



tema 10

Aprendizaje bayesiano

la probabilidad de un evento se puede aproximar por la frecuencia de que ocurra (estadística pero para bayesianos no sirve)

Yo tengo una idea en la cabeza pero los datos luego me darán la razón o no

¿cuántos días va a llover este año?

70% porque lo digo yo, y luego los datos me harán cambiar de opinión

tienen la ventaja de que necesita muy pocos datos

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

si es difícil calcularlo de esta forma, puede que al darle la vuelta sea más fácil de calcularlo (el 95% de los casos)

Prior no informativa, asigna la misma probabilidad a cada uno de los posibles conjuntos de parámetros

¿Cuál es la probabilidad de unos parámetros conocidos unos datos (esto es lo que vamos a calcular)?

Importante, pregunta segura (el nombre de estas)

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

$$p(X|\theta)$$

$p(X)$ is called the *likelihood function*

$p(\theta)$ is called the *marginal distribution of data*

as mentioned, it is called the *prior*

- funcion de verosimilitud
- distribución marginal de los datos
- Probabilidad a priori
- Probabilidad a posteriori

$$p(\theta|X)$$

Ahora pillo datos de un bd y los pongo para ver como cambia (en el numerador)

$p(X)$ suele ser difícil de calcular entonces nos la vamos a cargar

no me importa el número en concreto de una probabilidad, me interesa saber si con un conjunto de datos da una mayor probabilidad que dado otro conjunto sobre los mismos datos

$$p(\theta_1|X) = \frac{p(X|\theta_1)p(\theta_1)}{p(X)} \quad \text{vs.} \quad p(\theta_2|X) = \frac{p(X|\theta_2)p(\theta_2)}{p(X)}$$

miro si lo de la izquierda es mayor que lo de la derecha, entonces, me da igual el denominador

procedimiento máximo a posteriori (MAP)

$$\theta_{MAP} = \underset{\theta}{\operatorname{argmax}} p(\theta|X) = \underset{\theta}{\operatorname{argmax}} p(X|\theta)p(\theta)$$

(no tiene porque dar la suma de estos 1)

naive bayes

independencia condicional dada la clase

$$p(X_1, X_2|Y) = p(X_1|Y)p(X_2|Y)$$

esto es lo que vamos a utilizar (que no es lo mismo que seria en estadística)

ejemplo

the following data

long	sweet	yellow	fruit	TOTAL
400	350	450	banana	500
0	150	300	orange	300
100	150	50	other	200

- Assume that we find one piece of fruit is long, yellow and sweet, what kind of fruit is it?

- Note that we have to compute:

$$p(\text{banana} | \text{long, sweet, yellow})$$

$$p(\text{orange} | \text{long, sweet, yellow})$$

$$p(\text{other} | \text{long, sweet, yellow})$$

15

- We have $p(\text{banana}|\text{long, sweet, yellow}) =$

$$\frac{p(\text{long}|\text{banana})p(\text{sweet}|\text{banana})p(\text{yellow}|\text{banana})p(\text{banana})}{p(\text{long})p(\text{sweet})p(\text{yellow})} = \frac{0.8 \times 0.7 \times 0.9 \times 0.5}{0.5 \times 0.65 \times 0.8} = 0.96$$

$$p(\text{orange}|\text{long, sweet, yellow}) =$$

$$\frac{p(\text{long}|\text{orange})p(\text{sweet}|\text{orange})p(\text{yellow}|\text{orange})p(\text{orange})}{p(\text{long})p(\text{sweet})p(\text{yellow})} = \frac{0 \times 0.7 \times 0.9 \times 0.3}{0.5 \times 0.65 \times 0.8} = 0$$

$$p(\text{other}|\text{long, sweet, yellow}) =$$

$$\frac{p(\text{long}|\text{other})p(\text{sweet}|\text{other})p(\text{yellow}|\text{other})p(\text{other})}{p(\text{long})p(\text{sweet})p(\text{yellow})} = \frac{0.5 \times 0.75 \times 0.25 \times 0.2}{0.5 \times 0.65 \times 0.8} = 0.07$$

16

da igual los denominadores y no deberian estar puestos

- We have $p(\text{banana}|\text{long, sweet, yellow}) =$

$$\frac{p(\text{long}|\text{banana})p(\text{sweet}|\text{banana})p(\text{yellow}|\text{banana})p(\text{banana})}{p(\text{long})p(\text{sweet})p(\text{yellow})} =$$

$$0.8 \times 0.7 \times 0.9 \times 0.5$$

calculo esto para todos y me quedo con el que tenga mayor probabilidad

ventajas de este método, reduce muchísimo la carga computacional

Tenemos el problema de que si una probabilidad es 0, todo va a ser 0

no tengo naranjas largas ya que tengo una muestra muy pequeña que no tiene ese caso

Esto es el problema de 0 probabilidad

solucion:

Laplacian correction

- For example, assume that $X=\{x_1, x_2\}$, $\theta=\{\theta_1, \theta_2\}$ and we have the sample

X	θ	$p(x, \theta)$
1	1	$p(1,1)=3/6$
1	2	$p(1,2)=1/6$
2	1	$p(2,1)=2/6$
2	1	$p(2,2)=0$
1	1	$p(\theta=1)=5/6$
1	1	$p(\theta=2)=1/6$

23

me invento datos pero de forma justa

- Alternatively, using *Laplace smoothing* ($k=1$):

X	θ	
1	1	$p(1,1)=4/10$
1	2	$p(1,2)=2/10$
2	1	$p(2,1)=3/10$
2	1	$p(2,2)=1/10$
1	1	$p(X=1)=6/10$
1	1	$p(X=2)=4/10$
1	1	
1	2	
2	1	
2	2	

ventajas

- Advantages:
 1. Easy to understand
 2. Can be applied even with small datasets
 3. If conditional independence is true, Naive Bayes is the BEST classifier one may build
 4. Excellent performance in most of the cases
 5. Computationally very efficient (*eager*)
 6. Insensitive to irrelevant features
 7. Robust to outliers (particularly in the discrete case)
 8. Robust to missing data
 9. Compared to other algorithms it does not need to be finely “tuned”, there are no hiper parameters (except the “prior”)
 10. New data can easily be incorporated without changing the structure of the model

26

desventajas

- Disadvantages:
 1. Zero probability problem requires artificially modifying the distribution
 2. The performance degrades when there is a high interaction among the variables
 3. It does not provide “true” estimated probabilities
 4. For continuous variables one needs to either use some model or discretize the variable
 5. Can not be used in regression problems
 6. It can not learn interactions between features
 7. Different priors will produce different models