

Probabilistic Models

In real world, there are a lot of scenarios where certainty of something is not confirmed.

In order to represent uncertain knowledge where we are not sure about the predicates, we need uncertain reasoning or probabilistic reasoning

When there are ① unpredictable outcomes

② too large probabilities

③ unknown errors / events

In order to tackle uncertainty, probabilistic reasoning is used.

Probabilistic models represent the relationships between events in a probabilistic way

Then inference uses the probability to take decisions.

Few exs of Probabilistic models are BBN, Naive Bayes

Naive Bayes Classifier

Bayesian classifier works on probability
gives probabilistic prediction

(Unlike SVM & KNN which are based on distance, Bayesian & Decision tree are based on probability)

Handles independent events only [Assumption of independent events]

Used to make Bayesian belief networks

Very simple Model

Useful for spam filtering, sentiment analysis and medical diagnosis

Works on categorical data (Unlike SVM)

Used for classification

Called as Naive classifier because

- ① Very simple
- ② Empirically useful
- ③ Scales very well

Assumption \rightarrow all variables are independent

\rightarrow Dependant events \rightarrow H & T of coin

if there is Head, it means Tails is not there definitely

\rightarrow Independent events \rightarrow 2 coins tossed simultaneously

If one has head, we don't know what is the event of first coin

"Naive" because it considers events independent

In actual sense, even if the events are dependant Naive Bayes assumes that the events are independent

Bayesian classifier principle \rightarrow

"If it walks like a duck, quacks like a duck, then it probably is a duck"

In many applications, the relationship between the attributes set and class is **non-deterministic**.

This means that a test cannot be classified to a class label with certainty

eg age income Bay

Young	High	Yes	— ①	} same tuple different outcomes
Young	Low	No		
Young	High	No	— ②	
Young	High	Yes	— ③	
Old	High	Yes		

∴ No inference rules can be drawn.

It is like two points are coinciding with one another in an SVM

But, we can probabilistically say that given tuple young high, the probability

$$\text{Yes} : \frac{2}{3} \qquad \text{No} : \frac{1}{3}$$

Bayes classifier is an approach for modelling such **probabilistic relationships**.

Example Training data →

id	age	income	student	crediting	buy Car
1	Youth	high	no	fair	no
2	y	h	no	excellent	no
3	m-a	h	no	f	y
4	S	medium	no	f	y
5	S	low	yes	f	y
6	S	L	yes	e	n
7	m-a	L	yes	e	y
8	y	m	no	f	n
9	y	L	y	f	y
10	S	m	y	f	y
11	y	m	y	e	y
12	ma	m	n	e	y
13	ma	h	y	f	y
14	S	m	n	e	n

A young student with fair credit score,
has medium income, will he buy car?

Here income & age might be dependant
but we assume them to be
independant for naive bayes

* H - Hypothesis (eg buy or not buy)

* x - evidence (data type)

* prior probability $P(H)$ \rightarrow knowledge is advanced. that is used

Example buy computer event has a probability 0.7 is known in advanced by historical data

* Posterior probability $P(H|x)$

ie posterior probability of H conditioned on x

* Goal \rightarrow Calculate $P(H|x)$

$$P(H|x) = P(H) \frac{P(x|H)}{P(x)}$$

Depending on values of income, age, etc will student buy the laptop?

* Posterior probability = Prior probability \times likelihood

* likelihood $P(x|H)$

Since variables are independant

$$P(x|H) = \prod P(x_i|H)$$

* Marginal probability $P(x)$

It is the overall probability of observing feature vector x

eg $x = [\text{age} = y, \text{income} = m, \dots]$

$P(x)$ is constant for all H , and since we only want to compare,

$$P(H|x) \propto P(H) P(x|H)$$

We need to calculate $P(x)$ & divide

We want to find out the probability that the student buys computer given age = youth, income = medium, etc

ie $P(H = \text{yes} \mid \text{age} = y, \text{income} = m, \text{student} = s, \text{creditscore} = f)$

$$P(H \mid X) = \frac{P(H)}{P(X)} P(X \mid H)$$

(From Bayes theorem)

Since we assume age, income, student etc to be independent variables, we can use the Bayes theorem to split it as

$$P(X \mid H) = \prod_{i=1}^n P(x_i \mid H)$$

Here H is the output ie buy computer = true
 X is the parameters in testing tuple
 ie $X = [\text{age} = y, \text{income} = m, \dots]$

$$\therefore P(H \mid X) = P(H = \text{yes}) \times P(\text{age} = y \mid H = \text{yes}) \times P(\text{income} = m \mid H = \text{yes}) \dots$$

In order to predict if computer buy or not, we calculate $P(H = \text{yes} \mid X)$ & $P(H = \text{no} \mid X)$ and compare them

$P(H) \rightarrow$ Probability student buys computer

$$P(H = \text{yes}) = \text{yes} : 9/14$$

$$P(H = \text{no}) = \text{no} : 5/14$$

$P(X_1 \mid H) \rightarrow$ Probability that the age is youth when computer is bought

$$P(\text{age} = \text{youth} \mid H = \text{yes}) = \frac{2}{9}$$

$$P(\text{age} = \text{youth} \mid H = \text{no}) = \frac{3}{5}$$

Here we considered age = youth because it is in our testing tuple

$$P(X_2 \mid H)$$

$$P(\text{income} = \text{medium} \mid H = \text{yes}) = \frac{4}{9}$$

$$P(\text{income} = \text{medium} \mid H = \text{no}) = \frac{2}{5}$$

$$P(X_3 \mid H)$$

$$P(\text{student} = \text{yes} \mid H = \text{yes}) = 6/9$$

$$P(\text{student} = \text{yes} \mid H = \text{no}) = 1/5$$

$$P(X_4 \mid H)$$

$$P(\text{creditscore} = \text{fair} \mid H = \text{yes}) = 6/9$$

$$P(\text{creditscore} = \text{fair} \mid H = \text{no}) = 2/5$$

Using these probabilities,

$P(X \mid \text{buy computer} = \text{yes})$ is the probability of a customer being a youth student with medium income

$$P(X \mid H = \text{yes}) = P(X_1 \mid H = \text{yes}) \times P(X_2 \mid H = \text{yes}) \times P(X_3 \mid H = \text{yes}) \times P(X_4 \mid H = \text{yes})$$

$$= \frac{2}{9} \times \frac{4}{9} \times \frac{6}{9} \times \frac{6}{9}$$

$$= 0.044$$

$$P(X \mid H = \text{no}) = \frac{3}{5} \times \frac{2}{5} \times \frac{1}{5} \times \frac{2}{5} = 0.019$$

$$P(H = \text{yes} \mid X) = \frac{P(X \mid H = \text{yes}) P(H = \text{yes})}{P(X)} = \frac{0.044 \times \frac{9}{14}}{\frac{1}{14}} = \frac{0.028}{P(X)}$$

$$P(H = \text{no} \mid X) = \frac{P(X \mid H = \text{no}) P(H = \text{no})}{P(X)} = \frac{0.019 \times \frac{5}{14}}{\frac{1}{14}} = \frac{0.006}{P(X)}$$

$$P(H = \text{yes}) > P(H = \text{no})$$

\therefore The tuple can buy computer

$P(X)$ is same for all H (both buy & not buy), so we can ignore it without calculation

Q2

Color	Cloth	Brand	On sale
Blue	Jeans	Y	Yes
Blue	Shirt	X	No
Blue	J	Y	Y
Black	J	Y	Y
Black	S	X	N
Black	S	X	Y
Black	S	Y	Y
Blue	J	Y	N
Black	J	X	Y
Blue	S	X	Y

Classify tuple Blue, jeans brand X would be on sale

$$\rightarrow P(H = \text{yes}) = P(\text{Onsale}) = \frac{4}{10}$$

$$P(H = \text{No}) = \frac{3}{10}$$

$$P(\text{Color} = \text{Blue} \mid H = \text{yes}) = \frac{3}{7}$$

$$P(\text{Color} = \text{Blue} \mid H = \text{No}) = \frac{2}{3}$$

$$P(\text{cloth} = \text{Jeans} \mid H = \text{yes}) = \frac{4}{7}$$

$$P(\text{cloth} = \text{Jeans} \mid H = \text{No}) = \frac{1}{3}$$

$$P(\text{brand} = X \mid H = \text{yes}) = \frac{3}{7}$$

$$P(\text{brand} = X \mid H = \text{No}) = \frac{2}{3}$$

$$P(H = \text{yes} \mid X, B, J) = P(H = \text{yes}) \times P(\text{Color} = \text{Blue} \mid H = \text{yes}) \times \dots$$

$$= \frac{4}{10} \times \frac{3}{7} \times \frac{4}{7} \times \frac{3}{7} = 0.073$$

$$P(H = \text{No} \mid X, B, J) = P(H = \text{No}) \times P(\text{Color} = \text{blue} \mid H = \text{No}) \times \dots$$

$$= \frac{3}{10} \times \frac{2}{3} \times \frac{1}{3} \times \frac{2}{3} = 0.044$$

$$P(H = \text{yes}) > P(H = \text{No})$$

Since probability of Onsale is higher, we can say it will be on sale

Multinomial Naive Bayes

Used for numeric data

$n(\text{free})$	$n(\text{offer})$	$n(\text{money})$	email
3	0	1	spam
0	1	0	Not spam
2	1	1	Spam
0	1	0	Not spam
4	2	1	spam
0	0	0	Not spam

find class where $\left. \begin{array}{l} \text{"free"} : 2 \text{ times} \\ \text{"offer"} : 1 \text{ time} \\ \text{"money"} : 0 \text{ time} \end{array} \right\} \text{string occurrence}$

$$\rightarrow P(\text{spam}) = \frac{3}{6} = 0.5$$

$$P(\text{not spam}) = \frac{3}{6} = 0.5$$

Vocabulary size = 3 (3 features)

$$\text{Count}(\text{free}) \ \& \ \text{spam} = 9$$

$$\text{Count}(\text{offer}) \ \& \ \text{spam} = 3$$

$$\text{Count}(\text{money}) \ \& \ \text{spam} = 3$$

$$\text{Total words} = 15$$

$$P(\text{"free"} | \text{spam}) = \frac{9}{15}$$

$$P(\text{"offer"} | \text{spam}) = \frac{3}{15}$$

$$P(\text{"money"} | \text{spam}) = \frac{3}{15}$$

$$P(\text{spam} | \begin{array}{l} \text{free} : 2 \\ \text{offer} : 1 \\ \text{money} : 0 \end{array}) = P(\text{free} | \text{spam})^2 \times P(\text{offer} | \text{spam})^1 \times P(\text{money} | \text{spam})^0 \times P(\begin{array}{l} \text{free} : 2 \\ \text{offer} : 1 \\ \text{money} : 0 \end{array})$$

$$= \left(\frac{9}{15}\right)^2 \times \frac{3}{15} \times P(\begin{array}{l} \text{free} : 2 \\ \text{offer} : 1 \\ \text{money} : 0 \end{array})$$

$$= 0.072 \times P(\begin{array}{l} \text{free} : 2 \\ \text{offer} : 1 \\ \text{money} : 0 \end{array})$$

Similarly for non spam

$$P(\text{"free"} | \text{non spam}) = \frac{1}{2} = \frac{1}{2}$$

$$P(\text{"offer"} | \text{non spam}) = \frac{2}{2} = 1$$

$$P(\text{"money"} | \text{non spam}) = \frac{0}{2} = 0$$

$$P(\text{nonspam} | \begin{array}{l} \text{free} : 2 \\ \text{offer} : 1 \\ \text{money} : 0 \end{array}) = P(\text{free} | \text{non spam})^2 \times P(\text{offer} | \text{non spam})^1 \times P(\text{money} | \text{non spam})^0 \times P(\begin{array}{l} \text{free} : 2 \\ \text{offer} : 1 \\ \text{money} : 0 \end{array})$$

$$= \left(\frac{1}{2}\right)^2 \times 1 \times P(\begin{array}{l} \text{free} : 2 \\ \text{offer} : 1 \\ \text{money} : 0 \end{array})$$

$$= 0.25 \times P(\begin{array}{l} \text{free} : 2 \\ \text{offer} : 1 \\ \text{money} : 0 \end{array})$$

$$P(\text{spam}) < P(\text{non spam})$$

\therefore Not spam

General formulae for Multinomial

for class (c_k) , prior probability =

$$P(c_k) = \frac{N_k}{N}$$

No of training examples in c_k

total no of training samples

likelihood

$$P(x_i | c_k) = \frac{n_i}{N_k}$$

total (sum) of all feature counts in c_k

No of training examples in c_k

$$P(c_k | x_1^{a_1}, x_2^{a_2}, \dots) \propto P(c_k) \prod P(x_i | c_k)^{a_i}$$

$a_i \rightarrow$ count of feature in test tuple

Instead of raising probability floats to powers, we can simplify by taking log scale

$$\log P(c_k | x) \propto \log \left(P(c_k) \prod_{i=1}^n P(x_i | c_k)^{a_i} \right)$$



$$\propto \log(P(c_k)) + \sum_{i=1}^n a_i \log(P(x_i | c_k))$$

making it linear! No products needed.

We can now calculate the log probabilities instead during training time

And since we are only comparing, no need to take inverse

Note \rightarrow In multinomial we assume only in one direction.
eg too high occurrences = high occurrences

 & not  Possible

0-5 occurrences \rightarrow spam
eg 5-10 occurrences \rightarrow Not spam.
10+ occurrences \rightarrow spam } Not Possible

Note \rightarrow in some references

$$P(X | c_k) = \frac{(\sum n_i)!}{\prod n_i!} \cdot \prod_i P(x_i | c_k)^{n_i}$$

Here the $\frac{(\sum n_i)!}{\prod n_i!}$ represents multinomial coefficient, or no of ways the features in test tuple can be arranged.

Same for all classes (c_k) , hence neglected

matters only when the arrangements matter, not in Naive Bayes.

eg No of possible arrangement of

$\left. \begin{array}{l} \text{"free"} \times 2 \\ \text{"offer"} \times 1 \\ \text{"money"} \times 2 \end{array} \right\} \frac{5!}{2! \times 2!} = \frac{120}{4} = 30$

eg free free offer money
money free offer free
...

Zero Probability Smoothing

What if $(\text{buy} \mid \text{Brand} = z)$ asked and z not in dataset?

If a test data has a value never observed in dataset training phase, then $P = 0$ which may lead to incorrect classification

Entire product $P(x_1, x_2 \dots \mid \text{buy})$ will become 0

$$P(x_i \mid c_k) = \frac{\text{Count}(x_i, c_k)}{\text{Count}(c_k)}$$

$$\text{Count}(x_i, c_k) = 0 \therefore P = 0$$

To overcome this, smoothing is used

① Laplacian smoothing

$$P(x_i \mid c_k) = \frac{\text{Count}(x_i, c_k) + 1}{\text{Count}(c_k) + V}$$

V : total number of unique features
(Vocabulary size)

Probability of unseen feature is small, but still not zero

$$\text{eg } P(\text{"money"} \mid \text{spam}) = \frac{3 + 1}{15 + 3} = \frac{4}{18}$$

\uparrow
 $V = 3$

② Additive smoothing

$$P(x_i \mid c_k) = \frac{\text{Count}(x_i, c_k) + \alpha}{\text{Count}(c_k) + V}$$

α - small constant

Useful when feature set is large

Gaussian Naive Bayes

for continuous data we can use Gaussian Naive Bayes which assumes that the features are in normal distribution

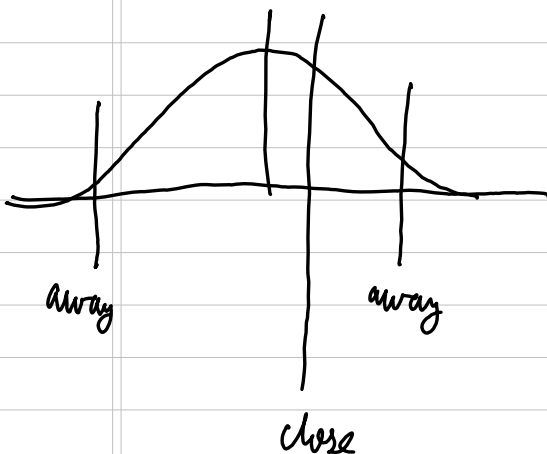
$$P(x_i = x | c_k) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(x_i - \mu)^2}{2\sigma^2}}$$

σ & μ calculated from the data

$$P(c_k | \begin{matrix} x_1 = a_1 \\ x_2 = a_2 \\ \vdots \end{matrix}) = P(c_k) \prod P(x_i = a_i)$$

↖ continuous variable

μ & σ for a class c_k calculated as training.



close to mean \rightarrow high prob
away from mean \rightarrow low prob

Two way can be accommodated

prob never 0, hence no need for regularization

Bernoulli Naive Bayes

Only 2 class data (binary) features are allowed

specialization of general case

eg income : high, low ✓
income : high, medium, low ✗

$$P(x|c_k) = \prod_{i=1}^n P(x_i=1|c_k)^{x_i} (1 - P(x_i=1|c_k))^{1-x_i}$$

$x_i \rightarrow$ 1 : Present
0 : Absent

same as the general case

when present $P(x=1|c_k)$

when absent $1 - P(x=1|c_k) = P(x=0|c_k)$

Generative vs Discriminative models

Generative	Discriminative
Learn joint probability $P(x \wedge y)$	Learn conditional prob $P(x y)$
Model how data is generated to make predictions	Find decision boundary between different classes
Estimate probability of both the input features x and output labels y	Directly map input x and output y without modelling how data is generated
Can generate new data	Can only classify or predict
eg N.B. GMM, HMM BBN	eg D.T. SVM, N.N., logistic regression

N.B. models the distribution for every class, thus making it generative model

BBNs can make new synthetic samples from the conditional probability table

N.B. is generally used for classification only, but here is how it can be used for generation of new samples.

1. sample class label based on $P(y)$
eg $P(y = \text{spam}) = 0.6$, then pick spam 60% time
 2. sample each feature independantly given y
 $P(\text{Free} | \text{spam}) = 0.9$, then pick word free 90% of time
 3. repeat for all features
- Can be used for data augmentation

Advantages → ① Simple

② Fast

③ Low Power

④ Can handle large datasets

Disadvantages → ① Assumption of independent predictors may not hold

② Cannot capture complex relationships in the data. Very simple "naive" model

③ Doesn't work well for small datasets because probability won't be correct.

Applications → ① Spam filtering
② Sentiment analysis
③ Topic classification
④ Recommendation systems
⑤ Anomaly detection
⑥ Credit scoring

Output of training of naive bayes on a dataset is the calculation of prior probabilities & likelihoods

Data Type	Model	Example
Binary	Bernouli N.B.	High/Low
Categorical	N.B.	High/Medium/Low
Discrete	Multinomial N.B.	0, 1, 2, 3
Continuous	Gaussian N.B.	0.1, -10.3, 4.3