

DMM, 4 Feb 2019

Supervised learning

Decision Tree

Logistic Regression

Class Association Rules

Regression

Probabilistic prediction

Bayes rule for conditional probability

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

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

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

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

A_1, A_2, \dots, A_k, C

By simple counting from given training data

$$P(A_1=a_1, A_2=a_2, \dots, A_n=a_n | C=c)$$

When we see $(a_1', a_2', \dots, a_n')$,
what is the category value for C ?

Compare $P(C=c_1 | a_1' \dots, a_n')$
 $P(C=c_2 | a_1' \dots, a_n')$

Use Bayes' rule to calculate

A	B	C
m	b	t
m	s	t
g	q	t
h	s	t
g	q	t
g	q	f
g	s	f
h	b	f
h	q	f
m	b	f

$$P(A=m, B=b | C=t) = \frac{1}{5}$$

$$P(A=m, B=b | C=f) = \frac{1}{5}$$

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

$$P(C=t) = \frac{5}{10}$$

$$P(A=m, B=b) = \frac{2}{10}$$

What is this capturing?

We assume a probabilistic model for
generating the data

Unless you can describe a generative
model, don't claim to do
probabilistic analysis

A	B	C
m	b	t
m	s	t
g	q	t
h	s	t
g	q	t
g	q	f
g	s	f
h	b	f
h	q	f
m	b	f

Here, a plausible model
is:

- ① Randomly chose $C = t/f$
with $P(C)$
- ② Given C , choose
attributes using
 $P(A, B|C)$

Given this generative model, reverse
engineer the probabilities from the data
- Parameter estimation

Digression

Given 7 heads in 10 tosses, we estimate $P(\text{heads})$ as $\frac{7}{10}$ Why?

Suppose $P(\text{heads})$ is some unknown value p

Probability of observed outcome: $\alpha \times p^7 (1-p)^3$

What value of p maximizes this probability
 $= \frac{7}{10}$

Maximum Likelihood Estimate MLE

Observation \rightarrow Parameter θ

$\theta \rightarrow$ likelihood $L(\theta)$ of observation

$$\arg \max_{\theta} L(\theta)$$

Back to Bayesian classification

A	B	C
m	b ✓	t
m	s	t
g ✓	q	t
h	s	t
g ✓	q	t
g ✓	q	f
g ✓	s	f
h	b ✓	f
h	q	f
m	b ✓	f

A=g occurs 4 times

B=b occurs 3 times

$\{A=g, B=b\}$?

$$P(Y) = 0$$

$$P(C=t | A=g, B=b) = \frac{P(g, b | t) \cdot P(t)}{P(g, b)}$$

$$P(C=f | A=g, B=b) = \frac{P(g, b | f) \cdot P(f)}{P(g, b)}$$

Naïve Bayes assumption

Attributes are independent

$$P(A \wedge B) = P(A) \cdot P(B)$$

$$P(A, B|C) = P(A|C) \cdot P(B|C) \quad \text{Conditional independent.}$$

Estimate $P(A|C)$, $P(B|C)$ -- from data

A	B	C
m	b	t
m	s	t
g	q	t
h	s	t
g	q	t
g	q	f
g	s	f
h	b	f
h	q	f
m	b	f

Can we compute a class
for $(A=g, B=b)$

$$P(g|t) = \frac{2}{5} \quad P(b|t) = \frac{1}{5}$$

$$P(g|f) = \frac{2}{5} \quad P(b|f) = \frac{2}{5}$$

$$P(t|g,b) = \frac{P(g|t)P(b|t)P(t)^2}{P(b)P(g)}$$

$$P(g) = \frac{4}{10} \quad P(L) = \frac{3}{10}$$

$$P(t) = P(f) = \frac{5}{10}$$

$$P(f|g,b) = \frac{P(g|f)P(b|f)P(f)^2}{P(b)P(g)}$$

$\Rightarrow 2 \cdot P(t|g,b)$
 \therefore answer is f.

No theoretical justification for naive Bayes
assumption — works well "in practice"

Very easy to implement — build a
model, one pass over data updating
counts

One very popular and successful use
case is spam filters for email

Application Text classification

Spam, or topic classification

Document Model?

Vocabulary V of "known" words

Simplest model - set of words

Boolean model - each $w \in V$ is
either present or absent

Generative model

1. Choose a topic $t_i \in T = \{t_1, \dots, t_m\}$
 - Equally likely? historical ratios?
 - $P(t_i)$
2. Given t_i , for each $w \in V$, include w with probability $P(w|t_i)$
 - Like tossing one coin per word

Given a training set of labelled documents

N documents

N_i documents of topic t_i , for each i

$$P(t_i) = \frac{N_i}{N}$$

For each topic t_i , for each $w \in V$

Among N_i documents, w_i documents contain the word w

$$P(w|t_i) = \frac{w_i}{N_i}$$

Given a document d — some subset of V

$$P(T=t_j | d)$$

$$d = d_1 \dots d_M$$

is a 0-1 vector, $|V|=M$

$$\prod_{i=1}^M P(d_i | t_j) \cdot P(t_j)$$

$$\prod_{i=1}^M P(d_i) \leftarrow \text{same for all } t_j$$

Back to zero counts

Full Bayes \rightarrow Naïve Bayes

May still have rare attribute values.

Suppose 'zebra' is not in training data

$$P(\text{'zebra'} | t_i) = 0 \text{ for all } t_i$$

$$P(t_i | \text{'zebra'}) = \dots P(\text{'zebra'} | t_i) \dots$$

All numerators are 0!

"Smoothing" — Laplace

$n_v = \#$ of times we see v

$n =$ total sample

$$P(v) = \frac{n_v}{n}$$

Suppose values are $\{v_1, \dots, v_m\}$

$$\frac{n_v}{n} \Rightarrow \frac{n_v + 1}{n + m}$$

Variations

$$\frac{n_v + \lambda}{n + \lambda m} \leftarrow \lambda = \frac{1}{n} \text{ is used}$$

Tomorrow - a slightly more sophisticated document model