



Universidad
Carlos III de Madrid



Jesús García Herrero

CLASIFICADORES BAYESIANOS

En esta clase se presentan los algoritmos Análisis de Datos para abordar tareas de aprendizaje de modelos predictivos.

Se particularizan las técnicas estadísticas vistas anteriormente para resolver tareas predictivas: por un lado la regresión (lineal o no lineal) para predecir valores numéricos, y su aplicación para predecir probabilidades y clasificar instancias mediante regresión logística, y los clasificadores bayesianos como modelo de aprendizaje de parámetros de distribuciones de probabilidad para predecir probabilidades de clases.

La clasificación mediante regresión permite estimar las fronteras de decisión entre clases, presentando la equivalencia del aprendizaje de fronteras de clasificación al de las funciones de estimación de pertenencia a la clase. Se puede ver la limitación de esta técnica a problemas linealmente separables.

Los clasificadores bayesianos parten del principio de probabilidad condicionada para estimar probabilidades a posteriori de pertenencia a clases de las instancia, una vez calculadas las probabilidades de los valores de los atributos (probabilidades a priori) en la fase de entrenamiento. Estos clasificadores permiten tratar con datos nominales y numéricos, en este último caso utilizando distribuciones normales para simplificar el proceso. La limitación está en el cálculo de dependencias entre los atributos, que requeriría un número de datos exponencial con la dimensión de éstos, problema habitualmente tratado con la simplificación de independencia condicional (método "naïve Bayes"). Se completa el tema revisando aspectos prácticos que surgen al aplicar técnicas de clasificación sobre datos reales: tratamiento de datos incompletos y datos insuficientes para estimar probabilidades muy pequeñas.

Clasificadores Bayesianos

Métodos probabilísticos y numéricos de clasificación

Jesús García Herrero

Universidad Carlos III de Madrid



Universidad
Carlos III de Madrid



Clasificación numérica

- Modelado de datos con atributos numéricos para su aplicación a Clasificación. Generalización
- Datos representados como vectores de atributos numéricos: patrones

A_1	A_2	...	A_F
x_{11}^1	x_{21}^1	...	x_{F1}^1
x_{11}^2	x_{21}^2	...	x_{F1}^2
x_{11}^3	x_{21}^3	...	x_{F1}^3
...
x_{11}^N	x_{21}^N	...	x_{F1}^N

→ $\{\vec{X}^1, \vec{X}^2, \dots, \vec{X}^N\}$

- Problemas: **dimensionalidad, sobreajuste.**

Clasificación numérica

Problema de Clasificación

- Clases: $\{C_1, \dots, C_M\}$
- Muestras: $E = \{\vec{X}_1^{(1)}, \dots, \vec{X}_{n_1}^{(1)}, \vec{X}_1^{(2)}, \dots, \vec{X}_{n_2}^{(2)}, \dots, \vec{X}_1^{(M)}, \dots, \vec{X}_{n_M}^{(M)}\}$

– Tamaño:
$$N = \sum_{j=1}^M n_j$$

- Para cada clase, C_i , hay n_i patrones, cada uno con F atributos: para cada clase C_i : $\{\vec{X}_1^{(i)}, \dots, \vec{X}_{n_i}^{(i)}\}$

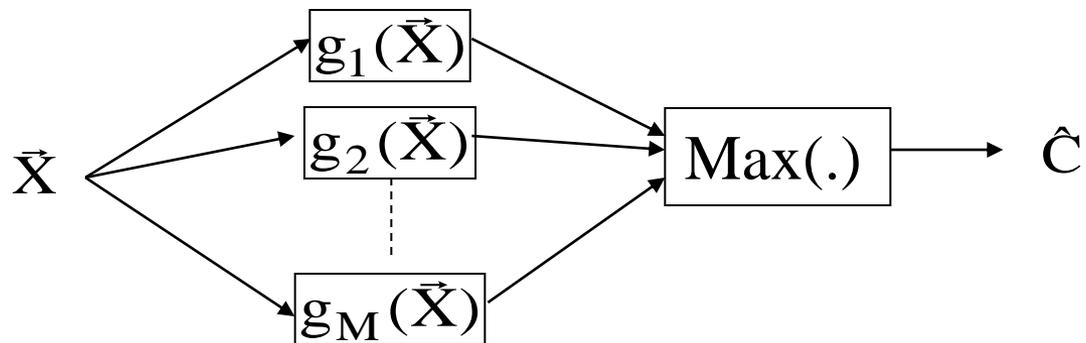
$$\vec{X}_j^{(i)} = \begin{bmatrix} x_{1j}^{(i)} \\ \vdots \\ x_{Fj}^{(i)} \end{bmatrix}; \quad j = 1, \dots, n_i$$

Clasificación numérica

$$g(\cdot): \mathbb{R}^F \longrightarrow C = \{C_1, \dots, C_M\}$$

$$\vec{X} \longrightarrow \hat{C} = g(\vec{X})$$

- Función discriminante de cada clase: $g_i(\vec{X}), i = 1, \dots, M$

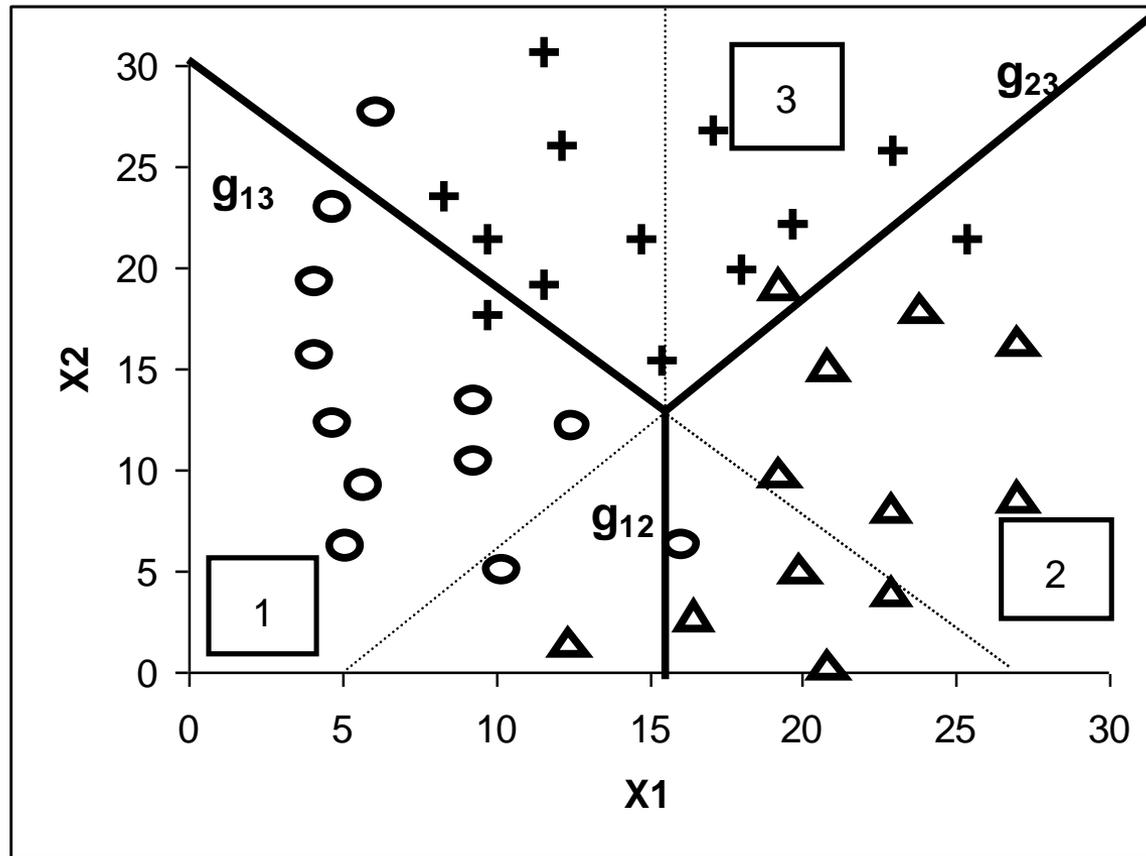


- Propiedad deseable para el diseño de $g_i(\cdot)$: sobre el conjunto de entrenamiento E , cada patrón de la clase C_i tiene un valor máximo con el discriminante $g_i(\cdot)$:

$$g_i(\vec{X}_j^{(i)}) = \max_{k=1, \dots, M} \{g_k(\vec{X}_j^{(i)})\}, \forall j = 1, \dots, n_i$$

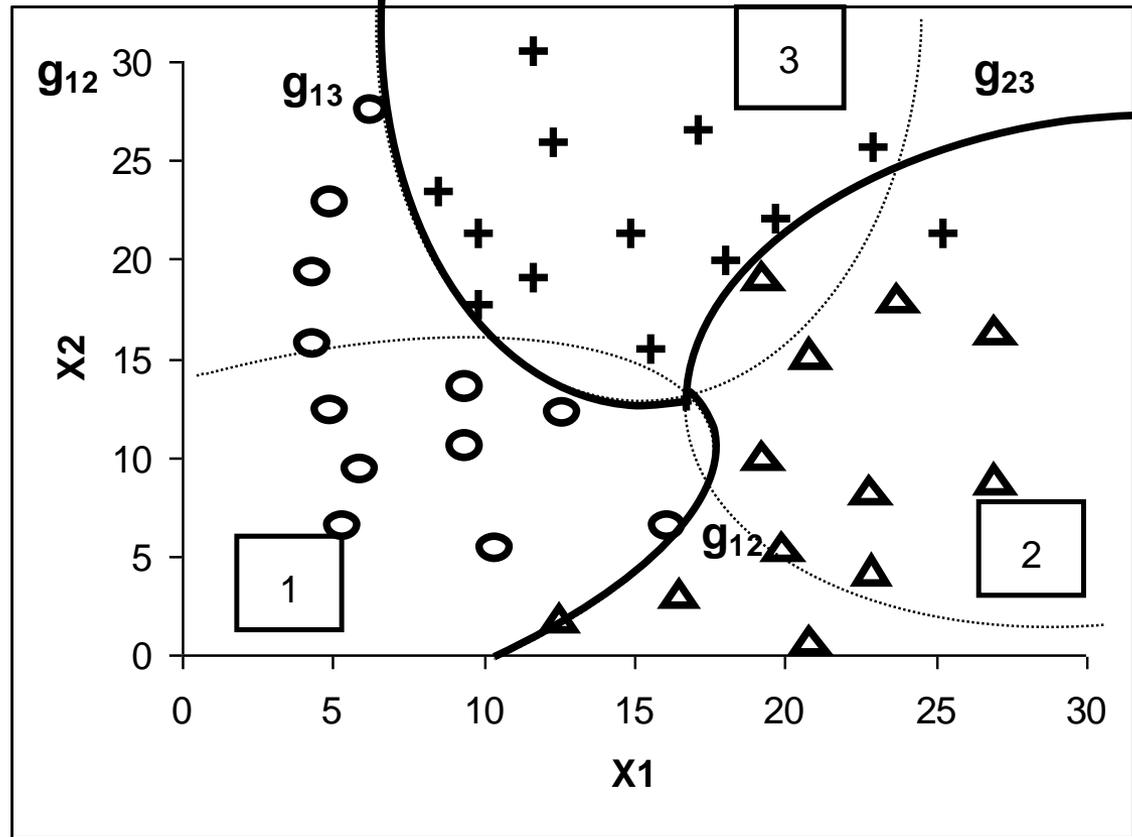
Fronteras de decisión

$g_{ij}(\vec{X})$: lineales



Fronteras de decisión

$g_{ij}(\vec{X})$: cuadraticas



Clasificación con Regresión Lineal

- Para cada clase se define la función de pertenencia g_i :

$$g_i(\vec{X}) = \begin{cases} 1; & \vec{X} \in C_i \\ 0; & \vec{X} \notin C_i \end{cases}$$

- Se construye una función lineal que “aproxime” g_i :

1^s en los patrones

de C_i

0^s en resto

$$\vec{y}_i = \begin{bmatrix} 1 \\ \vdots \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

$$H_i = \begin{bmatrix} 1 & (\vec{X}_1^{(i)})^t \\ \vdots & \vdots \\ 1 & (\vec{X}_{n_i}^{(i)})^t \\ 1 & (\vec{X}_1^{(1)})^t \\ \vdots & \vdots \\ 1 & (\vec{X}_{n_I}^{(I)})^t \end{bmatrix};$$

todos los datos

$$\vec{A}_i = [H_i^t H_i]^{-1} H_i^t \vec{y}_i$$

- Hay que “aprender” M funciones g_i

Clasificación bayesiana: aplicación de modelos estadísticos

- Clasificación con modelo de estructura probabilística conocida

Clases: $\{C_1, \dots, C_M\}$. Se conoce a priori:

- Probabilidades de clase: $P(C_i)$
- Distribuciones de probabilidad condicionadas (parámetros constantes)

$$F_{\vec{X}}(x_1, \dots, x_I | C_i) = P(X_1 \leq x_1, \dots, X_I \leq x_I | C_i) = \frac{P(X_1 \leq x_1, \dots, X_I \leq x_I, C_i)}{P(C_i)}$$

- densidad

$$f_{\vec{X}}(x_1, \dots, x_I | C_i) = \frac{\partial F_{\vec{X}}(x_1, \dots, x_I | C_i)}{\partial x_1 \dots \partial x_I}$$

Ej.: distribución normal multivariada

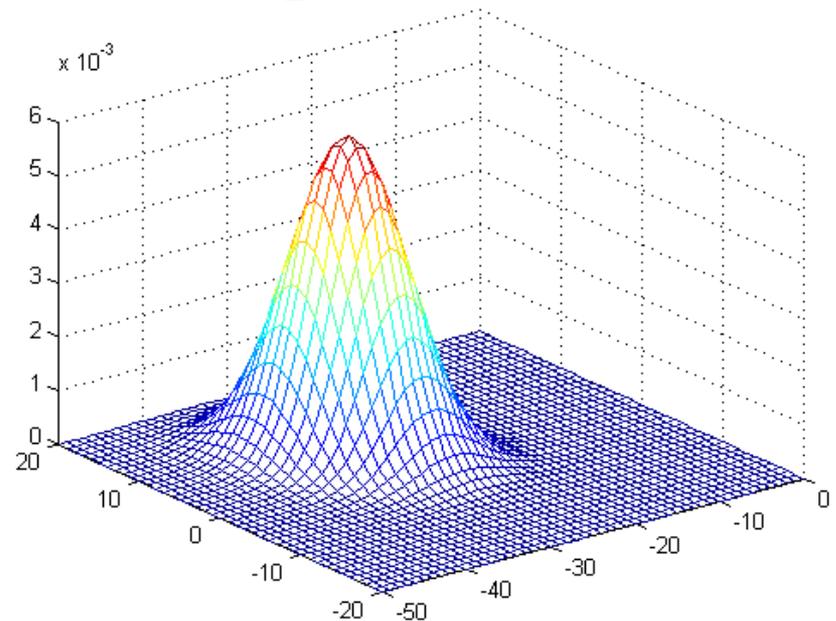
- Parámetros: vector de medias y matriz covarianzas

$$f(\bar{x}) = \frac{1}{(2\pi)^{n/2} \sqrt{|S|}} \exp\left[-\frac{1}{2}(\bar{x} - \bar{\mu})^t S^{-1}(\bar{x} - \bar{\mu})\right]$$

$$\bar{\mu} = \begin{bmatrix} \mu_1 \\ \vdots \\ \mu_n \end{bmatrix}; \quad S = \begin{bmatrix} \sigma_{x_1 x_1}^2 & \sigma_{x_1 x_2} & \cdots & \sigma_{x_1 x_F} \\ \vdots & \vdots & \vdots & \vdots \\ \sigma_{x_F x_1} & \sigma_{x_n x_2} & \cdots & \sigma_{x_F}^2 \end{bmatrix}$$

- Ejemplo

$$\bar{\mu} = \begin{bmatrix} -30 \\ 5 \end{bmatrix}; \quad S = \begin{bmatrix} 21 & -6 \\ -6 & 21 \end{bmatrix}$$



Teorema de Bayes aplicado a clasificación

$$P(C_i | \vec{X}) = \frac{f(\vec{X} | C_i)p(C_i)}{f(\vec{X})}$$

- Probabilidad a posteriori: es la probabilidad de que el ejemplo tenga clase C_i : $P(C_i | \vec{X})$
- Probabilidad a priori: $P(C_i)$ es la probabilidad total de cada clase
- Verosimilitud: $f(\vec{X} | C_i)$: es la distribución de C_i aplicada a \vec{X}
- Densidad total: $f(\vec{X}) = f(\vec{X} | C_1)P(C_1) + \dots + f(\vec{X} | C_M)P(C_M)$

Criterio de clasificación MAP:

$$\text{Clase}(\vec{X}) = \underset{i}{\text{m\u00e1ximo}} \{P(C_i | \vec{X})\} = \underset{i}{\text{m\u00e1ximo}} \{f(\vec{X} | C_i)p(C_i)\}$$

- función discriminante de C_i : proporcional a su prob a posteriori:

$$g_i(\vec{X}) = f(\vec{X} | C_i)p(C_i)$$

- la clase es la de aquella que maximiza el discriminante

Clasificación bayesiana y distrib normal

- Distribuciones condicionales gaussianas. Para cada clase C_i hay una función discriminante de parámetros $\mu_{ij}, \sigma_{ij}, j=1 \dots l$

$$g_i(\bar{x}) = \log(P(C_i)f(\bar{x} | C_i)) = \log \frac{P(C_i)}{(2\pi)^{n/2} \sigma_{1i} \sigma_{2i} \dots \sigma_{Fi}} - \frac{1}{2}(\bar{x} - \bar{\mu})^t S^{-1}(\bar{x} - \bar{\mu})$$

$$\text{simplificaiion : } K - \frac{1}{2} \sum_{i=1}^F (x_j - \mu_{ij})^2 / \sigma_{ij}^2$$

- Parámetros de distribución condicionada a cada clase
- Regiones de decisión:
 - Funciones cuadráticas (hipérbolas) dadas por diferencias:
$$g_{ij}(\bar{x}) = g_i(\bar{x}) - g_j(\bar{x})$$
 - Si son iguales, y diagonales: regiones lineales (caso particular)

Ejemplo con distribución normal

- C1 :

$$P_1 = 0.3; \quad C_1 = [-30 \quad 5]^t; \quad R_1 = \begin{bmatrix} 21 & -6 \\ -6 & 21 \end{bmatrix}$$

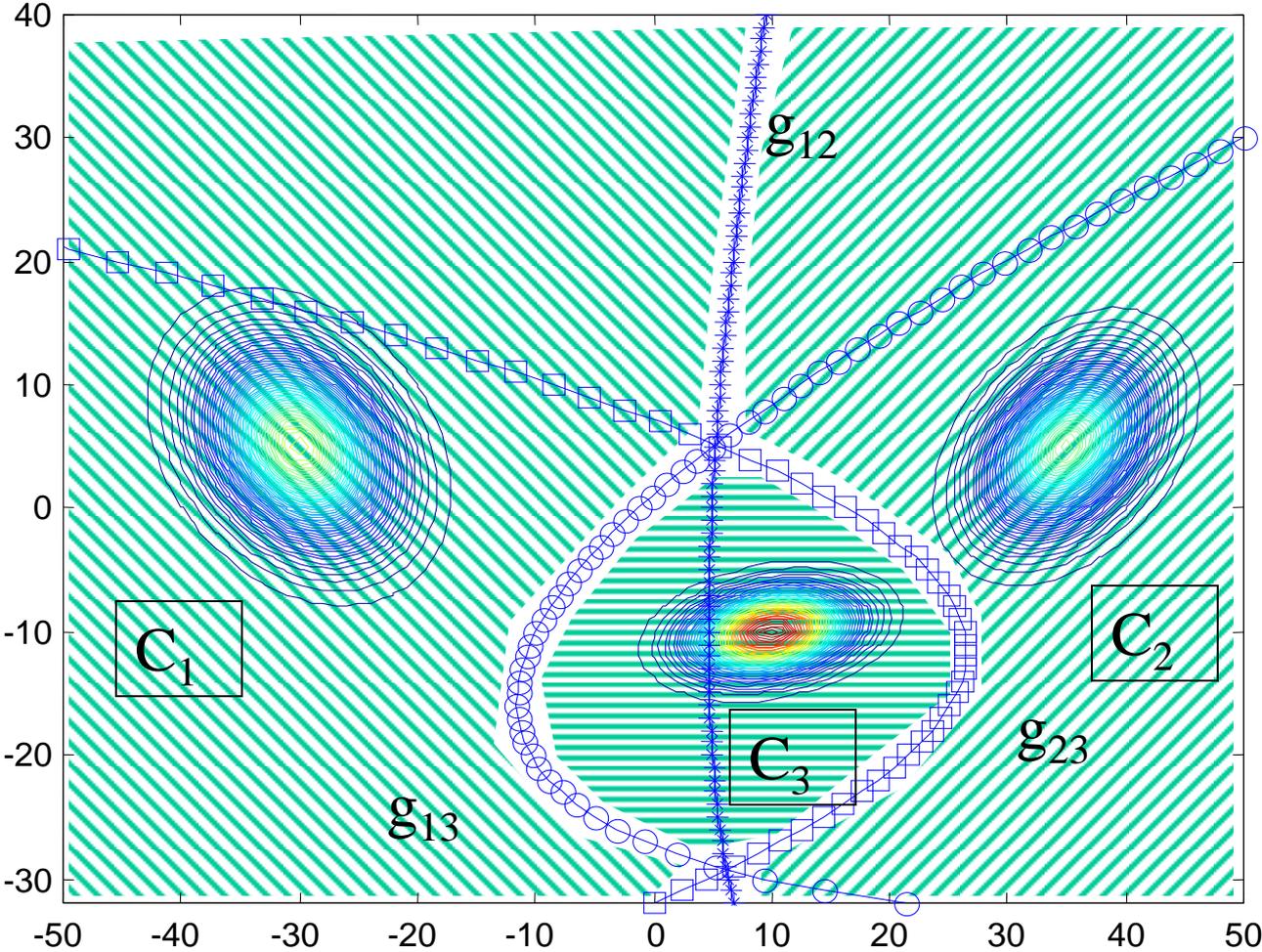
- C2 :

$$P_2 = 0.2; \quad C_2 = [35 \quad 5]^t; \quad R_2 = \begin{bmatrix} 16 & 6 \\ 6 & 16 \end{bmatrix}$$

- C3 :

$$P_3 = 0.5; \quad C_3 = [10 \quad -10]^t; \quad R_3 = \begin{bmatrix} 16 & 2 \\ 2 & 4 \end{bmatrix}$$

Ejemplo



Resumen clasificador bayesiano numérico

- Algoritmo:

- Estimar parámetros de cada clase C_i (entrenamiento)

$$C_i : \{\bar{X}_1^{(i)}, \dots, \bar{X}_{n_i}^{(i)}\} \longrightarrow \bar{\mu}_i, \quad C_i$$

$$\hat{\mu}_i = \frac{1}{n} \sum_{j=1}^{n_i} \bar{x}_j^{(i)}$$

- Estimar probabilidad de cada clase

$$C_i = \frac{1}{n_i} \sum_{j=1}^{n_i} (\bar{x}_i - \mu_i)(\bar{x}_i - \mu_i)^t$$

$$\hat{P}(C_i) = \frac{n_i}{N}; \quad n = \sum_{i=1}^M n_i$$

- Obtener regiones de decisión: $g_{ij}(\cdot)$

Clasificación Bayesiana con Atributos Nominales

Atributos nominales con valores discretos

- $A_i = \{V_1, \dots, V_{n_i}\}$: atributo con n_i valores posibles
- Pasamos de densidades a probabilidades:
probabilidad a priori: $p(A_i = V_j | C_k)$?
- Estimación “contando” el número de casos:

$$p(A_i = V_j | C_k) = \frac{\text{n}^\circ \text{ de ejemplos de clase } C_k \text{ con } A_i = V_j}{\text{n}^\circ \text{ de ejemplos de clase } C_k}$$

Clasificación Bayesiana con Atributos Nominales

- Simplificación: independencia de atributos (“Naive Bayes”):
la probabilidad conjunta de varios atributos se pone como producto

$$X_i = (A_1 = V_1, A_2 = V_2, \dots, A_I = V_I)$$

$$p(X_i | C_k) = p(A_1 = V_1 | C_k) * p(A_2 = V_2 | C_k) * \dots * p(A_I = V_I | C_k)$$

- Clasificación:

$$p(C_k | X_i) = \frac{p(X_i | C_k) * p(C_k)}{p(X_i)} =$$

$$\frac{p(A_1 = V_1 | C_k) * p(A_2 = V_2 | C_k) * \dots * p(A_F = V_F | C_k) * p(C_k)}{p(X_i)}$$

Ejemplo con atributos nominales

SALARIO	CLIENTE	EDAD	HIJOS	CRÉDITO
Poco	Sí	Joven	Uno	NO
Mucho	Si	Joven	Uno	SI
Mucho	Si	Joven	Uno	SI
Poco	Si	Joven	Uno	NO
Mucho	Si	Joven	Dos	SI
Poco	Si	Joven	Dos	NO
Mucho	Si	Adulto	Dos	SI
Mucho	Si	Adulto	Dos	SI
Poco	No	Adulto	Dos	NO
Mucho	Si	Adulto	Dos	SI
Medio	No	Adulto	Tres	NO
Mucho	Si	Adulto	Dos	SI
Medio	Si	Adulto	Dos	SI
Medio	No	Adulto	Tres	NO
Medio	No	Adulto	Dos	SI
Mucho	No	Mayor	Tres	NO
Poco	No	Mayor	Tres	SI
Poco	No	Mayor	Tres	SI
Mucho	No	Mayor	Tres	NO
Mucho	No	Mayor	Tres	SI

$$p(SI) = 12/20$$

$$p(NO) = 8/20$$

	Crédito	No	Sí
Salario			
Poco		4/8	2/12
Mucho		2/8	8/12
Medio		2/8	2/12
Cliente			
Sí		3/8	8/12
No		5/8	4/12
Edad			
Joven		3/8	3/12
Adulto		3/8	6/12
Mayor		2/8	3/12
Hijos			
Uno		2/8	2/12
Dos		2/8	7/12
Tres		4/8	3/12

Ejemplo con atributos nominales

- Ej.: (salario=poco, cliente=si, edad=adulto, hijos=tres)

$$p(\text{SI} | X_i) =$$

$$p(s = \text{poco} | \text{SI}) * p(c = \text{si} | \text{SI}) * p(e = \text{adulto} | \text{SI}) * p(h = \text{tres} | \text{SI}) * p(\text{SI}) / p(X_i) = \\ 2/12 * 8/12 * 6/12 * 3/12 * 12/20 / p(X_i) = 0.0083 / p(X_i)$$

$$p(\text{NO} | X_i) =$$

$$p(s = \text{poco} | \text{NO}) * p(c = \text{si} | \text{NO}) * p(e = \text{adulto} | \text{NO}) * p(h = \text{tres} | \text{NO}) * p(\text{NO}) / p(X_i) = \\ 4/8 * 3/8 * 3/8 * 4/8 * 8/20 / p(X_i) = 0.0141 / p(X_i)$$

Atributos sin valores

- Si el ejemplo a clasificar no tiene un atributo, simplemente se omite.
 - Ej.: (**salario=poco, cliente=si, edad=?, hijos=3**)

$$p(\text{SI} | X_i) =$$

$$p(s = \text{poco} | \text{SI}) * p(c = \text{si} | \text{SI}) * p(h = \text{tres} | \text{SI}) * p(\text{SI}) / p(X_i) =$$

$$2/12 * 8/12 * 3/12 * 12/20 / p(X_i) = 0.0167 / p(X_i)$$

$$p(\text{NO} | X_i) =$$

$$p(s = \text{poco} | \text{NO}) * p(c = \text{si} | \text{NO}) * p(h = \text{tres} | \text{NO}) * p(\text{NO}) / p(X_i) =$$

$$4/8 * 3/8 * 4/8 * 8/20 / p(X_i) = 0.0375 / p(X_i)$$

- Si hay faltas en la muestra de entrenamiento, no cuentan en la estimación de probabilidades de ese atributo

Faltas en atributo EDAD

SALARIO	CLIENTE	EDAD	HIJOS	CRÉDITO
Poco	Sí	Joven	Uno	NO
Mucho	Si	Joven	Uno	SI
Mucho	Si	Joven	Uno	SI
Poco	Si	?	Uno	NO
Mucho	Si	?	Dos	SI
Poco	Si	?	Dos	NO
Mucho	Si	?	Dos	SI
Mucho	Si	Adulto	Dos	SI
Poco	No	Adulto	Dos	NO
Mucho	Si	Adulto	Dos	SI
Medio	No	Adulto	Tres	NO
Mucho	Si	Adulto	Dos	SI
Medio	Si	Adulto	Dos	SI
Medio	No	Adulto	Tres	NO
Medio	No	Adulto	Dos	SI
Mucho	No	Mayor	Tres	NO
Poco	No	Mayor	Tres	SI
Poco	No	Mayor	Tres	SI
Mucho	No	Mayor	Tres	NO
Mucho	No	Mayor	Tres	SI

$$p(SI) = 12/20$$

$$p(NO) = 8/20$$

	Crédito	No	Sí
Salario			
Poco		4/8	2/12
Mucho		2/8	8/12
Medio		2/8	2/12
Cliente			
Sí		3/8	8/12
No		5/8	4/12

	Crédito	No	Sí
Edad			
Joven		1/6	2/10
Adulto		3/6	5/10
Mayor		2/6	3/10

	Crédito	No	Sí
Hijos			
Uno		2/8	2/12
Dos		2/8	7/12
Tres		4/8	3/12

Atributos no representados.

Laplace

- Problema: con muestra poco representativa, puede ocurrir que en alguna clase, un valor de atributo no aparezca: $p(A_i=V_j|C_k)=0$
 - Cualquier ejemplo X con $A_i=V_j$ generará $P(C_k|X)=0$, independientemente de los otros atributos!
- Se suele modificar la estimación de las probabilidades a priori con un factor que elimina los ceros.

- Ej.: $P(\text{Edad}|\text{Crédito}=\text{NO}) = \left\{ \text{Joven: } \frac{3}{8}, \text{ Adulto: } \frac{3}{8}, \text{ Mayor: } \frac{2}{8} \right\}$

- Ley μ : $\left\{ \text{Joven: } \frac{3+\mu/3}{8+\mu}, \text{ Adulto: } \frac{3+\mu/3}{8+\mu}, \text{ Mayor: } \frac{2+\mu/3}{8+\mu} \right\}$

- A veces simplemente se inicializan las cuentas a 1 en vez de 0:

$$\left\{ \text{Joven: } \frac{3+1}{8+3}, \text{ Adulto: } \frac{3+1}{8+3}, \text{ Mayor: } \frac{2+1}{8+3} \right\}$$

Atributos mixtos

- Independencia de atributos (“Naive Bayes”)

$$p(X_i | C_k) = p(A_1 = V_1 | C_k) * p(A_2 = V_2 | C_k) * \dots * p(A_F = V_F | C_k)$$

- Atributos discretos: probabilidades a priori con cada clase C_k

$$p(A_i = V_j | C_k) = \frac{\text{n}^\circ \text{ de ejemplos de clase } C_k \text{ con } A_i = V_j}{\text{n}^\circ \text{ de ejemplos de clase } C_k}$$

- Atributos continuos: densidades de clase C_k : normales de parámetros

μ_k, σ_k

$$p(A_i = V_j | C_k) \rightarrow f_{A_i}(V_j | C_k) = \frac{1}{\sqrt{2\pi}\sigma_{ik}} \exp\left[-\frac{1}{2} \frac{(V_j - \mu_{ik})^2}{\sigma_{ik}^2}\right]$$

Ejemplo con atributos mixtos

SALARIO	CLIENTE	EDAD	HIJOS	CRÉDITO
525	Sí	Joven	1	NO
2000	Si	Joven	1	SI
2500	Si	Joven	1	SI
470	Si	Joven	1	NO
3000	Si	Joven	2	SI
510	Si	Joven	2	NO
2800	Si	Adulto	2	SI
2700	Si	Adulto	2	SI
550	No	Adulto	2	NO
2600	Si	Adulto	2	SI
1100	No	Adulto	3	NO
2300	Si	Adulto	2	SI
1200	Si	Adulto	2	SI
900	No	Adulto	3	NO
800	No	Adulto	2	SI
800	No	Mayor	3	NO
1300	No	Mayor	3	SI
1100	No	Mayor	3	SI
1000	No	Mayor	3	NO
4000	No	Mayor	3	SI

$$p(SI) = 12/20$$

$$p(NO) = 8/20$$

	Crédito	No	Sí
Salario			
Media		732	2192
Desv Estándar		249	942

	Crédito	No	Sí
Cliente			
Sí		3/8	8/12
No		5/8	4/12

	Crédito	No	Sí
Edad			
Joven		3/8	3/12
Adulto		3/8	6/12
Mayor		2/8	3/12

	Crédito	No	Sí
Hijos			
Media		2.25	2.08
Desv Estándar		0.89	0.67

Ejemplo con atributos mixtos

- Ej.: (salario=700, cliente=si, edad=adulto, hijos=3)

$$p(\text{SI} | X_i) =$$

$$\begin{aligned} & f_S(s = 700 | \text{SI}) * p(c = \text{si} | \text{SI}) * p(e = \text{adulto} | \text{SI}) * f_H(h = 3 | \text{SI}) * p(\text{SI}) / p(X_i) = \\ & \frac{1}{\sqrt{2\pi}942} \exp\left[-\frac{1}{2} \frac{(700 - 2192)^2}{942^2}\right] * 8/12 * 6/12 * \frac{1}{\sqrt{2\pi}0.67} \exp\left[-\frac{1}{2} \frac{(3 - 2.08)^2}{0.67^2}\right] * 12/20 * 1/P(X_i) = \\ & = 5.61e - 6 / p(X_i) \end{aligned}$$

$$p(\text{NO} | X_i) =$$

$$\begin{aligned} & f_S(s = 700 | \text{NO}) * p(c = \text{si} | \text{NO}) * p(e = \text{adulto} | \text{NO}) * f_H(h = 3 | \text{NO}) * p(\text{NO}) / p(X_i) = \\ & \frac{1}{\sqrt{2\pi}249} \exp\left[-\frac{1}{2} \frac{(700 - 732)^2}{249^2}\right] * 3/8 * 3/8 * \frac{1}{\sqrt{2\pi}0.89} \exp\left[-\frac{1}{2} \frac{(3 - 2.25)^2}{0.89^2}\right] * 8/20 * 1/P(X_i) = \\ & = 2.81e - 5 / p(X_i) \end{aligned}$$

Clasificación con costes

- MAP proporciona clasificación con mínima prob. de Error
 - Coste de decisión $D_i | \vec{X}$: prob. Error total = $1 - P(C_i | \vec{X})$
- Con frecuencia los costes son asimétricos, y unos errores son más graves que otros. Matriz de costes

$$\begin{array}{c}
 \text{Clase} \\
 \text{real} \\
 \downarrow
 \end{array}
 \begin{array}{c}
 \xrightarrow{\text{Clasificado como}} \\
 \left(\begin{array}{ccc}
 0 & c_{12} & c_{13} \\
 c_{21} & 0 & c_{23} \\
 c_{31} & c_{32} & 0
 \end{array} \right)
 \end{array}$$

- Costes de cada decisión. Criterio de mínimo coste medio: dada una decisión, promedio los costes de cada equivocación y su coste:

$$\text{coste}(D_1 | \vec{X}) = c_{21}p(C_2 | \vec{X}) + c_{31}p(C_3 | \vec{X})$$

$$\text{coste}(D_2 | \vec{X}) = c_{12}p(C_1 | \vec{X}) + c_{32}p(C_3 | \vec{X})$$

$$\text{coste}(D_3 | \vec{X}) = c_{13}p(C_1 | \vec{X}) + c_{23}p(C_2 | \vec{X})$$

Ejemplo de clasificación con costes

- Clasificación de setas con dos atributos, (X, Y) y tres categorías: *Venenosa, Mal sabor, comestible*: {V, MS, C}

Clasificado como

		V	MS	C
<i>Clase real</i>	V	→		
	MS	⎧	⎩	
	C			

$$\begin{pmatrix} 0 & 1000 & 1000 \\ 1 & 0 & 10 \\ 1 & 1 & 0 \end{pmatrix}$$

$$V: \mu_1 = [-5 \quad -5]^t; \quad C_1 = \begin{bmatrix} 71 & -50 \\ -50 & 71 \end{bmatrix}$$

$$C: \mu_2 = [5 \quad 5]^t; \quad C_2 = \begin{bmatrix} 71 & -40 \\ -40 & 71 \end{bmatrix}$$

$$MS: \mu_3 = [20 \quad -20]^t; \quad C_3 = \begin{bmatrix} 51 & 45 \\ 45 & 51 \end{bmatrix}$$

