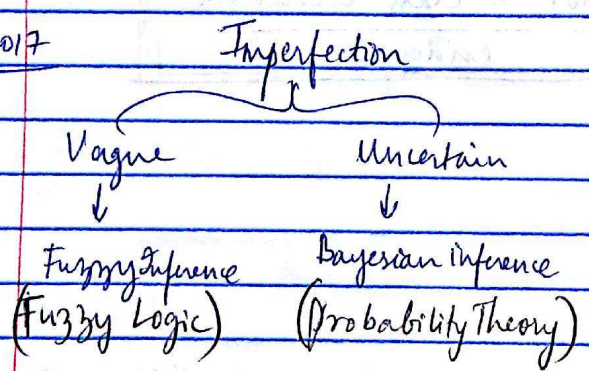


March 9<sup>th</sup>, 2017



$U$  = Universe of discourse  
 $A$  = set (event)  $\subset U$   
 $x$  = variable  $\in A$

$P(A)$	$\mu_A(x)$
probability that an <u>ill-known</u> variable $x$ ranging on $U$ hits the <u>well-known</u> set $A$ .	Membership of <u>well-known</u> variable $x$ ranging on $U$ hits the <u>ill-known</u> set $A$ .

before it happens → event → after it happened

Eg. prob that next person entering room is neary red.      Eg. intensity of red?

Probability based on <u>Measure Theory</u> * Domain = $2^{101}$ * Based on Boolean Algebra	Fuzzy comes from <u>Set theory</u> * <u>logic</u> Domain = $[0, 1]$ * Cannot be Boolean Algebra.
--	--



Two major probabilistic methods

- ① Naive Bayes Learning
- ② Bayesian Belief Method.

Bayes Theorem

(What is the probability that my hypothesis (solution) is the right solution given the training data (D))

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

$$\text{posterior} = \frac{\text{prior} \times \text{class likelihood}}{\text{evidence}}$$

$h \in H$

all possible hypotheses/solutions

Bayes Rule Example

Meningitis can cause the patient to have a stiff neck. 50% of the time. Also,  $P(\text{meningitis}) = \frac{1}{50,000}$

$$P(\text{Meningitis}/s) = \frac{P(m/s)}{P(s)} = \frac{P(s/m) \times P(m)}{P(s)}$$

$$P(s/m) = 0.5$$

$$P(m) = 1/50,000$$

$$P(s) = 1/20$$

$$P(m/s) = 0.5 \times \frac{1}{50,000} = \frac{0.00001}{20}$$

$$0.00002$$

Choosing Hypotheses

Maximum a posteriori hypothesis

$$h_{\text{MAP}} = \arg\max_{h \in H} P(h/D) = \arg\max_{h \in H} \frac{P(D/h) P(h)}{P(D)} = \arg\max_{h \in H} P(D/h) P(h)$$



## Basic Rules for Probabilities.

- Product rule  $P(A \cap B) = P(A/B) \cdot P(B)$   
 $= P(B/A) \cdot P(A)$

- Sum rule  $P(A \cup B) = P(A) + P(B) - P(A \cap B)$

- Theorem of Total Probabilities  
If events  $A_1, A_2, \dots, A_n$  are mutually exclusive with  $\sum_{i=1}^n P(A_i) = 1$ , then  $P(B) = \sum_{i=1}^n P(B/A_i) P(A_i)$

### Brute Force MAP

is very slow. (intractable)  
 $h_{MAP} = \underset{h \in H}{\operatorname{argmax}} P(h/D)$

If the size of the problem  $< 32$  bits, then brute force works.

### Minimum Description Length (MDL)

Use the length ( $h$ ) + length (misclassifications).

$v_j \in V$  (set of classification cases) e.g.  $\{+, -\}$

$$\underset{v_j \in V, h_i \in H}{\operatorname{argmax}} \sum P(v_j / h_i) P(h_i / D)$$

$$P(h_1 / D) = 0.4$$

$$P(h_2 / D) = 0.3$$

$$P(h_3 / D) = 0.3$$

$$P(+ / h_1) = 0$$

$$P(+ / h_2) = 1$$

$$P(+ / h_3) = 1$$

(Prob of "no" classification given  $h_i$  solution)

$$P(- / h_1) = 1$$

$$P(- / h_2) = 0$$

$$P(- / h_3) = 0$$



(Optimal Bayes)

$$\sum_{h_i \in H} P(\oplus | h_i) P(h_i | D) = 0.4$$

$$\sum_{h_i \in H} P(\ominus | h_i) P(h_i | D) = 0.6$$

The decision is  $\ominus$

(Practical)

### Practical Naive Bayes Classifier (NBC)

NBC is, along with Decision Trees and ANNs, among the most commonly used ML methods.

when to use? ① Moderate  $\rightarrow$  large training data.

② Attributes are conditionally independent given classification.  
("naively" assumes that conditions are independent while they aren't.)

Successful applications of NBC

① Diagnosis.

② Classify text document.

Assume the target function  $f: X \rightarrow Y$ , where each instance  $x$  is described by attributes  $\langle a_1, a_2, \dots, a_n \rangle$ . The most probable value of  $f(x)$  is  $V_{NB} = \arg \max_{V_j \in V} P(V_j | a_1, a_2, \dots, a_n)$ .

$$= \arg \max_{V_j \in V} \frac{P(a_1, a_2, \dots, a_n | V_j) P(V_j)}{P(a_1, a_2, \dots, a_n)} = \arg \max_{V_j \in V} P(a_1, a_2, \dots, a_n | V_j) * P(V_j)$$

Naive Bayes Assumption:

attributes are conditionally independent. This means:

$$P(a_1, a_2, \dots, a_n | V_j) = \prod P(a_i | V_j)$$

This is a naive assumption

★ NBC  $V_{NB} = \arg \max_{V_j \in V} P(V_j) \prod P(a_i | V_j)$

Naive Bayes Learn (examples)

For each target value  $v_j$

$\hat{P}(v_j) \leftarrow \text{estimate } P(v_j)$

For each attribute value  $a_i$  of  $a$

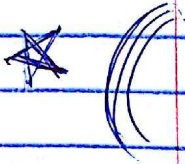
$\hat{P}(a_i | v_j) \leftarrow \text{estimate } (P(a_i | v_j))$

Save this



Classify New Instance (x)

$$V_{NB} = \underset{V_j \in V}{\operatorname{argmax}} \hat{P}(V_j) \prod \hat{P}(a_i | V_j)$$



|V| should be small.  
(# of classification states)

Play Tennis

$$\begin{aligned} P(y) P(\text{sun}/y) P(\text{cool}/y) P(\text{high}/y) P(\text{strong}/y) &= 0.005 \\ P(n) P(\text{sun}/n) P(\text{cool}/n) P(\text{high}/n) P(\text{strong}/n) &= 0.021 \end{aligned}$$

$\hookrightarrow V_{AB} = n.$

even if one probability is 0, then the whole probability = 0.

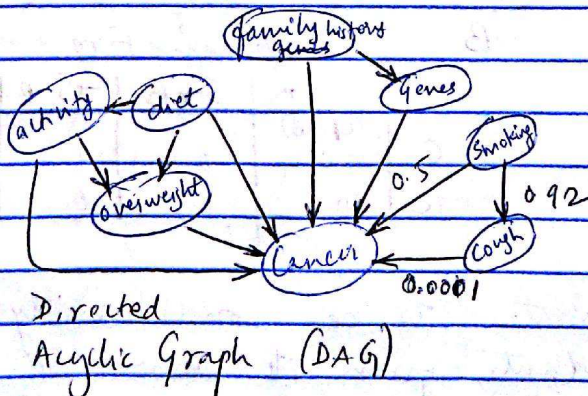
So use lots of data so that no probability is 0.

If it is 0, then we use an  $m$ -estimator to find some value to estimate a non zero value when you encounter it in N.B.

$$m \rightarrow \text{estimator } P(a_i | V_j) = \frac{n_c + m}{n + m}$$

= 0 ?

March 14<sup>th</sup>, 2017



Belief network  
how to infer?

Bayesian Belief Network

- Naive Bayes is too restrictive
- The problem may prove to be intractable without such assumptions.
- Bayesian Belief Networks describe conditional independencies among subsets of variables.
- allows combining the prior knowledge about (in)dependencies