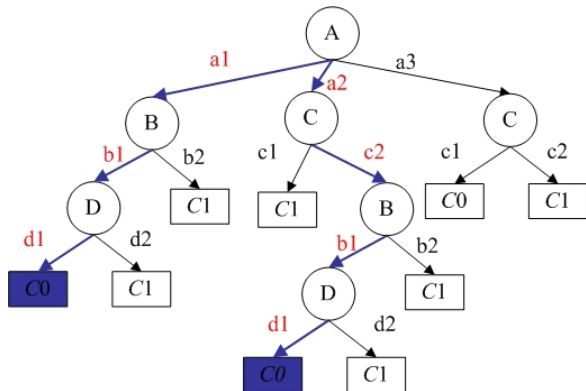
**Figure:** Extracting rules about  $C_0$  from decision tree

## Cont. (Extracting rules about $C_0$ )



**Figure:** Extracting rules about  $C_0$  from decision tree

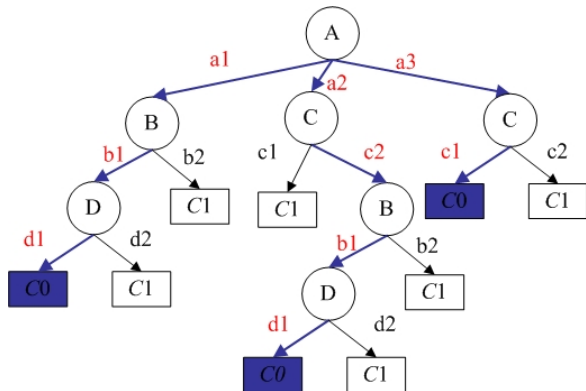
## Cont. (Extracting rules about $C_0$ )



**Figure:** Extracting rules about  $C_0$  from decision tree



## Cont. (Extracting rules about $C_0$ )



**Figure:** Extracting rules about  $C_0$  from decision tree

## Cont. (Extracted rules)

### Extracted Rules for Class $C_0$ from DT

- $Rule1a : a_1 \wedge b_1 \wedge d_1 \rightarrow C_0.$
- $Rule1b : a_2 \wedge c_2 \wedge b_1 \wedge d_1 \rightarrow C_0.$
- $Rule2 : a_3 \wedge c_1 \rightarrow C_0.$

### Explored the whole decision tree when extracting

- Why  $Rule1a, 1b$ ? Irrelevant attributes are added as a term to them.
- May cause serious problem, for example, a medical diagnose case which might requires an unnecessary surgery.

What causes this problem in ID3? (the inherent weakness)

## Information Entropy in ID3

- 1 The problem: ID3 Prefers an attribute which minimizes the **average Entropy**.

### Entropy

- 

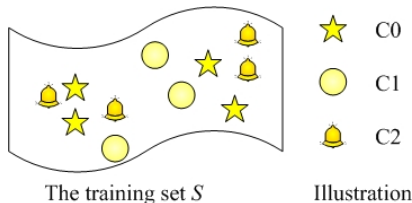
$$H(S) = - \sum_i^n p(C_i) \log_2(c_i) \text{ bits}$$

- $S, n, p(C_i)$  is the probability of occurrence of  $C_i$ .
- Entropy measures the uncertainty of current set of instances.

What causes this problem in ID3? (the inherent weakness)

## Why we say average Entropy?

- 1 Calculate the Entropy of a given set  $S$ .



**Figure:** The distribution of instances of  $S$

- 2  $H(S) = -p(C0)\log_2 p(C0) - p(C1)\log_2 p(C1) - p(C2)\log_2 p(C2)$ .
- 3 Measures the uncertainty in **Average**.
  - We added them to calculate the uncertainty.
  - Using  $H(S)$ , means consider all three,  $C0$ ,  $C1$ ,  $C2$ .

What causes this problem in ID3? (the inherent weakness)

## What about the uncertainty after knowing an Attribute?

- 1 ID3 chooses the attribute that contributed maximum information to lower the uncertainty.
- 2 But, that information measures in **average**.

### Information Gain



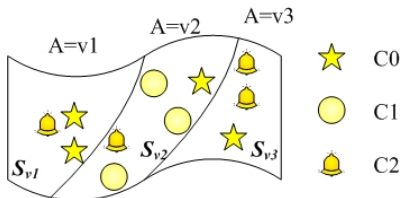
$$Gain(S, A) = H(S) - \sum_i \frac{|S_{vi}|}{|S|} H(S_{vi}) \text{ bits}$$

- Average entropy **Before** – **After** (knowing A).
- the second part is the info. A contributed.
- The second part measures the **average information** of all the branches of A.

What causes this problem in ID3? (the inherent weakness)

## Why the info. contributed by an attribute measures in average?

- 1 When choose attribute  $A$  ( $Gain(S, A)$  has max. value).
- 2 A partitions  $S$  into three branches,  $S_{v1}$ ,  $S_{v2}$ ,  $S_{v3}$ .



**Figure:** The training set  $S$  is partitioned by  $A$

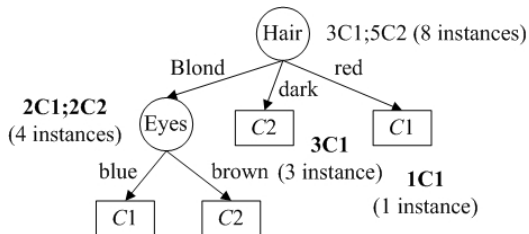
3

$$\sum_i \frac{|S_{vi}|}{|S|} H(S_{vi}) \text{ bits} = \frac{|S_{v1}|}{|S|} \text{Entropy}(\text{Branch } v1) \\ + \frac{|S_{v2}|}{|S|} \text{Entropy}(\text{Branch } v2) + \frac{|S_{v3}|}{|S|} \text{Entropy}(\text{Branch } v3)$$

What causes this problem in ID3? (the inherent weakness)

## Average dose not mean Good

- 1 An example: sometimes it would be worse for a branch



- 2 The average uncertainty of  $A$  is low.



$$\sum_1^3 \frac{|S_{v_i}|}{|S|} H(S_{v_i}) = 0.25 \text{ bits}$$

- 3 Uncertainty some branches of  $A$  is low, some rather high
  - Branch  $Hair = Blond$  is 0.5. **high**
  - Branch  $Hair = dark, Hair = red$  is 0. (low)

What causes this problem in ID3? (the inherent weakness)

## A short summary of the inner weakness of ID3

### ID3

- ID3 is attribute oriented.
- Selecting an attribute, then all the sub-branches are consider in average.
- ID3 measures the average information entropy.
- Average doesn't mean good to each rule.

### ID3 doesn't consider following cases

- An attribute might be highly **relevant** to only one classification and **irrelevant** to the others.
- Sometimes only one value of the attribute is **relevant**.



## How does PRISM fix this problem?

### The strategy of PRISM

- A branch could be considered as an attribute-value pair.
- Consider the relevance between an attribute-value pair and the specific classification.
- Choose the attribute-value pair that contributes maximum information as the term of a rule for one specific classification.

## An Information theoretic approach: PRISM

- 1 The task of PRISM.

**Find the  $\alpha_x$  that contributes maximum Information about  $C_i$ .**

- An attribute-value pair,  $\alpha_x$ .
- A specific classification,  $C_i$ .

- 2 The amount of Information about occurrence of  $C_i$  given  $\alpha_x$  is told:

$$I(C_i, \alpha_x)$$

$$= \log_2 \left( \frac{\text{Probability of occurrence of } C_i \text{ after knowing } \alpha_x}{\text{Probability of occurrence of } C_i \text{ before knowing } \alpha_x} \right) \text{bits}$$

$$= \log_2 \left( \frac{p(C_i | \alpha_x)}{p(C_i)} \right) \text{bits}$$

## Cont.

$$1 \quad I(C_i, \alpha_x) = \log_2\left(\frac{p(C_i|\alpha_x)}{p(C_i)}\right) \text{ bits}$$

$$2 \quad p(C_i|\alpha_x) = \frac{\text{Number of instances labeled } C_i}{|S_{\alpha_x}|}$$

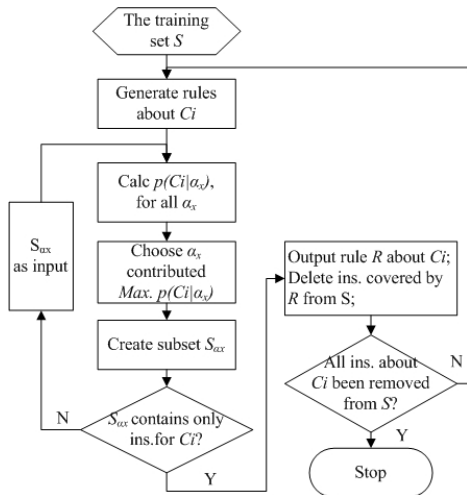
- The **After**.
- The probability of occurrence of  $C_i$  in  $S_{\alpha_x}$ .
- $S_{\alpha_x}$  is the subset of instances contain  $\alpha_x$ .

$$3 \quad p(C_i) = \frac{\text{Number of instances labeled } C_i}{|S|}$$

- The **Before**.
- The probability of occurrence of  $C_i$  in  $S$ .
- For all the  $\alpha_x$ , it's the same.
- Thus, we only calculate the  $p(C_i|\alpha_x)$ .

## PRISM algorithm: the basic steps

- Steps for generating rules about  $C_i$ , like  $C_1$ .



## Cont.(steps in detail)

- 1 Calculate the probability of occurrence,  $p(C_i|\alpha_x)$ , of the classification  $C_i$  for each attribute-value pair.
- 2 Select the attribute-value pair  $\alpha_x$  for which  $p(C_i|\alpha_x)$  is maximum, and create a subset,  $S_{\alpha_x}$ , that contains instances with  $\alpha_x$ .
- 3 Repeat step 1 and 2 for the subset, until it contains only instances for classification  $C_i$ . The induced rule is a conjunction of all the attribute-value pairs used in creating the subset.
- 4 remove all instances covered by this rule from the training set  $S$ .
- 5 Repeat Steps 1-4 until all instances of class  $C_i$  have been removed.

**Note. (For those steps)**

- 1  $p(C_i|\alpha_x)$  measures the contribution of  $\alpha_x$ .
- 2 Trying to find all rules about one specific classification  $C_i$ .

**Rules about Class  $C_1$** 

- *Rule1* :  $b_1 \wedge d_1 \rightarrow C_1$ .
- *Rule2* :  $a_3 \wedge c_1 \rightarrow C_1$ .

- 3 A rule is the conjunction of attribute-value pairs.

**Generating a rule about Class  $C_1$** 

- $\alpha_1$  : *Hair = Blond*. (1st attribute-value pair, term)
- $\alpha_2$  : *Eyes = Blue*. (2nd pair, term)
- *Rule1* :  $(\textit{Hair} = \textit{Blond} \wedge \textit{Eyes} = \textit{Blue}) \rightarrow C_1$

**Note. (For those steps)**

- 1  $p(C_i|\alpha_x)$  measures the contribution of  $\alpha_x$ .
- 2 Trying to find all rules about one specific classification  $C_i$ .

**Rules about Class  $C_1$** 

- *Rule1* :  $b_1 \wedge d_1 \rightarrow C_1$ .
- *Rule2* :  $a_3 \wedge c_1 \rightarrow C_1$ .

**Then  $C_2, \dots$** 

- *Rule3* :  $p_3 \wedge q_7 \rightarrow C_2$ .
- *Rule4* :  $k_2 \wedge t_5 \rightarrow C_2$ .

- 3 A rule is the conjunction of attribute-value pairs.

**Generating a rule about Class  $C_1$** 

- $\alpha_1$  : *Hair = Blond*. (1st attribute-value pair, term)
- $\alpha_2$  : *Eyes = Blue*. (2nd pair, term)
- *Rule1* :  $(\textit{Hair} = \textit{Blond} \wedge \textit{Eyes} = \textit{Blue}) \rightarrow C_1$

## An example for calculation

- ① Current training set  $S = \{1, 2, 3, 4, 5, 6, 7, 8\}$ .

Object	Height	Hair	Eyes	Class
O1	short	blond	blue	C1
O2	short	blond	brown	C2
O3	tall	red	blue	C1
O4	tall	dark	blue	C2
O5	tall	dark	blue	C2
O6	tall	blond	blue	C1
O7	tall	dark	brown	C2
O8	short	blond	brown	C2



An example for basic steps

## Generate rules for $C1$

- 1 Find 1st rule about  $C1$  ( $\rightarrow C1$ )

An example for basic steps

# Generate rules for $C1$

- 1 Find 1st rule about  $C1$  ( $\rightarrow C1$ )
- 2 Calculate all the  $p(C1|\alpha_x)$  for all  $\alpha_x$

$\alpha_x$	$C1$ (instances)	$S_{\alpha_x}$	$p(C1 \alpha_x)$
Height=short	{1}	{1,2,8}	1/3=0.333
Height=tall	{3,6}	{3,4,5,6,7}	2/5=0.4
Hair=blond	{1,6}	{1,2,6,8}	2/4=0.5
<u>Hair=red</u>	<u>{3}</u>	<u>{3}</u>	<u>1/1=1</u>
Hair=dark	{}	{4,5,7}	0
Eyes=blue	{1,3,6}	{1,3,4,5,6}	3/5=0.6
Eyes=brown	{}	{2,7,8}	0

**Figure:** Probability of occurrence of  $C1$  with each pair

An example for basic steps

**Calculate**  $p(C1 | Hair = blond)$

- 1 Probability of occurrence of  $C1$  with  $\alpha_x : Hair = blond$ .

An example for basic steps

## Calculate $p(C1 | Hair = blond)$

- 1 Probability of occurrence of  $C1$  with  $\alpha_x$  :  $Hair = blond$ .

Object	Height	Hair	Eyes	Class
O1	short	blond	blue	C1
O2	short	blond	brown	C2
O6	tall	blond	blue	C1
O8	short	blond	brown	C2

- 2  $p(C1 | \alpha_x) = p(C1 | Hair = blond) = \frac{|\{1,6\}|}{|\{1,2,6,8\}|} = \frac{2}{4} = 0.5$ .

## Output the *Rule1*

- 1 Choose  $\alpha_x : \textit{Hair} = \textit{red}$  as the first term for *Rule1* :  $(\textit{Hair} = \textit{red}) \wedge (\dots) \rightarrow C1$ .

## Output the *Rule1*

- 1 Choose  $\alpha_x : \textit{Hair} = \textit{red}$  as the first term for *Rule1* :  $(\textit{Hair} = \textit{red}) \wedge (\dots) \rightarrow C1$ .
- 2 Create subset  $S_{\alpha_x} = S_{\textit{Hair}=\textit{red}} = \{3\}$

## Output the *Rule1*

- 1 Choose  $\alpha_x : \textit{Hair} = \textit{red}$  as the first term for *Rule1* :  $(\textit{Hair} = \textit{red}) \wedge (\dots) \rightarrow C1$ .
- 2 Create subset  $S_{\alpha_x} = S_{\textit{Hair}=\textit{red}} = \{3\}$
- 3  $S_{\textit{Hair}=\textit{red}} = \{3\}$  contains only instance *Object3* labeled by *C1*.

## Output the *Rule1*

- 1 Choose  $\alpha_x : \textit{Hair} = \textit{red}$  as the first term for *Rule1* :  $(\textit{Hair} = \textit{red}) \wedge (\dots) \rightarrow C1$ .
- 2 Create subset  $S_{\alpha_x} = S_{\textit{Hair}=\textit{red}} = \{3\}$
- 3  $S_{\textit{Hair}=\textit{red}} = \{3\}$  contains only instance *Object3* labeled by *C1*.
- 4 Output the *Rule1* :  $(\textit{Hair} = \textit{red}) \rightarrow C1$ .



## Delete *Object3* from the training set

- 1 Delete *Object3* from  $S$ , thus  $S = \{1, 2, 4, 5, 6, 7, 8\}$ .

An example for basic steps

## Delete *Object3* from the training set

- 1 Delete *Object3* from  $S$ , thus  $S = \{1, 2, 4, 5, 6, 7, 8\}$ .
- 2 Current training set  $S = \{1, 2, 4, 5, 6, 7, 8\}$ .

Object	Height	Hair	Eyes	Class
O1	short	blond	blue	C1
O2	short	blond	brown	C2
O3	tall	red	blue	C1
O4	tall	dark	blue	C2
O5	tall	dark	blue	C2
O6	tall	blond	blue	C1
O7	tall	dark	brown	C2
O8	short	blond	brown	C2

An example for basic steps

## Repeat to find the *Rule2* about *C1*

- 1 Recalculate the  $p(C1|\alpha_x)$  for all  $\alpha_x$ .

$\alpha_x$	$C1(\text{instances})$	$S_{\alpha_x}$	$p(C1 \alpha_x)$
Height=short	{1}	{1,2,8}	1/3=0.333
Height=tall	{6}	{4,5,6,7}	1/4=0.25
<u>Hair=blond</u>	<u>{1,6}</u>	<u>{1,2,6,8}</u>	<u>2/4=0.5</u>
Hair=dark	{}	{4,5,7}	0
<u>Eyes=blue</u>	<u>{1,6}</u>	<u>{1,4,5,6}</u>	<u>2/4=0.5</u>
Eyes=brown	{}	{2,7,8}	0

**Figure:** Selecting the first term of *Rule2* about *C1*

- 2 *Hair = blond*, *Eyes = blue* have the equal value.
- 3 Choose *Hair = blond* as 1st term for *Rule2*.

An example for basic steps

## The second term of *Rule2* about *C1*

- 1 Create the subset  $S_{\alpha_x} = S_{\text{Hair}=\text{blond}} = \{1, 2, 6, 8\}$
- 2 *Object2* and *Object8* are labeled with *C2*.
- 3 Take  $S_{\alpha_x} = S_{\text{Hair}=\text{blond}} = \{1, 2, 6, 8\}$  as the current set.  
Trying to find second term.

**Table** The subset  $S_{\alpha_x} = S_{\text{Hair}=\text{blond}}$

Object	Height	Hair	Eyes	Class
O1	short	blond	blue	<b>C1</b>
O2	short	blond	brown	C2
O6	tall	blond	blue	<b>C1</b>
O8	short	blond	brown	C2

An example for basic steps

## The second term of *Rule2* about *C1*

- 1 Create the subset  $S_{\alpha_x} = S_{\text{Hair=blond}} = \{1, 2, 6, 8\}$
- 2 *Object2* and *Object8* are labeled with *C2*.
- 3 Take  $S_{\alpha_x} = S_{\text{Hair=blond}} = \{1, 2, 6, 8\}$  as the current set.  
Trying to find second term.

$\alpha_x$	<i>C1</i> (instances)	$S_{\alpha_x}$	$p(C1   \alpha_x)$
Height=short	{1}	{1,2,8}	1/3=0.333
Height=tall	{1}	{1}	1/1=1
<u>Eyes=blue</u>	<u>{1,6}</u>	<u>{1, 6}</u>	<u>2/2=1</u>
Eyes=brown	{}	{2,8}	0

## Cont.

- 1 Choose the *Eyes = blue* as the second term (consistent).

## Cont.

- 1 Choose the  $Eyes = blue$  as the second term (consistent).
- 2 Create subset  $S_{\alpha'_x} = S_{Hair=blond \wedge Eyes=blue} = \{1, 6\}$ .

## Cont.

- 1 Choose the  $Eyes = blue$  as the second term (consistent).
- 2 Create subset  $S_{\alpha'_x} = S_{Hair=blond \wedge Eyes=blue} = \{1, 6\}$ .
- 3  $\{1, 6\}$  are all labeled with  $C1$ , output *Rule2*.



## Cont.

- 1 Choose the  $Eyes = blue$  as the second term (consistent).
- 2 Create subset  $S_{\alpha'_x} = S_{Hair=blond \wedge Eyes=blue} = \{1, 6\}$ .
- 3  $\{1, 6\}$  are all labeled with  $C1$ , output *Rule2*.
- 4 *Rule2 : (Hair = blond  $\wedge$  Eyes = blue)  $\rightarrow$  C1.*

## Cont.

- 1 Choose the  $Eyes = blue$  as the second term (consistent).
- 2 Create subset  $S_{\alpha'_x} = S_{Hair=blond \wedge Eyes=blue} = \{1, 6\}$ .
- 3  $\{1, 6\}$  are all labeled with  $C1$ , output  $Rule2$ .
- 4  $Rule2 : (Hair = blond \wedge Eyes = blue) \rightarrow C1$ .
- 5 Delete Object 1, 6 from current training set.

## Cont.

- 1 Choose the  $Eyes = blue$  as the second term (consistent).
- 2 Create subset  $S_{\alpha'_x} = S_{Hair=blond \wedge Eyes=blue} = \{1, 6\}$ .
- 3  $\{1, 6\}$  are all labeled with  $C1$ , output *Rule2*.
- 4 ***Rule2 : (Hair = blond  $\wedge$  Eyes = blue)  $\rightarrow$  C1.***
- 5 Delete Object 1, 6 from current training set.
- 6 No others instances labeled with  $C1$ , stop.

## Cont.

- 1 Choose the  $Eyes = blue$  as the second term (consistent).
- 2 Create subset  $S_{\alpha'_x} = S_{Hair=blond \wedge Eyes=blue} = \{1, 6\}$ .
- 3  $\{1, 6\}$  are all labeled with  $C1$ , output  $Rule2$ .
- 4  $Rule2 : (Hair = blond \wedge Eyes = blue) \rightarrow C1$ .
- 5 Delete Object 1, 6 from current training set.
- 6 No others instances labeled with  $C1$ , stop.
- 7 Repeat above steps for  $C2$ .

## The results by PRISM and ID3

### Results by PRISM

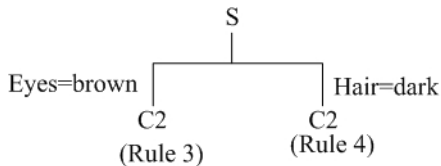
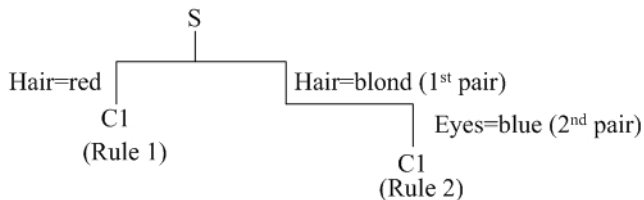
- $(\text{Hair} = \text{red}) \rightarrow C1.$
- $(\text{Hair} = \text{blond} \wedge \text{Eyes} = \text{blue}) \rightarrow C1).$
- $(\text{Eyes} = \text{brown}) \rightarrow C2.$
- $(\text{Hair} = \text{dark}) \rightarrow C2.$

### Results by ID3

- $(\text{Hair} = \text{red}) \rightarrow C1.$
- $(\text{Hair} = \text{blond} \wedge \text{Eyes} = \text{blue}) \rightarrow C1).$
- $(\text{Hair} = \text{blond} \wedge \text{Eyes} = \text{brown}) \rightarrow C2.$
- $(\text{Hair} = \text{dark}) \rightarrow C2.$

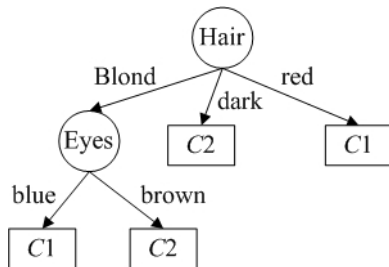
# Cont.

## 1 'Decision Tree' by PRISM



# Cont.

## 1 Decision Tree by ID3



## Summary

### ID3

- Greedy algorithm.
- Measures average information an attribute contributed.
- Attribute-oriented.
- Rules might contain irrelevant attributes.

### PRISM

- Greedy algorithm.
- Measures the attribute-value pair in determination of the classification.
- Attribute-value-oriented.
- More general and less rules.



## Q.&A.

# Any questions?