



SNS COLLEGE OF TECHNOLOGY

Coimbatore – 35

An Autonomous Institution

Accredited by NBA – AICTE and Accredited by NAAC – UGC

with 'A++' Grade

Approved by AICTE, New Delhi & Affiliated to Anna

University, Chennai



Introduction to Bayesian Learning



Overview

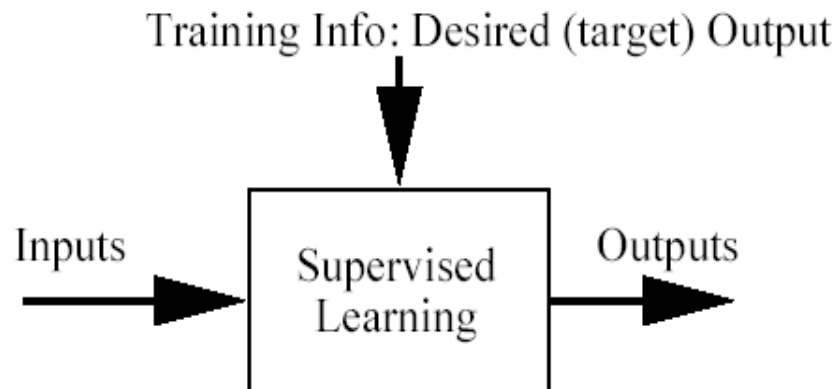
Today we learn about:

- Bayes rule & turn this into a classifier
 - E.g. How to decide if a patient is ill or healthy, based on
 - A probabilistic model of the observed data
 - Prior knowledge



Classification problem

- Training data: examples of the form $(d, h(d))$
 - where d are the data objects to classify (inputs)
 - and $h(d)$ are the correct class info for d , $h(d) \in \{1, \dots, K\}$
- Goal: given d_{new} , provide $h(d_{\text{new}})$



$$\text{Error} = (\text{target output} - \text{actual output})$$



A word about the Bayesian framework

- Allows us to combine observed data and prior knowledge
- Provides practical learning algorithms
- It is a generative (model based) approach, which offers a useful conceptual framework
 - This means that any kind of objects (e.g. time series, trees, etc.) can be classified, based on a probabilistic model specification



Bayes' Rule

$$P(h|d) = \frac{P(d|h)P(h)}{P(d)}$$

Understanding Bayes' rule

d = data

h = hypothesis

Proof. Just rearrange :

$$P(h|d)P(d) = P(d|h)P(h)$$

$$P(d, h) = P(d, h)$$

the same joint probability

on both sides

Who is who in Bayes' rule

$$P(h|d) = \frac{P(d|h)P(h)}{P(d)}$$



Probabilities – auxiliary slide for memory refreshing

- Have two dice h_1 and h_2
- The probability of rolling an i given die h_1 is denoted $P(i|h_1)$. This is a conditional probability
- Pick a die at random with probability $P(h_j)$, $j=1$ or 2 . The probability for picking die h_j and rolling an i with it is called joint probability and is $P(i, h_j)=P(h_j)P(i| h_j)$.
- For any events X and Y , $P(X,Y)=P(X|Y)P(Y)$
- If we know $P(X,Y)$, then the so-called marginal probability $P(X)$ can be computed as $P(X)=\sum_Y P(X,Y)$
- Probabilities sum to 1. Conditional probabilities sum to 1 **provided that their conditions are the same.**



Does patient have cancer or not?

- A patient takes a lab test and the result comes back positive. It is known that the test returns a correct positive result in only 98% of the cases and a correct negative result in only 97% of the cases. Furthermore, only 0.008 of the entire population has this disease.
 1. What is the probability that this patient has cancer?
 2. What is the probability that he does not have cancer?
 3. What is the diagnosis?



$\left. \begin{array}{l} \text{hypothesis } H_1: \text{cancer} \\ \text{hypothesis } H_2: \text{no-cancer} \end{array} \right\} \text{ hypothesis } H$
 -data '+'

$$1. P(\text{cancer} | +) = \frac{P(+ | \text{cancer})P(\text{cancer})}{P(+)} = \frac{\dots\dots\dots}{\dots\dots\dots} = \dots\dots\dots$$

$$P(+ | \text{cancer}) = 0.98$$

$$P(\text{cancer}) = 0.008$$

$$P(+)=P(+ | \text{cancer})P(\text{cancer})+P(+ | \text{no-cancer})P(\text{no-cancer})$$

$$= \dots\dots\dots$$

$$P(+ | \text{no-cancer}) = 0.03$$

$$P(\text{no-cancer}) = \dots\dots\dots$$

$$2. P(\text{no-cancer} | +) = \dots\dots\dots$$

3. Diagnosis?



Choosing Hypotheses

- *Maximum Likelihood* hypothesis:

$$h_{MIL} = \underset{h \in H}{\operatorname{argmax}}$$

- Generally we want the most probable hypothesis given training data. This is the *maximum a posteriori* hypothesis:

$$h_{MAP} = \underset{h \in H}{\operatorname{argmax}}$$

- Useful observation: it does not depend on the denominator $P(d)$



The Naïve Bayes Classifier

- What can we do if our data d has several attributes?
- Naïve Bayes assumption: Attributes that describe data instances are conditionally independent given the classification hypothesis



- it is a simplifying assumption, obviously it may be violated in reality
- in spite of that, it works well in practice
- The Bayesian classifier that uses the Naïve Bayes assumption and computes the MAP hypothesis is called Naïve Bayes classifier
- One of the most practical learning methods
- Successful applications:
 - Medical Diagnosis
 - Text classification



Example. 'Play Tennis' data

Day	Outlook	Temperature	Humidity	Wind	Play Tennis
<i>Day1</i>	Sunny	Hot	High	Weak	<i>Nb</i>
<i>Day2</i>	Sunny	Hot	High	Strong	<i>Nb</i>
<i>Day3</i>	Overcast	Hot	High	Weak	<i>Yes</i>
<i>Day4</i>	Rain	Mild	High	Weak	<i>Yes</i>
<i>Day5</i>	Rain	Cool	Normal	Weak	<i>Yes</i>
<i>Day6</i>	Rain	Cool	Normal	Strong	<i>Nb</i>
<i>Day7</i>	Overcast	Cool	Normal	Strong	<i>Yes</i>
<i>Day8</i>	Sunny	Mild	High	Weak	<i>Nb</i>
<i>Day9</i>	Sunny	Cool	Normal	Weak	<i>Yes</i>
<i>Day10</i>	Rain	Mild	Normal	Weak	<i>Yes</i>
<i>Day11</i>	Sunny	Mild	Normal	Strong	<i>Yes</i>
<i>Day12</i>	Overcast	Mild	High	Strong	<i>Yes</i>
<i>Day13</i>	Overcast	Hot	Normal	Weak	<i>Yes</i>
<i>Day14</i>	Rain	Mild	High	Strong	<i>Nb</i>



Naïve Bayes solution

Classify any new datum instance $\mathbf{x}=(a_1, \dots, a_T)$ as:

$$\underset{h}{\text{argmax}} \prod_{t=1}^T P(a_t | h)$$

- To do this based on training examples, we need to estimate the parameters from the training examples:

- For each target value (hypothesis) h

$$\hat{P}(h)$$

$$\hat{P}(a_t | h)$$

- For each attribute value a_t of each datum instance



Based on the examples in the table, classify the following datum x :

$x=(Outl=Sunny, Temp=Cool, Hum=High, Wind=strong)$

- That means: Play tennis or not?



- Working:

RP Play Tennis 0.154064

RP Play Tennis 0.15406

RW Windr | Prlguy Tennis 0.3903

RW Windr | Prlguy Tennis 0.35060

etc

Ry | Sunny | Cool | High | Strong 0.00

Rn | Sunny | Cool | High | Strong 0.020

=ans: Play Tennis

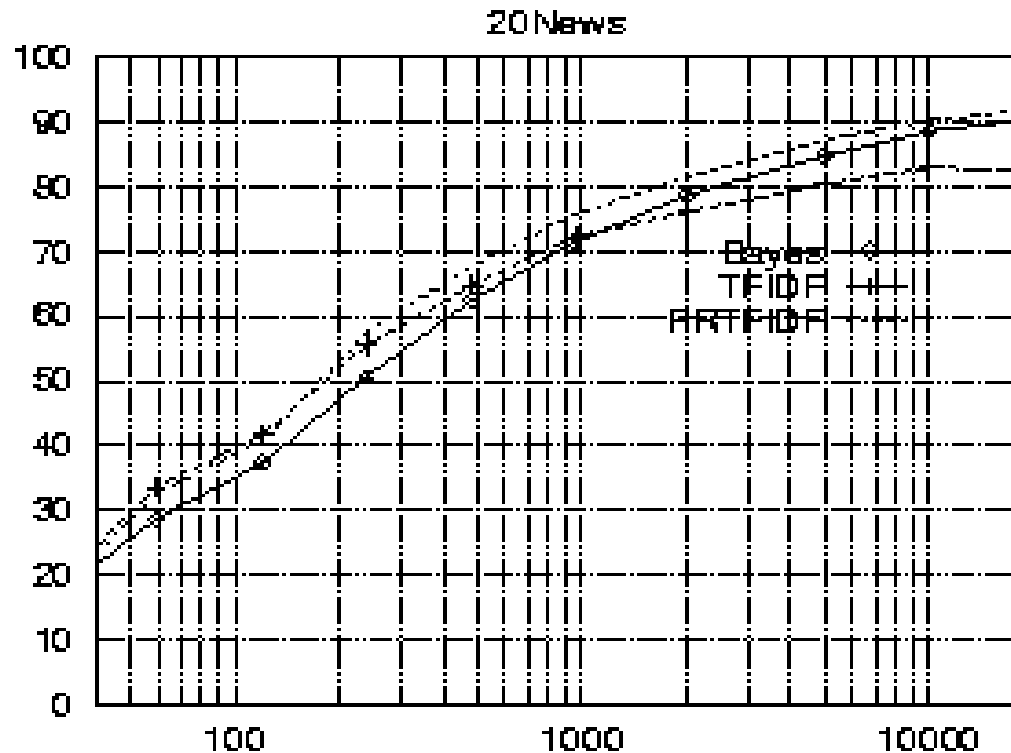


Learning to classify text

- Learn from examples which articles are of interest
- The attributes are the words
- Observe the Naïve Bayes assumption just means that we have a random sequence model within each class!
- NB classifiers are one of the most effective for this task
- Resources for those interested:
 - Tom Mitchell: Machine Learning (book) Chapter 6.



Results on a benchmark text corpus



Accuracy vs. Training set size (1/3 withheld for test.)



Remember

- Bayes' rule can be turned into a classifier
- Maximum A Posteriori (MAP) hypothesis estimation incorporates prior knowledge; Max Likelihood doesn't
- Naive Bayes Classifier is a simple but effective Bayesian classifier for vector data (i.e. data with several attributes) that assumes that attributes are independent given the class.
- Bayesian classification is a generative approach to classification



Resources

- Textbook reading (contains details about using Naïve Bayes for text classification):
Tom Mitchell, Machine Learning (book), Chapter 6.
- Further reading for those interested to learn more:
<http://www-2.cs.cmu.edu/~tom/NewChapters.html>