

Métodos de Inteligencia Artificial

L. Enrique Sucar (INAOE)

esucar@inaoep.mx

ccc.inaoep.mx/esucar

Tecnologías de Información

UPAEP

Agentes que Aprenden: Clasificador Bayesiano

- Clasificación
- Clasificador bayesiano
- Clasificador bayesiano simple
- Extensiones: TAN y BAN
- Mejora estructural
- Discretización

Clasificación

- El concepto de clasificación tiene dos significados:
 - No supervisada: dado un conjunto de datos, establecer clases o agrupaciones (*clusters*)
 - Supervisada: dadas ciertas clases, encontrar una regla para clasificar una nueva observación dentro de las clases existentes

Clasificación

- El problema de clasificación (supervisada) consiste en obtener el valor más probable de una variable (hipótesis) dados los valores de otras variables (evidencia, atributos)

$$\text{Arg}_H [\text{Max } P(H | E_1, E_2, \dots E_N)]$$

$$\text{Arg}_H [\text{Max } P(H | \mathbf{E})]$$

$$\mathbf{E} = \{E_1, E_2, \dots E_N\}$$

Tipos de Clasificadores

- Métodos estadísticos clásicos
 - Clasificador bayesiano simple (*naive Bayes*)
 - Discriminadores lineales
- Modelos de dependencias
 - Redes bayesianas
- Aprendizaje simbólico
 - Árboles de decisión, reglas, ...
- Redes neuronales, SVM, ...

Clasificación

- Consideraciones para un clasificador:
 - Exactitud – proporción de clasificaciones correctas
 - Rapidez – tiempo que toma hacer la clasificación
 - Claridad – que tan comprensible es para los humanos
 - Tiempo de aprendizaje – tiempo para obtener o ajustar el clasificador a partir de datos

Regla de Bayes

- La probabilidad posterior se puede obtener en base a la regla de Bayes:

$$P(H | \mathbf{E}) = P(H) P(\mathbf{E} | H) / P(\mathbf{E})$$

$$P(H | \mathbf{E}) = P(H) P(\mathbf{E} | H) / \sum_i P(\mathbf{E} | H_i) P(H_i)$$

- Normalmente no se requiere saber el valor de probabilidad, solamente el valor más probable de H

Regla de Bayes

- Para el caso de 2 clases $H: \{0, 1\}$, la regla de decisión de Bayes es:

$$H^*(E) = \begin{cases} 1 & \text{si } P(H=1 | \mathbf{E}) > 1/2 \\ 0 & \text{de otra forma} \end{cases}$$

- Se puede demostrar que la regla de Bayes es óptima

Valores Equivalentes

- Se puede utilizar cualquier función monótonica para la clasificación:

$$\text{Arg}_H [\text{Max } P(H | \mathbf{E})]$$

$$\text{Arg}_H [\text{Max } P(H) P(\mathbf{E} | H) / P(\mathbf{E})]$$

$$\text{Arg}_H [\text{Max } P(H) P(\mathbf{E} | H)]$$

$$\text{Arg}_H [\text{Max } \log \{P(H) P(\mathbf{E} | H)\}]$$

$$\text{Arg}_H [\text{Max } (\log P(H) + \log P(\mathbf{E} | H))]$$

Clasificador bayesiano simple

- Estimar la probabilidad: $P(\mathbf{E} | H)$ es complejo, pero se simplifica si se considera que los atributos son independientes dada la hipótesis:

$$P(E_1, E_2, \dots, E_N | H) = P(E_1 | H) P(E_2 | H) \dots P(E_N | H)$$

- Por lo que la probabilidad de la hipótesis dada la evidencia puede estimarse como:

$$P(H | E_1, E_2, \dots, E_N) = \frac{P(H) P(E_1 | H) P(E_2 | H) \dots P(E_N | H)}{P(\mathbf{E})}$$

- Esto se conoce como el clasificador bayesiano simple

Clasificador bayesiano simple

- Como veíamos, no es necesario calcular el denominador:

$$P(H | E_1, E_2, \dots, E_N) \sim P(H) P(E_1 | H) P(E_2 | H) \dots P(E_N | H)$$

- $P(H)$ se conoce como la *probabilidad a priori*, $P(E_i | H)$ es la probabilidad de los atributos dada la hipótesis (*verosimilitud*), y $P(H | E_1, E_2, \dots, E_N)$ es la *probabilidad posterior*

Ejemplo

- Para el caso del golf, cuál es la acción más probable (jugar / no-jugar) dado el ambiente y la temperatura?



Hoja de cálculo de
Microsoft Excel

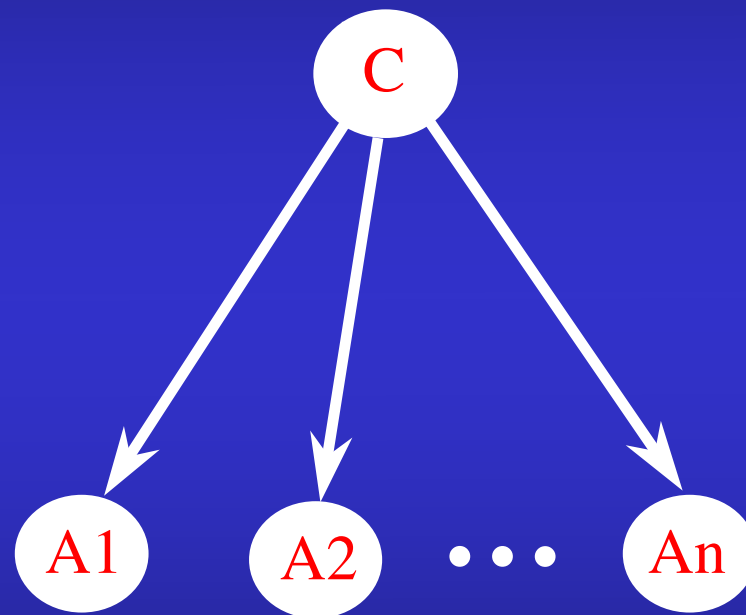
Ventajas

- Bajo tiempo de clasificación
- Bajo tiempo de aprendizaje
- Bajos requerimientos de memoria
- “Sencillez”
- Buenos resultados en muchos dominios

Limitaciones

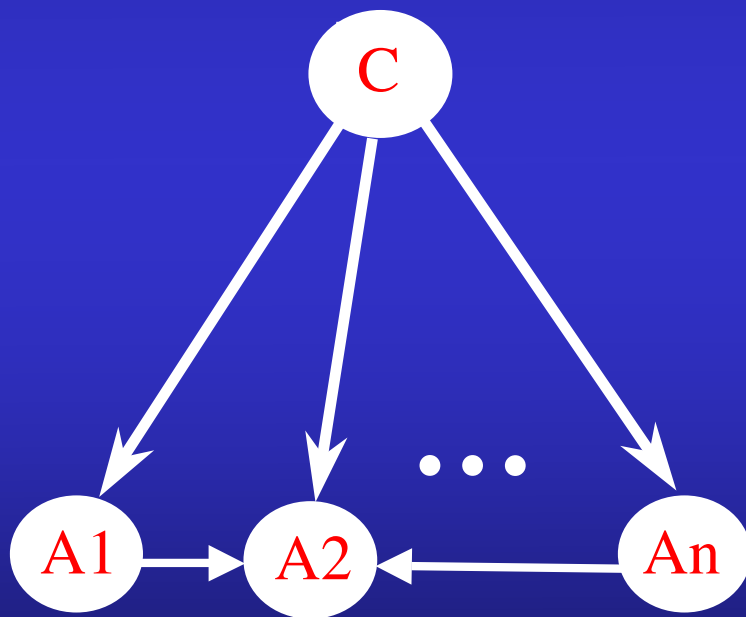
- En muchas ocasiones la suposición de independencia condicional no es válida
- Para variables continuas, existe el problema de discretización
- Alternativas – dependencias:
 - Estructuras que consideran dependencias
 - Mejora estructural del clasificador
- Alternativas – variables continuas:
 - Discriminador lineal (variables gaussianas)
 - Técnicas de discretización

CBS – modelo gráfico



Extensiones

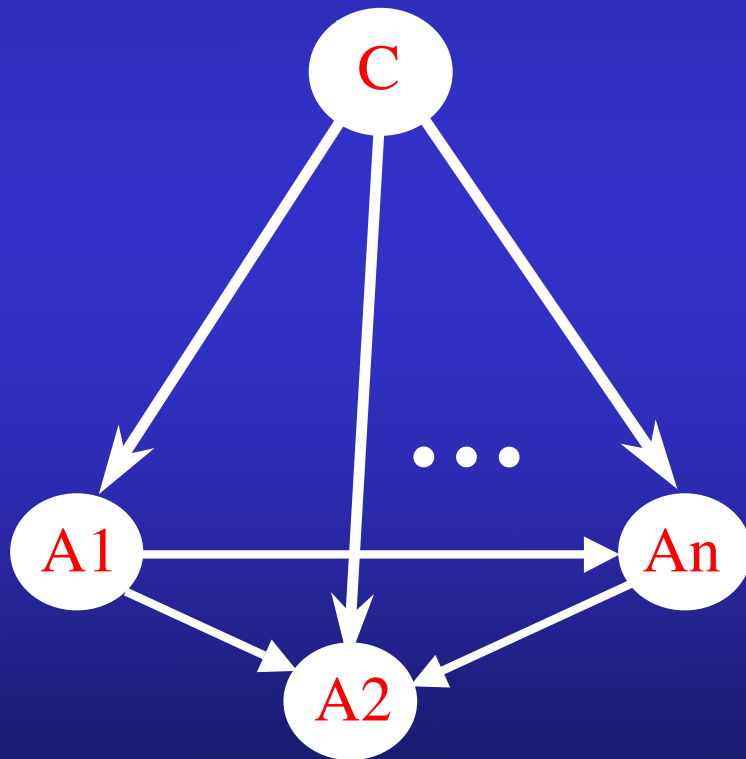
- TAN



Se incorpora algunas dependencias entre atributos mediante la construcción de un “árbol” entre ellos (técnica de Chow-Liu)

Extensiones

- BAN



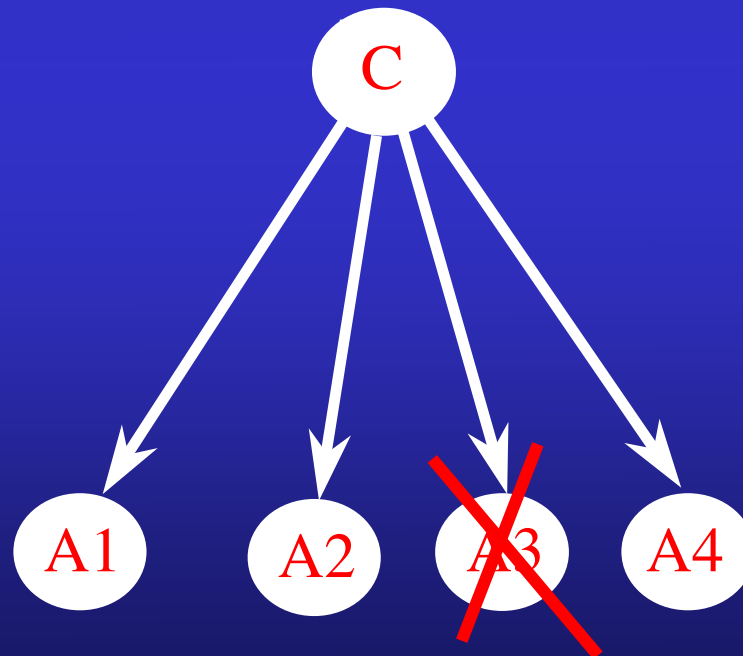
Se incorpora una “red” para modelar las dependencias entre atributos (aprendizaje de redes bayesianas).

Mejora estructural

- Otra alternativa para mejorar el CBS es partir de una estructura “simple” y modificarla mediante:
 - Eliminación de atributos irrelevantes (selección de atributos)
 - Verificación de las relaciones de independencia entre atributos y alterando la estructura:
 - Eliminar nodos
 - Combinar nodos
 - Insertar nodos

Eliminación de atributos

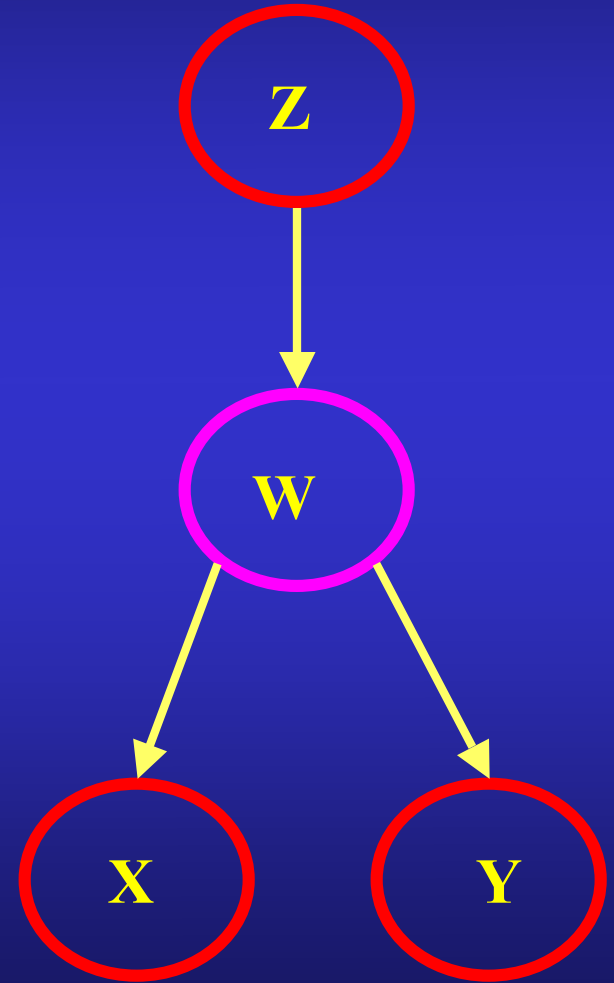
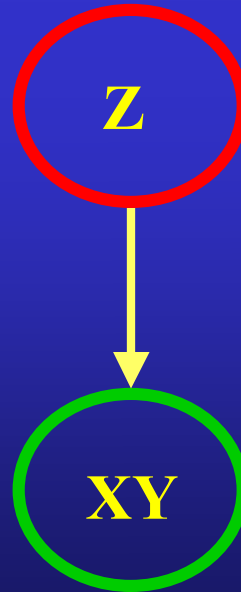
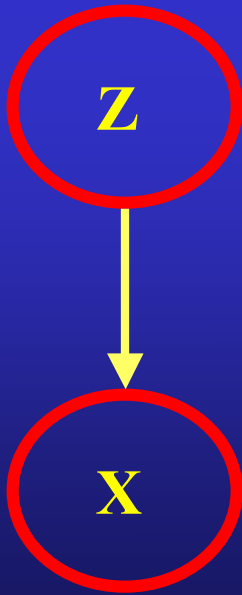
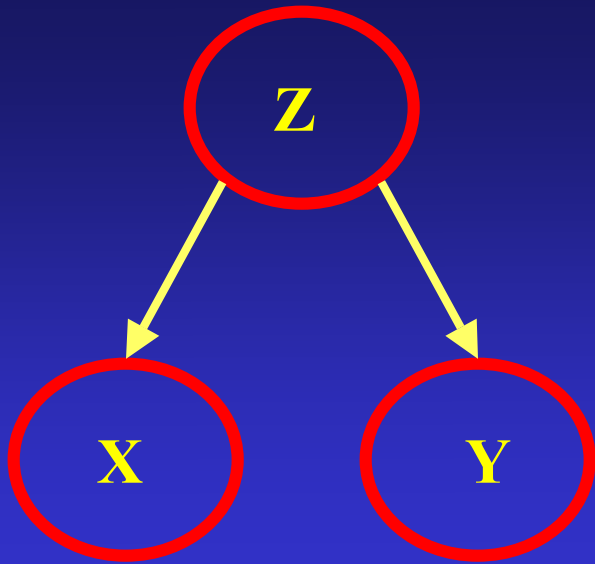
- Medir la “dependencia” entre la clase y atributos (por ejemplo con la información mutua), y eliminar aquellos con “poca” aportación



Mejora estructural

- Medir la dependencia entre pares de atributos dada la clase (por ejemplo mediante la información mutua condicional), alterar la estructura si hay 2 dependientes:
 1. Eliminación: quitar uno de los dos (redundantes)
 2. Unión: juntar los 2 atributos en uno, combinando sus valores
 3. Inserción: insertar un atributo “virtual” entre la clase y los dos atributos que los haga independientes.

Mejora Estructural



Atributos redundantes

- Prueba de dependencia entre cada atributo y la clase
- Información mutua:

$$I(C, A_i) = \sum P(C, A_i) \log [P(C, A_i) / P(C) P(A_i)]$$

- Eliminar atributos que no provean información a la clase

Atributos dependientes

- Prueba de independencia de cada atributo dada la clase
- Información mutua condicional

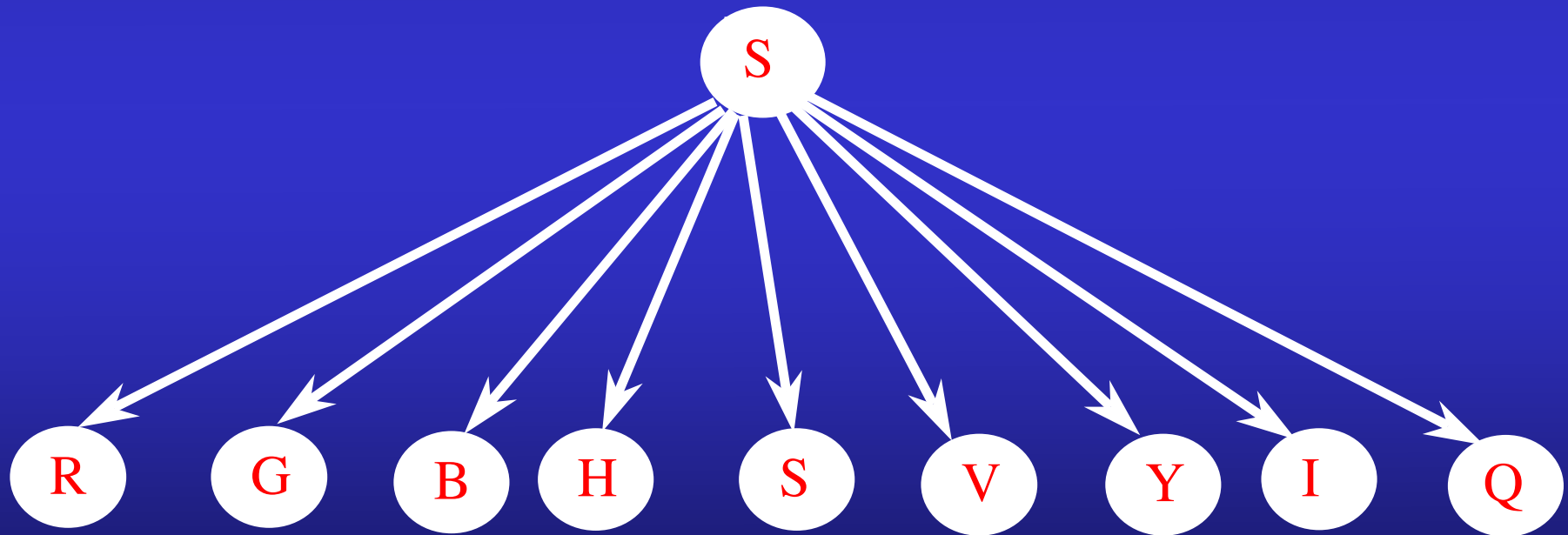
$$I(A_i, A_j | C) =$$

$$\sum P(A_i, A_j | C) \log [P(A_i, A_j | C) / P(A_i | C) P(A_j | C)]$$

- Eliminar, unir o (insertar) atributos

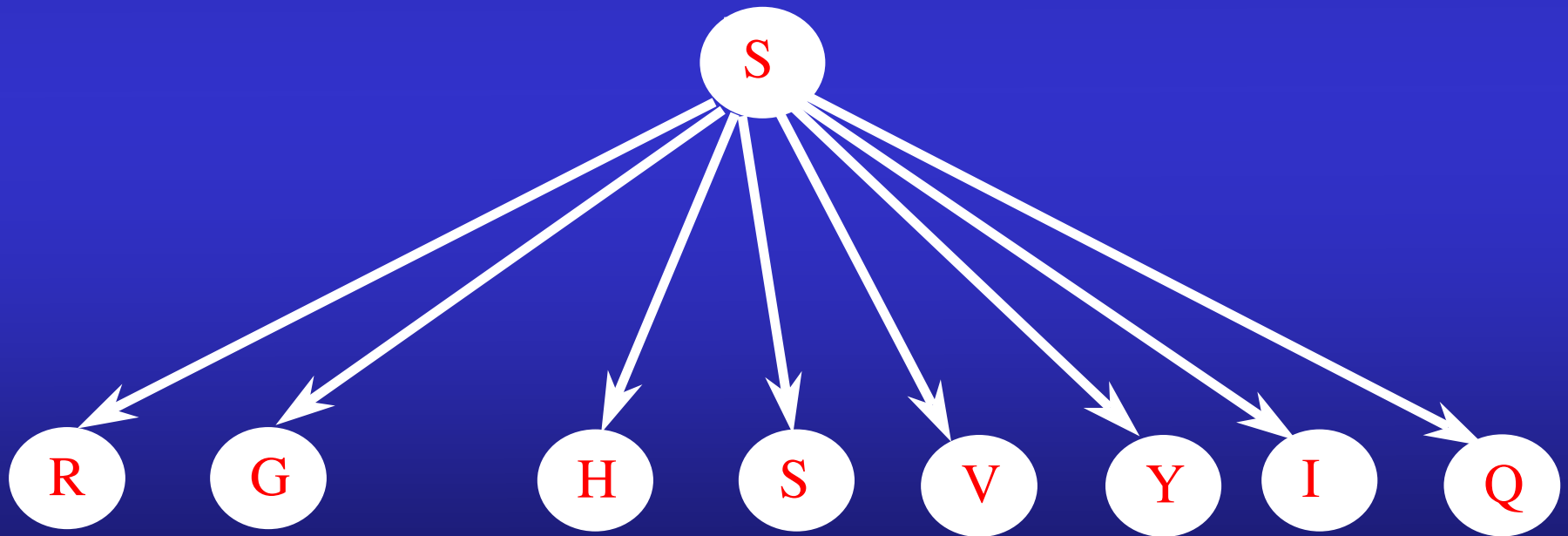
Ejemplo: clasificación de piel

- 9 atributos - 3 modelos de color: RGB, HSV, YIQ

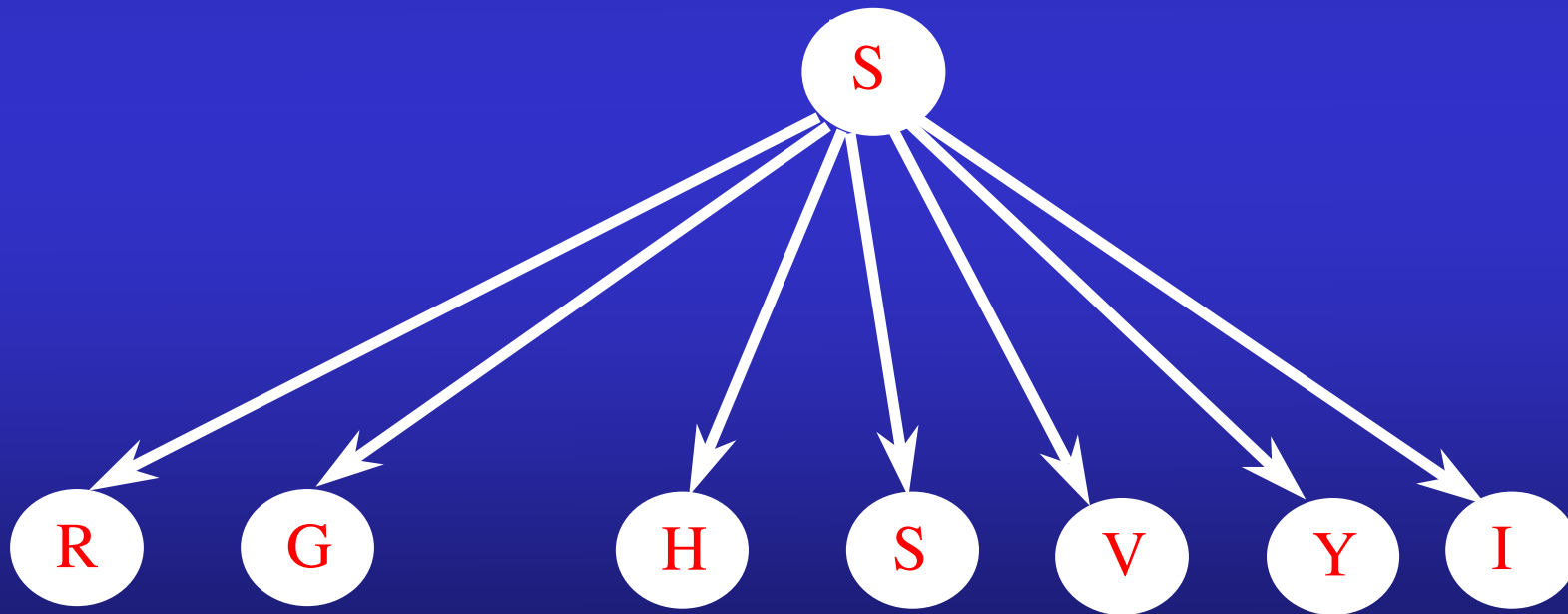


Mejora estructural

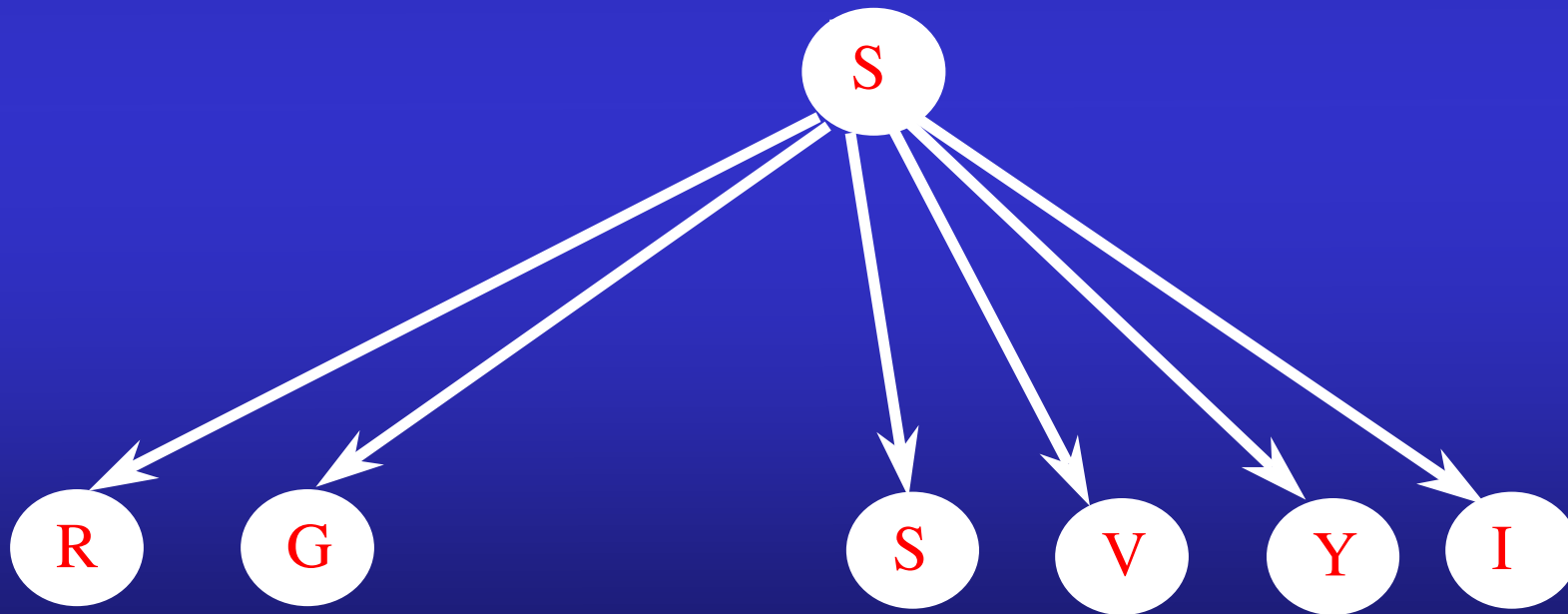
Elimina B



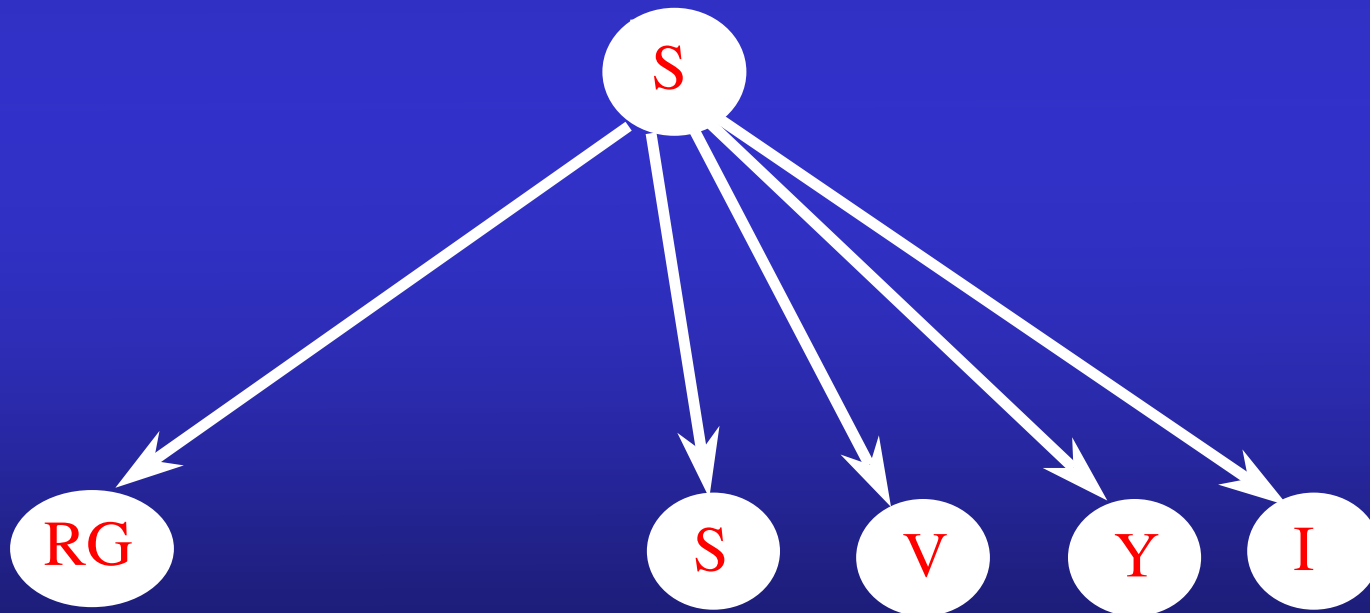
Elimina Q



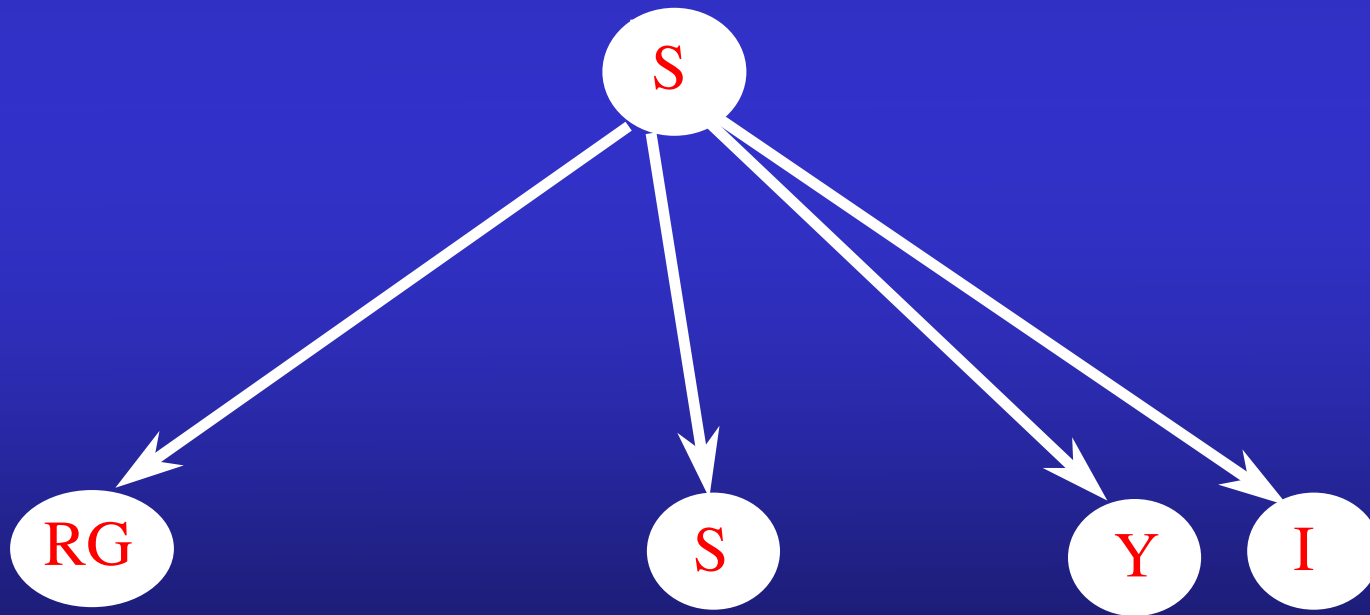
Elimina H



Unir RG

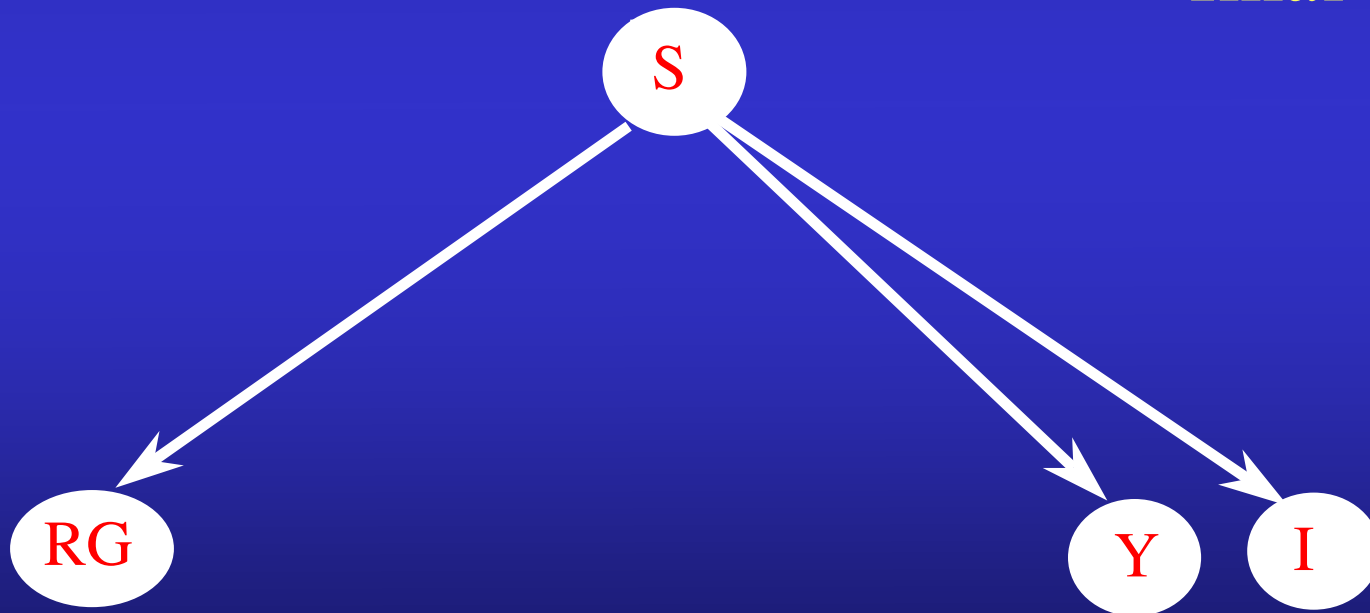


Elimina V



Elimina S

Exactitud: inicial 94%
final 98%



Discretización

- Si los atributos no siguen una distribución gaussiana, la alternativa es convertirlos a discretos agrupando los valores en un conjunto de rangos o intervalos
- Dos tipos de técnicas de discretización:
 - No supervisada: no considera la clase
 - Supervisada: en base a la clase

Discretización no supervisada

- Intervalos iguales
- Intervalos con los mismos datos
- En base al histograma

Discretización supervisada

- Considerando los posibles “cortes” entre clases:
 - Probar clasificador (con datos diferentes)
 - Utilizar medidas de información (p. ej., reducir la entropía)
- Problema de complejidad computacional

Costo de mala clasificación

- En realidad, no sólo debemos considerar la clase más probable si no también el costo de una mala clasificación
 - Si el costo es igual para todas las clases, entonces es equivalente a seleccionar la de mayor probabilidad
 - Si el costo es diferente, entonces se debe minimizar el costo esperado

Costo de mala clasificación

- El costo esperado (para dos clases, + y -) está dado por la siguiente ecuación:

$$CE = FN p(-) C(-|+) + FP p(+) C(+|-)$$

FN: razón de falsos negativos

FP: razón de falsos positivos

p: probabilidad de negativo o positivo

C(-|+): costo de clasificar un positivo como negativo

C(+|-): costo de clasificar un negativo como positivo

- Considerando esto y también la proporción de cada clase, existen técnicas más adecuadas para comparar clasificadores como la *curva ROC* y las *curvas de costo*

Tarea

- Leer Capítulo 19 de Russell
- Práctica de Redes Bayesianas (en la página del curso)