

[BAYES CLASSIFICATION]

—(1)

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

maximum likelihood (ML) hypothesis

$$h_{ML} = \operatorname{argmax}_{h_i \in H} P(D/h_i)$$

$$h_{MAP} = \operatorname{argmax}_{h \in H} P(h/D)$$

$$= \operatorname{argmax}_{h \in H} \frac{P(D/h) P(h)}{P(D)}$$

$$= \operatorname{argmax}_{h \in H} P(D/h) \cdot P(h)$$

Does patient have cancer or not?

A patient takes a lab test and the result comes back positive. The test returns a correct positive result (+) in only 98% of the cases in which the disease is actually present, and a correct negative result (-) in only 97% of the cases in which the disease is not present.

Furthermore 0.008 of the entire population have this cancer.

②

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

$$P(\neg \text{cancer}) = 0.992$$

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

$$P(+|\text{cancer}) = 0.02$$

$$P(+|\neg \text{cancer}) = 0.03$$

$$P(-|\neg \text{cancer}) = 0.97$$

$$P(+|\text{cancer}) \cdot P(\text{cancer}) = 0.98 \times 0.008 = \underline{0.0078}$$

$$P(+|\neg \text{cancer}) \cdot P(\neg \text{cancer}) = 0.03 \cdot 0.992 = \underline{0.0298}$$

$$\underline{h_{\text{MAP}} = \neg \text{cancer}}$$

• Naive assumption: Attribute independence

$$P(x_1, \dots, x_k | c) = P(x_1 | c) \cdot \dots \cdot P(x_k | c)$$

• If i -th attribute is categorical:

$P(x_i | c)$ is estimated as the relative freq. of samples having value x_i as i -th attribute in class c .

• If i -th attribute is continuous

$P(x_i | c)$ is estimated thru a Gaussian density function

Play - Tennis example; estimating $P(x_i|c)$ ③

$P(\text{play})$, $n(\text{don't play})$

$P(p) = 9/14$	$P(n) = 5/14$
---------------	---------------

Outlook

$$P(\text{Sunny}|p) = 2/9$$

$$P(\text{Sunny}|n) = 3/5$$

$$P(\text{Overcast}|p) = 4/9$$

$$P(\text{Overcast}|n) = 0$$

$$P(\text{rain}|p) = 3/9$$

$$P(\text{rain}|n) = 2/5$$

Temperature

$$P(\text{hot}|p) = 2/9$$

$$P(\text{hot}|n) = 2/5$$

$$P(\text{mild}|p) = 4/9$$

$$P(\text{mild}|n) = 2/5$$

$$P(\text{cool}|p) = 3/9$$

$$P(\text{cool}|n) = 1/5$$

Humidity

$$P(\text{high}|p) = 3/9$$

$$P(\text{high}|n) = 4/5$$

$$P(\text{normal}|p) = 6/9$$

$$P(\text{normal}|n) = 2/5$$

Windy

$$P(\text{true}|p) = 3/9$$

$$P(\text{true}|n) = 3/5$$

$$P(\text{false}|p) = 6/9$$

$$P(\text{false}|n) = 2/5$$

Play-tennis example: classify X

④

$X = \langle \text{rain}, \text{hot}, \text{high}, \text{false} \rangle$

$$\underline{P(X|p) \cdot P(p)} = \underline{P(\text{rain}|p) \cdot P(\text{hot}|p) \cdot P(\text{high}|p)} \\ \cdot \underline{P(\text{false}|p) \cdot P(p)}$$

$$= \frac{3}{9} \cdot \frac{2}{9} \cdot \frac{3}{9} \cdot \frac{6}{9} \cdot \frac{9}{14} = \underline{\underline{0.010582}}$$

$$P(X|n) \cdot P(n) = P(\text{rain}|n) \cdot P(\text{hot}|n) \cdot P(\text{high}|n) \\ \cdot P(\text{false}|n) \cdot P(n)$$

$$= \frac{2}{5} \cdot \frac{2}{5} \cdot \frac{4}{5} \cdot \frac{2}{5} \cdot \frac{5}{14} = \underline{\underline{0.018286}}$$

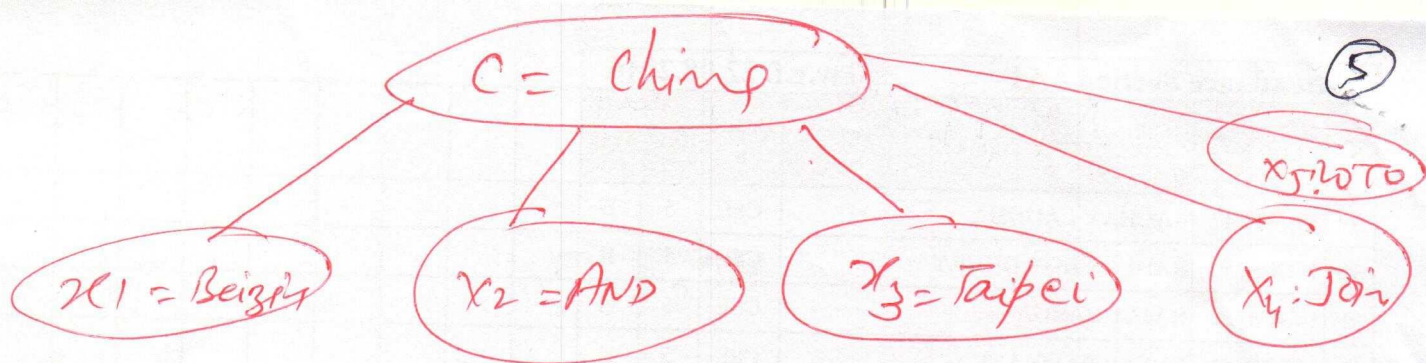
• Sample X is classified in class n
(don't play)

Training set
doc ID

inc = China

1	Chinese	Beijing	Chinese	Yes
2	Chinese	Chinese	Shanghai	Yes
3	Chinese	Macao		Yes
4	Tokyo	Japan	Chinese	no

Tests: Chinese Chinese Chinese Tokyo Japan = ?



Beijing joins the WTO; chinese

$$P(\text{Yes}) = \frac{3}{4}$$

$$P(\text{No}) = \frac{1}{4}$$

$$P(\text{chinese} | \text{yes}) = \frac{(5+1)}{8+6} = \frac{6}{14} = \frac{3}{7}$$

$$P(\text{tokyo} | \text{yes}) = \frac{(0+1)}{14} =$$

$$P(\text{Japan} | \text{yes}) = \frac{(0+1)}{14}$$

$$P(\text{chinese} | \text{No}) = \frac{(1+1)}{(3+6)} = \frac{2}{9}$$

$$P(\text{Japan} | \text{No}) = \frac{(1+1)}{(3+6)} = \frac{2}{9}$$

$$P(\text{tokyo} | \text{No}) = \frac{(1+1)}{(3+6)} = \frac{2}{9}$$

Using Laplace smoothing

$$\hat{P}(t|c) = \frac{T_{ct} + 1}{\left(\sum_{t' \in V} (T_{ct'} + 1) \right)}$$

$$= \frac{T_{ct} + 1}{\left(\sum_{t' \in V} T_{ct'} \right) + B'}$$

$B' = |V|$ is the no of terms in the vocabulary

$T_{ct} \Rightarrow$ no of occurrences of t in training doc from class c .

$$P(\text{yes} | ds) = \frac{3}{4} \times \frac{3}{7} \times \frac{3}{7} \times \frac{3}{7} \cdot \frac{1}{14} \times \frac{1}{14} \quad (6)$$

$$= \underline{\underline{0.0003}}$$

$$P(\text{no} | ds) = \frac{1}{4} \cdot \frac{2}{9} \times \frac{2}{9} \times \frac{2}{9} \cdot \frac{2}{9} \cdot \frac{2}{9}$$

$$= \underline{\underline{0.0001}}$$

So the DS belongs to class chins.