

Bayes Classifier

- A probabilistic framework for solving classification problems

- Conditional Probability:

$$P(C | A) = \frac{P(A, C)}{P(A)}$$

$$P(A | C) = \frac{P(A, C)}{P(C)}$$

- Bayes theorem:

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

The Good Reverend Bayes



Example of Bayes Theorem

- Given:

- A doctor knows that meningitis causes stiff neck 50% of the time
- Prior probability of any patient having meningitis is 1/50,000
- Prior probability of any patient having stiff neck is 1/20

- If a patient has stiff neck, what's the probability he/she has meningitis?

$$P(M | S) = \frac{P(S | M)P(M)}{P(S)} = \frac{0.5 \times 1/50000}{1/20} = 0.0002$$

Bayesian Classifiers

- Consider each attribute and class label as random variables
- Given a record with attributes (A_1, A_2, \dots, A_n)
 - Goal is to predict class C
 - Specifically, we want to find the value of C that maximizes $P(C | A_1, A_2, \dots, A_n)$
- Can we estimate $P(C | A_1, A_2, \dots, A_n)$ directly from data?

Bayesian Classifiers

- Approach:

- compute the posterior probability $P(C | A_1, A_2, \dots, A_n)$ for all values of C using the Bayes theorem

$$P(C | A_1 A_2 \dots A_n) = \frac{P(A_1 A_2 \dots A_n | C) P(C)}{P(A_1 A_2 \dots A_n)}$$

- Choose value of C that maximizes $P(C | A_1, A_2, \dots, A_n)$
- Equivalent to choosing value of C that maximizes $P(A_1, A_2, \dots, A_n | C) P(C)$

- How to estimate $P(A_1, A_2, \dots, A_n | C)$?

Naïve Bayes Classifier

- Assume independence among attributes A_i when class is given:

- $P(A_1, A_2, \dots, A_n | C) = P(A_1 | C) P(A_2 | C) \dots P(A_n | C)$
- Can estimate $P(A_i | C_j)$ for all A_i and C_j .
- New point is classified to C_j if $P(C_j) \prod P(A_i | C_j)$ is maximal.

How to Estimate Probabilities from Data?

Tid	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

- Class: $P(C) = N_C/N$

– e.g., $P(\text{No}) = 7/10$,
 $P(\text{Yes}) = 3/10$

- For discrete attributes:

$$P(A_i | C_k) = |A_{ik}| / N_{C_k}$$

– where $|A_{ik}|$ is number of instances having attribute A_i and belongs to class C_k

– Examples:

$$P(\text{Status}=\text{Married}|\text{No}) = 4/7$$

$$P(\text{Refund}=\text{Yes}|\text{Yes})=0$$

How to Estimate Probabilities from Data?

- For continuous attributes:

– **Discretize** the range into bins

- ◆ one binary attribute per bin
- ◆ violates independence assumption

– **Two-way split:** $(A < v)$ or $(A > v)$

- ◆ choose only one of the two splits as new attribute

– **Probability density estimation:**

- ◆ Assume attribute follows a normal distribution
- ◆ Use data to estimate parameters of distribution (e.g., mean and standard deviation)
- ◆ Once probability distribution is known, can use it to estimate the conditional probability $P(A_i|c)$

How to Estimate Probabilities from Data?

Tid	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

- Normal distribution:

$$P(A_i | c_j) = \frac{1}{\sqrt{2\pi\sigma_{ij}^2}} e^{-\frac{(A_i - \mu_{ij})^2}{2\sigma_{ij}^2}}$$

- One for each (A_i, c_i) pair

- For (Income, Class=No):

- If Class=No

- sample mean = 110
- sample variance = 2975

$$P(\text{Income} = 120 | \text{No}) = \frac{1}{\sqrt{2\pi}(54.54)} e^{-\frac{(120-110)^2}{2(2975)}} = 0.0072$$

Example of Naïve Bayes Classifier

Given a Test Record:

$X = (\text{Refund} = \text{No}, \text{Married}, \text{Income} = 120\text{K})$

naïve Bayes Classifier:

$P(\text{Refund}=\text{Yes}|\text{No}) = 3/7$
 $P(\text{Refund}=\text{No}|\text{No}) = 4/7$
 $P(\text{Refund}=\text{Yes}|\text{Yes}) = 0$
 $P(\text{Refund}=\text{No}|\text{Yes}) = 1$
 $P(\text{Marital Status}=\text{Single}|\text{No}) = 2/7$
 $P(\text{Marital Status}=\text{Divorced}|\text{No}) = 1/7$
 $P(\text{Marital Status}=\text{Married}|\text{No}) = 4/7$
 $P(\text{Marital Status}=\text{Single}|\text{Yes}) = 2/7$
 $P(\text{Marital Status}=\text{Divorced}|\text{Yes}) = 1/7$
 $P(\text{Marital Status}=\text{Married}|\text{Yes}) = 0$

For taxable income:

If class=No: sample mean=110
 sample variance=2975
 If class=Yes: sample mean=90
 sample variance=25

- $P(X|\text{Class}=\text{No}) = P(\text{Refund}=\text{No}|\text{Class}=\text{No})$
 $\times P(\text{Married}|\text{Class}=\text{No})$
 $\times P(\text{Income}=120\text{K}|\text{Class}=\text{No})$
 $= 4/7 \times 4/7 \times 0.0072 = 0.0024$
- $P(X|\text{Class}=\text{Yes}) = P(\text{Refund}=\text{No}|\text{Class}=\text{Yes})$
 $\times P(\text{Married}|\text{Class}=\text{Yes})$
 $\times P(\text{Income}=120\text{K}|\text{Class}=\text{Yes})$
 $= 1 \times 0 \times 1.2 \times 10^{-9} = 0$

Since $P(X|\text{No})P(\text{No}) > P(X|\text{Yes})P(\text{Yes})$

Therefore $P(\text{No}|X) > P(\text{Yes}|X)$
 $\Rightarrow \text{Class} = \text{No}$

Naïve Bayes Classifier

- If one of the conditional probability is zero, then the entire expression becomes zero
- Probability estimation:

$$\text{Original : } P(A_i | C) = \frac{N_{ic}}{N_c}$$

$$\text{Laplace : } P(A_i | C) = \frac{N_{ic} + 1}{N_c + c}$$

$$\text{m - estimate : } P(A_i | C) = \frac{N_{ic} + mp}{N_c + m}$$

c: number of classes

p: prior probability

m: parameter

Example of Naïve Bayes Classifier

Name	Give Birth	Can Fly	Live in Water	Have Legs	Class
human	yes	no	no	yes	mammals
python	no	no	no	no	non-mammals
salmon	no	no	yes	no	non-mammals
whale	yes	no	yes	no	mammals
frog	no	no	sometimes	yes	non-mammals
komodo	no	no	no	yes	non-mammals
bat	yes	yes	no	yes	mammals
pigeon	no	yes	no	yes	non-mammals
cat	yes	no	no	yes	mammals
leopard shark	yes	no	yes	no	non-mammals
turtle	no	no	sometimes	yes	non-mammals
penguin	no	no	sometimes	yes	non-mammals
porcupine	yes	no	no	yes	mammals
eel	no	no	yes	no	non-mammals
salamander	no	no	sometimes	yes	non-mammals
gila monster	no	no	no	yes	non-mammals
platypus	no	no	no	yes	mammals
owl	no	yes	no	yes	non-mammals
dolphin	yes	no	yes	no	mammals
eagle	no	yes	no	yes	non-mammals

A: attributes

M: mammals

N: non-mammals

$$P(A | M) = \frac{6}{7} \times \frac{6}{7} \times \frac{2}{7} \times \frac{2}{7} = 0.06$$

$$P(A | N) = \frac{1}{13} \times \frac{10}{13} \times \frac{3}{13} \times \frac{4}{13} = 0.0042$$

$$P(A | M)P(M) = 0.06 \times \frac{7}{20} = 0.021$$

$$P(A | N)P(N) = 0.004 \times \frac{13}{20} = 0.0027$$

Give Birth	Can Fly	Live in Water	Have Legs	Class
yes	no	yes	no	?

$$P(A|M)P(M) > P(A|N)P(N)$$

=> Mammals

Naïve Bayes (Summary)

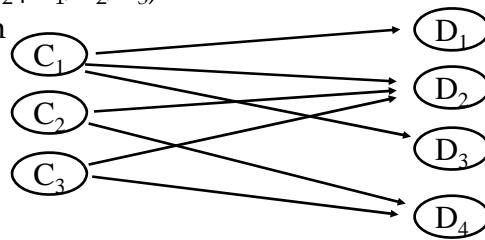
- Robust to isolated noise points
- Handle missing values by ignoring the instance during probability estimate calculations
- Robust to irrelevant attributes
- Independence assumption may not hold for some attributes
 - Use other techniques such as Bayesian Belief Networks (BBN)

Sequential Bayesian Inference

- Consider data attributes one by one
 - Prior probabilities $P(C_i)$
 - Observe data D_j
 - Updates priors using Bayes Rule:
 - Repeat for other attributes using the resulting posterior probability as the new prior
- If attributes are conditionally independent, same as doing it all at once
- Allows choice of what attribute to observe (test to perform) next in terms of cost/benefit.

Bipartite Graphs

- Multiple attributes, multiple classifications
- Classifications are probabilistically independent
- Attributes are conditionally independent
- Attribute probabilities depend only the classes exhibiting them
- Attributes with multiple exhibiting classes require joint probabilities $P(D_2 | C_1, C_2, C_3)$
- Information explosion

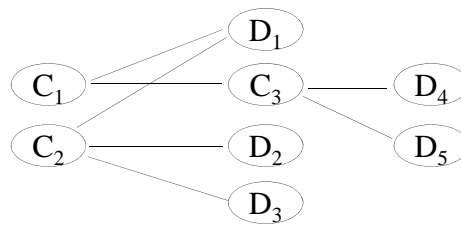


Noisy OR

- Simplify by assuming only one classification holds at a time
- Probability that all classifications exhibit the attribute is just the probability that at least one does
- Thus an attribute is absent only if no class exhibits it
$$1 - P(D_2 | C_1, C_2, C_3) = (1 - P(D_2 | C_1)) (1 - P(D_2 | C_2)) (1 - P(D_2 | C_3))$$
- Use class probabilities for the basic data
 $P_c(D_2 | C_1) = P(D_2 | C_1)$, all other C_i absent
- Reduces probability table size: if N classes and K attributes, from $N2^K$ to NK

Polytrees

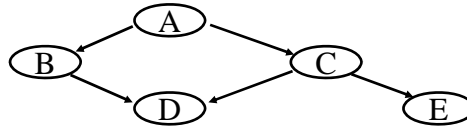
- What if classes are interrelated?
- Polytrees
 - At most one path between any two nodes
- Efficient sequential updating possible



The independence hypothesis...

- ... makes computation possible
- ... yields optimal classifiers when satisfied
- ... but is seldom satisfied in practice, as attributes (variables) are often correlated.
- Attempts to overcome this limitation:
 - Bayesian networks, that combine Bayesian reasoning with causal relationships between attributes
 - Decision trees, that reason on one attribute at the time, considering most important attributes first

Bayesian networks



- Directed acyclic graphs
- Absence of link implies conditional independence

$$P(X_1, \dots, X_n) = \text{Product } P(X_i | \text{parents}(X_i))$$

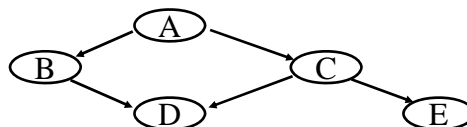
- Specify joint probability tables over parents for each node
- Probability A, B, C, D, E all present:

$$P(A, B, C, D, E) = P(A) * P(B | A) * P(C | A) * P(D | B, C) * P(E | C)$$

- Probability A, C, D present and B, E absent:

$$P(A, \neg B, C, D, \neg E) = P(A) * P(\neg B | A) * P(C | A) * P(D | \neg B, C) * P(\neg E | C)$$

Computing with partial information



- Probability that A present and E absent:

$$\begin{aligned}
 P(A | \bar{E}) &= \sum_{B, C, D} P(A, B, C, D, \bar{E}) \\
 &= \sum_{B, C, D} P(A) P(B | A) P(C | A) P(D | B, C) P(\bar{E} | C) \\
 &= P(A) \sum_C P(C | A) P(\bar{E} | C) \sum_B P(B | A) \sum_D P(D | B, C)
 \end{aligned}$$

- Graph separators (e.g. C) correspond to factorizations
- General problem of finding separators is NP-hard
 - $P(A | \neg E) = P(A, \neg E) / P(\neg E)$

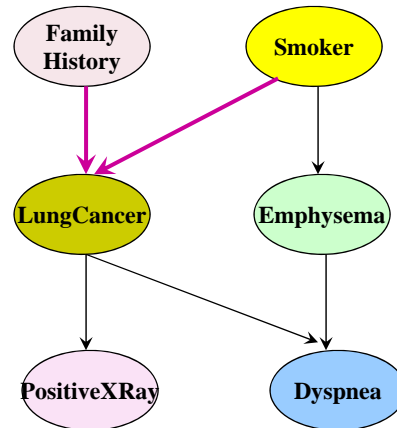
Bayesian networks

- Each node annotated with conditional probability table
 - Probability of node values given values of parent nodes

(FH, S) (FH, ~S) (~FH, S) (~FH, ~S)

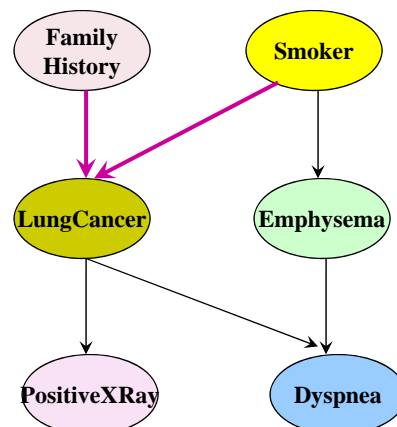
LC	0.8	0.5	0.7	0.1
~LC	0.2	0.5	0.3	0.9

- Conditional probability table for the variable LungCancer



Bayesian networks

- Table of joint probability distribution has $2^6 = 64$ entries
- Bayesian network tables have $8 + 4 + 4 + 8 = 24$ entries



Bayesian Belief Networks

- Bayesian belief networks allow subsets of the variables to be conditionally independent
- A graphical model of causal relationships
- Several methods for learning Bayesian belief networks
 - Given both network structure and all the variables: easy
 - Given network structure but only some variables
 - When the network structure is not known in advance