

Naïve Bayes

DRA. CONSUELO VARINIA GARCÍA MENDOZA

Clasificador Bayesiano

- Enfoque de clasificación a través de un modelo probabilístico
- Dada la probabilidad de que una instancia pertenezca a una clase es posible clasificar nuevas instancias

$$p(y|X)$$

Instance	Probability of class 1	Probability of class 2
1	0.8	0.6
2	0.4	0.7
3	0.6	0.6

Teorema de Bayes

- Describe la probabilidad de un evento (salida estimada), basado en conocimiento a priori de condiciones(características) que podrían estar relacionadas con el evento

$$p(y|X) = p(y) \frac{p(X|y)}{p(X)}$$

donde:

- $p(y|X)$ es la probabilidad condicional. Es la probabilidad de que ocurra y dadas las características X . Probabilidad actualizada después de considerar las pruebas (características)
- $p(y)$ es la probabilidad a priori. La probabilidad (de la salida) antes de considerar las pruebas
- $p(X|y)$ es la verosimilitud. Probabilidad de la evidencia (las características), dado que la creencia (salida) es verdadera
- $p(X)$ es la probabilidad marginal. La probabilidad de la evidencia (características), bajo cualquier circunstancia

Ejemplo

Instance	Features (X)	Class (y)
	Size	
1	big	pos
2	small	pos
3	small	pos
4	small	pos
5	small	neg
6	big	neg
7	big	neg
8	big	neg
9	small	pos
10	big	pos

- Teorema de Bayes

$$p(y|X) = p(y) \frac{p(X|y)}{p(X)}$$

- Probabilidad a priori $p(y)$

- $p(pos) = \frac{N_{pos}}{N_{all}} = \frac{6}{10}$

- $p(neg) = \frac{N_{neg}}{N_{all}} =$

Instance	Features (X)	Class (y)
	Size	
1	big	pos
2	small	pos
3	small	pos
4	small	pos
5	small	neg
6	big	neg
7	big	neg
8	big	neg
9	small	pos
10	big	pos

- Teorema de Bayes

$$p(y|X) = p(y) \frac{p(X|y)}{p(X)}$$

- Probabilidad marginal $p(X)$

- $p(big) = \frac{N_{big}}{N_{all}} = \frac{5}{10}$

- $p(small) = \frac{N_{small}}{N_{all}} =$

Instance	Features (X)	Class (y)
	Size	
1	big	pos
2	small	pos
3	small	pos
4	small	pos
5	small	neg
6	big	neg
7	big	neg
8	big	neg
9	small	pos
10	big	pos

- Teorema de Bayes

$$p(y|X) = p(y) \frac{p(X|y)}{p(X)}$$

- Verosimilitud $p(X|y)$

- $p(big|pos) = \frac{N_{big \cap pos}}{N_{pos}} = \frac{2}{6}$

- $p(small|pos) = \frac{N_{small \cap pos}}{N_{pos}} =$

- $p(big|neg) = \frac{N_{big \cap neg}}{N_{neg}} =$

- $p(small|neg) = \frac{N_{small \cap neg}}{N_{neg}} =$

- $p(pos) = 0.6$
- $p(neg) = 0.4$
- $p(big) = 0.5$
- $p(small) = 0.5$
- $p(big|pos) = 0.33$
- $p(big|neg) = 0.75$
- $p(small|pos) = 0.67$
- $p(small|neg) = 0.25$

Probabilidad a priori

- $p(pos|big) = p(pos) \frac{p(big|pos)}{p(big)} = 0.6 \cdot \frac{0.33}{0.5} = 0.4$
- $p(neg|big) = p(neg) \frac{p(big|neg)}{p(big)} =$
- $p(pos|small) = p(pos) \frac{p(small|pos)}{p(small)} =$
- $p(neg|small) = p(neg) \frac{p(small|neg)}{p(small)} =$

Teorema de Bayes (probabilidad a posteriori)

$$p(y|X) = p(y) \frac{p(X|y)}{p(X)} \rightarrow p(y)p(X|y)$$

Teorema de Bayes (probabilidad a posteriori)

$$p(y|X) = p(y) \frac{p(X|y)}{p(X)}$$

- $p(pos|big) = p(pos) \frac{p(big|pos)}{p(big)} = 0.4$
- $p(neg|big) = p(neg) \frac{p(big|neg)}{p(big)} = 0.6$
- $p(pos|small) = p(pos) \frac{p(small|pos)}{p(small)} = 0.8$
- $p(neg|small) = p(neg) \frac{p(small|neg)}{p(small)} = 0.2$

$$p(y|X) = p(y)p(X|y)$$

- $p(pos|big) = 0.2$
- $p(neg|big) = 0.3$
- $p(pos|small) = 0.4$
- $p(neg|small) = 0.1$

Hipótesis del Ingenuo Bayes

Teorema de Bayes

$$p(y|X) = p(y)p(X|y)$$

Cuando las instancias tienen más de una característica, cada característica contribuye a la clasificación

Para considerar la contribución de cada característica aplicamos la probabilidad conjunta

$$p(X|y_j) = \prod_{i=1}^n p(x_i|y_j)$$

El teorema de Bayes para n características a través de la probabilidad conjunta se expresa como

$$p(y_j|X) = p(y_j) \prod_{i=1}^n p(x_i|y_j)$$

Para aplicar esta regla, debemos suponer que cada atributo es mutuamente independiente

Esta suposición no suele estar justificada, por lo que se denomina a este clasificador como ingenuo Bayes

Ejemplo 2

Instance	Features (X)		Class (y)
	Size	Color	
1	big	green	pos
2	small	red	pos
3	small	red	pos
4	small	red	pos
5	small	red	neg
6	big	red	neg
7	big	green	neg
8	big	green	neg
9	small	green	pos
10	big	red	pos

- $p(pos) = 0.6$
- $p(neg) = 0.4$
- $p(big|pos) = 0.33$
- $p(big|neg) = 0.75$
- $p(small|pos) = 0.67$
- $p(small|neg) = 0.25$

- $p(red|pos) = \frac{N_{red \cap pos}}{N_{pos}} =$
- $p(green|pos) = \frac{N_{green \cap pos}}{N_{pos}} =$
- $p(red|neg) = \frac{N_{red \cap neg}}{N_{neg}} =$
- $p(green|neg) = \frac{N_{green \cap neg}}{N_{neg}} =$

$$p(y_j|X) = p(y_j) \prod_{i=1}^n p(x_i|y_j)$$

- $p(pos|big, red) = p(pos) \cdot p(big|pos) \cdot p(red|pos) =$
- $p(pos|big, green) = p(pos) \cdot p(big|pos) \cdot p(green|pos) =$
- $p(pos|small, red) = p(pos) \cdot p(small|pos) \cdot p(red|pos) =$
- $p(pos|small, green) = p(pos) \cdot p(small|pos) \cdot p(green|pos) =$
- $p(neg|big, red) = p(neg) \cdot p(big|neg) \cdot p(red|neg) =$
- $p(neg|big, green) = p(neg) \cdot p(big|neg) \cdot p(green|neg) =$
- $p(neg|small, red) = p(neg) \cdot p(small|neg) \cdot p(red|neg) =$
- $p(neg|small, green) = p(neg) \cdot p(small|neg) \cdot p(green|neg) =$

Distribuciones de probabilidad

Se pueden hacer distintas implementaciones del clasificador Naive considerando diversas distribuciones de probabilidad

- Bernoulli
- Multinomial
- Gaussiana

Distribución de Bernoulli

Una variable aleatoria Bernoulli es una variable aleatoria que sólo puede tomar dos valores posibles, normalmente 0 y 1

Esta variable aleatoria modela experimentos aleatorios que tienen dos posibles resultados, a veces denominados "éxito" y "fracaso"

La variante Bernoulli para el clasificador Naïve Bayes se utiliza cuando las características toman valores binarios o booleanos

- Género (Hombre = 1 o Mujer = 0)
- Tarjeta de crédito (Sí = 1 o No = 0)
- Vivienda propia (Sí = 1 o No = 0)

Distribución multinomial

La distribución multinomial se utiliza para encontrar probabilidades en experimentos en los que hay más de dos resultados

La variante multinomial para el clasificador Naive Bayes se utiliza cuando las características describen conteos de frecuencia discretos

- Calificaciones de películas (Número de estrellas 1-5)
- Conteo de palabras (número de veces que aparece la palabra "bueno" en la reseña de un producto)

Distribución Gaussiana

La distribución normal o gaussiana es un tipo de distribución de probabilidad continua para una variable aleatoria de valor real

La variante gaussiana para el clasificador Naive Bayes se utiliza cuando las características son continuas

- Peso
- Altura
- Presión arterial

Clasificación de textos

Dado un grupo de documentos, un modelo debe predecir la categoría a la que pertenece cada documento

Dataset	Document	Words	Class
Training	1	Chinese Beijing Chinese	c
	2	Chinese Chinese Shangai	c
	3	Chinese Macao	c
	4	Tokyo Japan Chinese	j
Test	5	Chinese Chinese Chinese Tokyo Japan	?

¿Qué variante de Naive Bayes debemos utilizar?

El teorema de Bayes para n características a través de la probabilidad conjunta se expresa como

$$p(y_j|X) = p(y_j) \prod_{i=1}^n p(x_i|y_j)$$

Distribución multinomial

$$p(w_i|y_j) = \frac{\text{count}(w_i, y_j) + \alpha}{\sum_{w \in V} [\text{count}(w, y_j) + \alpha]} = \frac{\text{count}(w_i, y_j) + \alpha}{(\sum_{w \in V} \text{count}(w, y_j)) + |V|}$$

donde

w_i : palabra i del vocabulario

y_j : clase j

V : vocabulario

α : suavizado de Laplace (Laplace smoothing)

Bayes y la distribución multinomial

$$p(y_j|d_k) = p(y_j) \prod_{i=1}^n p(w_i|y_j)$$

Dataset	Document	Words	Class
Training	1	Chinese Beijing Chinese	c
	2	Chinese Chinese Shangai	c
	3	Chinese Macao	c
	4	Tokyo Japan Chinese	j
Test	5	Chinese Chinese Chinese Tokyo Japan	?

$p(y_j)$

- $p(c) = \frac{N_c}{N_{all}} =$

- $p(j) = \frac{N_j}{N_{all}} =$

$$p(w_i|y_j) = \frac{\text{count}(w_i, y_j) + \alpha}{(\sum_{w \in V} \text{count}(w, y_j)) + |V|}$$

Dataset	Document	Words	Class
Training	1	Chinese Beijing Chinese	c
	2	Chinese Chinese Shangai	c
	3	Chinese Macao	c
	4	Tokyo Japan Chinese	j
Test	5	Chinese Chinese Chinese Tokyo Japan	?

$$p(\text{Chinese}|c) = \frac{\text{count}(\text{Chinese}, c) + \alpha}{\text{count}(w, c) + 6} = \frac{5+1}{8+6} = \frac{6}{14} = \frac{3}{7}$$

$$p(\text{Beijing}|c) = \frac{\text{count}(\text{Beijing}, c) + \alpha}{\text{count}(w, c) + 6} =$$

$$p(\text{Shangai}|c) = \frac{\text{count}(\text{Shangai}, c) + \alpha}{\text{count}(w, c) + 6} =$$

$$p(\text{Macao}|c) = \frac{\text{count}(\text{Macao}, c) + \alpha}{\text{count}(w, c) + 6} =$$

$$p(\text{Tokio}|c) = \frac{\text{count}(\text{Tokio}, c) + \alpha}{\text{count}(w, c) + 6} =$$

$$p(\text{Japan}|c) = \frac{\text{count}(\text{Japan}, c) + \alpha}{\text{count}(w, c) + 6} =$$

$$p(w_i|y_j) = \frac{\text{count}(w_i, y_j) + \alpha}{(\sum_{w \in V} \text{count}(w, y_j)) + |V|}$$

Dataset	Document	Words	Class
Training	1	Chinese Beijing Chinese	c
	2	Chinese Chinese Shangai	c
	3	Chinese Macao	c
	4	Tokyo Japan Chinese	j
Test	5	Chinese Chinese Chinese Tokyo Japan	?

$$p(\text{Chinese}|j) = \frac{\text{count}(\text{Chinese}, j) + \alpha}{\text{count}(w, j) + 6} = \frac{1+1}{3+6} = \frac{2}{9}$$

$$p(\text{Beijing}|j) = \frac{\text{count}(\text{Beijing}, j) + \alpha}{\text{count}(w, j) + 6} = \frac{0+1}{3+6} = \frac{1}{9}$$

$$p(\text{Shangai}|j) = \frac{\text{count}(\text{Shangai}, j) + \alpha}{\text{count}(w, j) + 6} =$$

$$p(\text{Macao}|j) = \frac{\text{count}(\text{Macao}, j) + \alpha}{\text{count}(w, j) + 6} =$$

$$p(\text{Tokio}|j) = \frac{\text{count}(\text{Tokio}, j) + \alpha}{\text{count}(w, j) + 6} =$$

$$p(\text{Japan}|j) = \frac{\text{count}(\text{Japan}, j) + \alpha}{\text{count}(w, j) + 6} =$$

$$p(c) = \frac{3}{4}$$

$$p(j) = \frac{1}{4}$$

$$p(\text{Chinese}|c) = \frac{3}{7}$$

$$p(\text{Chinese}|j) = \frac{2}{9}$$

$$p(\text{Beijing}|c) = \frac{1}{7}$$

$$p(\text{Beijing}|j) = \frac{1}{9}$$

$$p(\text{Shangai}|c) = \frac{1}{7}$$

$$p(\text{Shangai}|j) = \frac{1}{9}$$

$$p(\text{Macao}|c) = \frac{1}{7}$$

$$p(\text{Macao}|j) = \frac{1}{9}$$

$$p(\text{Tokio}|c) = \frac{1}{14}$$

$$p(\text{Tokio}|j) = \frac{2}{9}$$

$$p(\text{Japan}|c) = \frac{1}{14}$$

$$p(\text{Japan}|j) = \frac{2}{9}$$

Bayes y la distribución multinomial

$$p(y_j|d_k) = p(y_j) \prod_{i=1}^n p(w_i|y_j)$$

Dataset	Document	Words	Class
Training	1	Chinese Beijing Chinese	c
	2	Chinese Chinese Shangai	c
	3	Chinese Macao	c
	4	Tokyo Japan Chinese	j
Test	5	Chinese Chinese Chinese Tokyo Japan	?

$$p(c|d_5) = \frac{3}{4} \cdot \frac{3}{7} \cdot \frac{3}{7} \cdot \frac{3}{7} \cdot \frac{1}{14} \cdot \frac{1}{14} = 0.0003$$

$$p(j|d_5) = \frac{1}{4} \cdot \frac{2}{9} \cdot \frac{2}{9} \cdot \frac{2}{9} \cdot \frac{2}{9} \cdot \frac{2}{9} = 0.0001$$

Logaritmo es una función monótona, lo que significa que si $a > b$ es equivalente a $\log a > \log b$

Teorema de Bayes

$$p(y_j|d_k) = p(y_j) \prod_{i=1}^n p(x_i|y_j)$$

$$p(y_j|d_k) = \log p(y_j) + \sum_{i=1}^n \log p(x_i|y_j)$$