# Bayes Theorem and Concept Learning (6.3)

- Bayes theorem allows calculating the a posteriori probability of each hypothesis (classifier) given the observation and the training data
- This forms the basis for a straightforward learning algorithm
- Brute force Bayesian concept learning algorithm

# Example: Two categories, one binary-valued attribute

# Data D

Temp	Play Tennis
Hot	Yes
Cold	No

# **Bayes Concept Learning Approach**

Temp	Hypothesis 1	Hypothesis 2	Hypothesis 3	Hypothesis 4
Hot	No	No	Yes	Yes
Cold	No	Yes	No	Yes

# More Interesting Example: Two categories Three Binary Attributes

Task is to learn the output function by observing, D, For n binary inputs there are 2<sup>2<sup>n</sup></sup> possible hypotheses Not all the rows are available!

<b>X</b> 0	<b>x</b> <sub>1</sub>	<b>X</b> <sub>2</sub>	h0	h1	h2		h255
0	0	0	0	1	0		1
0	0	1	0	0	1		1
0	1	0	0	0	0		1
0	1	1	0	0	0		1
1	0	0	0	0	0		1
1	0	1	0	0	0		1
1	1	0	0	0	0		1
1	1	1	0	0	0		1

# **Bayes Concept Learning Approach**

- Best hypothesis:
  - Most probable hypothesis in hypothesis space *H* given training data *D*

- Bayes Theorem:
  - Method to calculate the posterior probability of *h* from the prior probability P(h) together with P(D) and P(D/h)

$$P(h \mid D) = \frac{P(D \mid h)P(h)}{P(D)}$$

# Maximum A Posteriori Probability (MAP) hypothesis

- A maximally probable hypothesis is called a maximum a posteriori (MAP) hypothesis
- Can use Bayes to calculate posterior probability of each candidate hypothesis
- $h_{MAP}$  is a MAP hypothesis provided

h

$$P_{MAP} \equiv \underset{h \in H}{\operatorname{arg\,max}} P(h \mid D)$$
$$= \underset{h \in H}{\operatorname{arg\,max}} \frac{P(D \mid h)P(h)}{P(D)}$$

 $= \arg \max_{h \in H} P(D \mid h) P(h)$ 

# Maximum Likelihood Hypothesis

- P(D|h) is called the likelihood of the data D given h
- If every hypothesis in H is equally probable a priori (P(h<sub>i</sub>) = P(h<sub>i</sub>) for all h<sub>i</sub> and h<sub>i</sub>
- Any hypothesis that maximizes P(D|h) is called a maximum likelihood (ML) hypothesis, h<sub>ML</sub>

$$h_{ML} \equiv \underset{h \in H}{\operatorname{arg\,max}} P(D \mid h)$$

# Brute-Force Bayes Concept Learning (6.3.1)

- Finite hypothesis space H
- To learn a target concept c: X --> {0,1}
- Training examples
  - <<x1, d1>, <x2,d2>,...<xm,dm>>
  - where xi is an instance from X
  - di is the target value of xi, ie, di = c(xi)
  - Simplify notation, D =(d1,.., dm)

# Brute-Force Bayes Concept Learning (6.3.1)

### Brute-Force MAP Learning Algorithm

 For each hypothesis h in H, calculate the posterior probability

 $P(h \mid D) = \frac{P(D \mid h)P(h)}{P(D)}$ 

• Output the hypothesis  $h_{MAP}$  with the highest posterior probability

$$h_{MAP} \equiv \underset{h \in H}{\operatorname{arg\,max}} P(h \mid D)$$

Need to calculate P(h/D) for each hypothesis.
 Impractical for larger hypothesis spaces!

# Choice of P(h) and P(D/h): Assumptions

- The training data *D* is noise-free (*i.e.*,  $d_i = c(x_i)$ ).
- The target concept *c* is contained in the hypothesis space *H*.
- We have no *a priori* reason to believe that any hypothesis is more probable than another.

# Choice of P(h) Given Assumptions

- Given no prior knowledge that one hypothesis (classifier) is more likely than another, same probability is assigned to every hypothesis h in H
- Since target concept is assumed to be contained in H, the prior probabilities should sum to 1
- We should choose,
  - For all h in H

$$P(h) = \frac{1}{\mid H \mid}$$

# Choice of P(D/h) Given Assumptions

- Probability of observing the target values  $D = \langle d1, ... dm \rangle$ for the fixed set of instances <x1,..,xm> given a world in which hypothesis h holds (ie, h is the correct description of the target concept c)
- Assuming noise-free training data

$$P(D \mid h) = \begin{cases} 1 \text{ if } d_i = h(x_i) \text{ for all } d_i \text{ in } D \\ 0 \text{ otherwise} \end{cases}$$

ie, Probability of Data D given hypothesis h is 1 if D is 0 consistent with h and 0 otherwise

# **Bayes Concept Learning Approach**

Temp	Hypothesis 0	Hypothesis 1	Hypothesis 2	Hypothesis 3
Hot	No	No	Yes	Yes
Cold	No	Yes	No	Yes
Prob(D/h)	0	0	1	0

$$P(h_0 | D) = \frac{P(D | h_0)P(h_0)}{P(D)} = \frac{0.\frac{1}{4}}{\frac{1}{4}} = 0$$

Similarly P(h1|D) = P(h3|D) = 0, P(h2|D) = 1

# Brute Force MAP Learning Algorithm

 First step: use Bayes rule to computer posterior probability P(h/D) for each hypothesis h given the training data D

$$P(h \mid D) = \frac{P(D \mid h)P(h)}{P(D)}$$

• If h is inconsistent with the training data D

$$P(h \mid D) = \frac{0.P(h)}{P(D)} = 0$$

# Brute Force MAP Learning Algorithm

• If h is inconsistent with the training data D

 $P(h \mid D) = \frac{1 \cdot \frac{1}{\mid H \mid}}{P(D)}$  $=\frac{1.\frac{1}{|H|}}{|VS_{H,D}|}$ |H| $=\frac{1}{|VS_{H,D}|}$ 

Where  $VS_{H,D}$  is the subset of hypotheses from H that are consistent with D (Version Space of H with respect to D)

# Brute Force MAP Learning Algorithm

## Deriving P(D) from the theorem of total probability

$$P(D) = \sum_{h_i \in H} P(D / h_i) P(h_i)$$
  
=  $\sum_{h_i \in VS_{H,D}} 1 \cdot \frac{1}{|H|} + \sum_{h_i \notin VS_{H,D}} 0 \cdot \frac{1}{|H|}$   
=  $\sum_{h_i \in VS_{H,D}} 1 \cdot \frac{1}{|H|}$   
=  $\frac{|VS_{H,D}|}{|H|}$ 

# Brute-Force MAP Learning Algorithm, continued

• In summary: Bayes theorem implies that the posterior probability P(h|D) under the assumed P(h) and P(D|h)

$$P(h \mid D) = \begin{cases} \frac{1}{\mid VS_{H,D} \mid} & \text{if } h \text{ is consistent with } D\\ 0 & \text{otherwise} \end{cases}$$

• where  $|VS_{H,D}|$  is the number of hypotheses from H consistent with D

# **Brute-Force Bayes Learning**





# MAP Hypotheses and Consistent Learners (6.3.2)

- A learning algorithm is a consistent learner if it outputs a hypothesis that commits zero errors over the training examples.
- Every consistent learner outputs a MAP hypothesis if
  - we assume a uniform prior probability distribution over H (i.e., P(h<sub>i</sub>)=P(h<sub>i</sub>) for all i, j) and
  - we assume deterministic, noise-free training data (*i.e.*, *P*(*D*|*h*)=1 if *D* and *h* are consistent and 0 otherwise).

# **Evolution of Posterior Probabilities**

#### • With increasing training data



- (a) uniform priors to each hypothesis
- (b) As training data increases first to D1
- (c) then to D1^D2 posterior probs for inconsistent hypotheses becomes zero

# Example: Two categories, binary-valued attribute

Temperature	Play Tennis
Hot	Yes
Hot	Yes
Hot	No
Cold	Yes
Hot	Yes
Cold	No
Cold	No
Cold	No
Cold	Yes

# **Bayes Optimal Decision Approach**

Temperature	Play Tennis
Hot	Yes
Hot	Yes
Hot	No
Cold	Yes
Hot	Yes
Cold	No
Cold	No
Cold	No
Cold	Yes

Prob (Hot/Yes) = 0.6Prob (Cold/No) = 0.75Prob (Yes) = 0.56

Prob (Hot) = Prob(Hot/Yes)P(Yes)
+Prob(Hot/No)Prob(No)
=0.6x0.56+0.25x0.44
=0.336+0.11=0.447

Bayes Optimal Decision Prob (Yes/Hot) = Prob (Hot/Yes) P(Yes)/P(Hot) =0.6x0.56/0.447 = 0.75

# Bayes Optimal Rule Example: Medical Diagnosis

- Two alternative hypotheses:
  - Patient has a particular form of cancer
  - Patient does not
- Available Data: a particular laboratory test
  - Lab-test is either +(positive) or negative (-)
- Prior knowledge:
- Over entire population only 0.008 have this disease

# An Example of using Bayes rule

# Known probabilities P(cancer)=.008, P(~cancer)=.992 P(+/cancer)=.98, P(-/cancer)=.02 P(+/~cancer)=.03, P(-/~cancer)=.97

# Statistical Hypothesis Testing Terminology

## Known Probabilities

- P(+/cancer)=.98, P(-/cancer)=.02
- P(+/~cancer)=.03, P(-/~cancer)=.97
- P(cancer)=.008, P(~cancer)=.992



# Bayes rule example (continued)

- Observed data: lab test is positive (+)
- P(+/cancer)P(cancer) = (.98).008 = .0078
- P(+/~cancer)P(~cancer) = (.03).992 = .0298
- Therefore  $h_{MAP} = \sim$  cancer
- Exact a posteriori probabilities
  - P(cancer/+) = .0078/(.0078 + .0298) = .21
  - P(~cancer/+) = .79
- The probability of cancer increased from .008 to .21 after the positive lab test,
  - but still it is still much more likely that it is not cancer

# Continuous valued lab test thresholded to yield positive and negative values



#### Relating Two Types of Error to continuous valued test Lab-test **Cancer-Present Cancer-Absent** Lab-test Positive **4 0.98** Positive 0.03 Negative 0.02 0.97 Negative Cancer Present Cancer Absent False Positive True Positive **Test Value Decision Threshold**

*ErrorRate* = (*FalsePositive* + *FalseNegative*)/2





# Bayes Optimal Classifier (6.7)

# **Bayes Optimal Classification**

$$\underset{v_j \in V}{\operatorname{arg\,max}} \sum_{h_i \in H} P(v_j \mid h_i) P(h_i \mid D)$$

# **Bayes Optimal Classifier**

- Instead of asking "What is the most probable hypothesis given the training data?", ask:
- "What is the most probable classification of the new instance given the training data?"

$$\begin{array}{c} x_{0} \\ x_{1} \\ x_{2} \end{array} \xrightarrow{V_{j}} \\ \end{array} \xrightarrow{V_{j}} \end{array}$$

 Instead of learning the function f<sub>i</sub>, the Bayes optimal classifier assigns any given input to the most likely output v<sub>i</sub>

# **Bayes Optimal Classifier**

 Instead of learning the function, the Bayes optimal classifier assigns any given input to the most likely output

$$x_{1}$$
  
 $x_{2}$   
 $f_{1}$ 

- Calculate a posteriori probabilities  $P(0 \mid x_0, x_1, x_2) = \frac{P(x_0, x_1, x_2 \mid 0)P(0)}{P(x_0, x_1, x_2)}$
- $P(x_0, x_1, x_2|0)$  is the class-conditional probability

## Example of Bayes Optimal Classifier



# **Bayes Optimal Classifier**

- To calculate a posteriori probabilities
  - Need to know Class-conditional probabilities
  - Each is a table of 2<sup>n</sup> different probabilities estimated from many training samples



# **Bayes Optimal Classifier**

Need to know Class-conditional probabilities

 $P(x_0, x_1, x_2 \mid 0)$   $P(x_0, x_1, x_2 \mid 1)$ 

- Tables have 2.2<sup>n</sup> entries in tables
- Will need many training samples:
  - need to see every instance many times in order to obtain reliable estimates
  - When number of attributes is large, impossible to even list all probabilities in a table

# **Bayes Optimal Classifier**

- Target function f(x)
- Takes any value from finite set V, eg 0,1
- Each instance x is composed of attribute values x<sub>1</sub>,x<sub>2</sub>,..,x<sub>n</sub>
- Most possible target value v<sub>map</sub>

$$v_{MAP} = \arg\max_{v_j \in V} P(v_j | x_1, x_2, ..., x_n)$$
  
= 
$$\arg\max\frac{P(x_1, x_2, ..., x_n | v_j) P(v_j)}{P(x_1, x_2, ..., x_n)}$$

0

39

# Most Probable Hypothesis vs Most Probable Classification

- Classification result can be different!
  - Suppose three hypotheses, f<sub>0</sub>,f<sub>1</sub>,f<sub>2</sub> have posterior probabilities given the training data as .3, .4, .3.
  - Therefore MAP hypothesis is f<sub>1</sub>
  - Instance x = <0,0,0> classified as 1 by  $f_1$  but as 0 by  $f_0$  and  $f_2$ 
    - $P(1|x,D)=P(1|f_0,x)P(f_0|D,x)+P(1|f_1,x)P(f_1|D,x)+P(1|f_2,x)P(f_2|D,x)$ 
      - =0..3 + 1..4 + 0..3 = .4
    - Similarly P(0|x,D) = .6
    - Therefore most probable classification of x is 0

