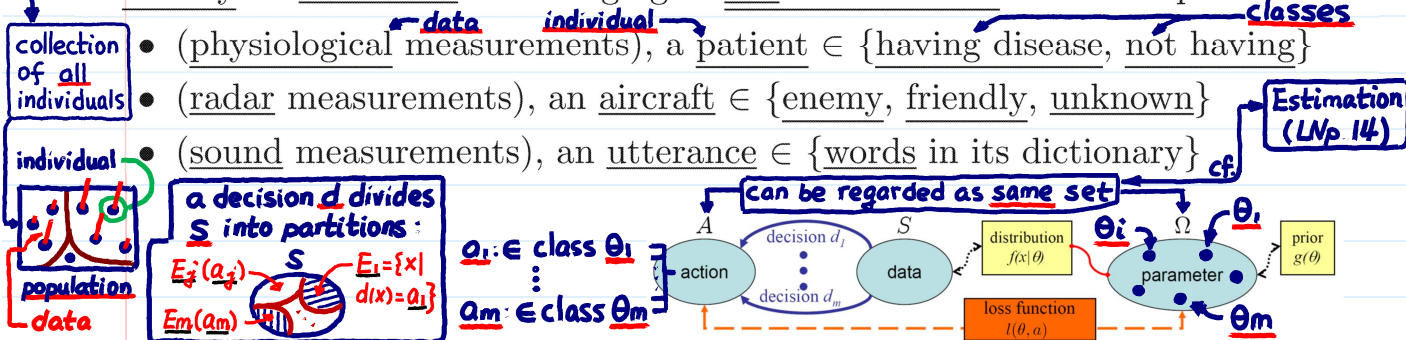


Application of Decision Theory: Classification

cf. Estimation Testing

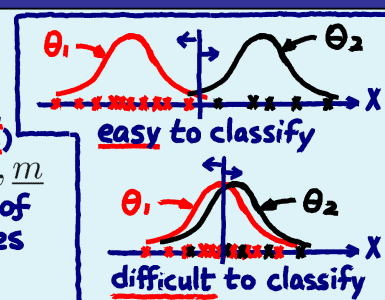
Data

A population consists of m classes, on the basis of some measurements, we wish to classify an individual as belonging to one of the classes. For example:



Formulation of Classification Problem (2nd Ed., TBp.581)

- Data: $X (\in S)$ = the measurements of one particular individual
- $\Omega = \{\theta_1, \theta_2, \dots, \theta_m\}$: class membership. usually, finite classes
- $\pi_1, \pi_2, \dots, \pi_m$: prior probabilities, e.g. may use proportions of individuals in each classes
- $\pi_i = P(\Theta = \theta_i), i = 1, 2, \dots, m$, with $\sum_{i=1}^m \pi_i = 1$.
- $f(x|\theta_i)$: assumed known (in practice, the key to the problem)
- l_{ij} : loss in classifying a member of class i to class j



Why?

$i=j \rightarrow$ correct classification
 $i \neq j \rightarrow$ misclassification

$l_{ij} = l(\theta_i, a_j): m \times m \text{ matrix}$

Theorem 10.8 (Bayes rule for Classification, 2nd Ed., TBp.581)

- Observe $X = x$, posterior distribution of Θ is

$$h(\theta_i | x) = P(\Theta = \theta_i | X = x) = \frac{\pi_i f(x|\theta_i)}{\sum_k \pi_k f(x|\theta_k)}$$

joint, e.g., $P(X=x, \Theta=\theta_i)$

update

marginal, e.g., $P(X=x)$

Q: How should we assign prior in classification?

Ans. If available, the proportion of individuals of each classes in the population is a good choice (check LNp.29)

- posterior risk of assigning an observation x to the j th class is

$$PR(j|x) = \sum_{i=1}^m l_{ij} h(\theta_i | x) = \frac{\sum_{i=1}^m l_{ij} \pi_i f(x|\theta_i)}{\sum_k \pi_k f(x|\theta_k)}$$

average loss of a_j over θ_i 's with posteriors as weights

- Bayes rule $d(x)$: choose the value of j for which $PR(j|x)$ is a minimum.
- map to a_j

Example 10.10 (0-1 loss, 2nd Ed., TBp. 581-582)

Q: Why not use L^2 -norm?

- 0-1 loss function:

average loss of θ_i over X with $X|\theta_i$ as weights

$$l_{ij} = \begin{cases} 0, & i = j \text{ } \leftarrow \text{correct classification} \\ 1, & i \neq j \text{ } \leftarrow \text{misclassification} \end{cases}$$

quantitative θ

qualitative θ

$i \backslash j$	θ_1	\dots	θ_m
a_1	0		1
\vdots			
a_m	1		0

- For any classification rule $d(X)$, the risk function is

$$R(i, d) = E_X[l(\theta_i, d(X))] = \sum_{j=1}^m l_{ij} P_{\theta_i}(d(X) = j) = \sum_{j \neq i} P_{\theta_i}(d(X) = j)$$

E_j (LNp.29)

θ_i

a_j

a_i

θ_i

$\sum_{j=1}^m P_{\theta_i}(d(X) = j) = 1$

which is the probability of misclassification in class i .

- Posterior risk of assigning an observation \underline{x} to the j th class:

Bayesian approach \rightarrow Q_j \rightarrow **fixed**

$$PR(j|\underline{x}) = \sum_{i=1}^m l_{ij} h(\theta_i|\underline{x}) = \sum_{i \neq j} h(\theta_i|\underline{x}) = \frac{\sum_{i \neq j} \pi_i f(\underline{x}|\theta_i)}{\sum_{k=1}^m \pi_k f(\underline{x}|\theta_k)}$$

posterior $\propto f(\underline{x}|\theta) g(\theta)$ = **joint**

same \rightarrow **same derivation as for (*) in LNp.30**

$$= P_{\Theta|\underline{X}}(\Theta \neq \theta_j|\underline{x}) = 1 - P_{\Theta|\underline{X}}(\Theta = \theta_j|\underline{x}) = 1 - h(\theta_j|\underline{x})$$

- $PR(j|\underline{x})$ is minimized for that value of j such that $P_{\underline{X}}(\underline{x})$: **marginal**

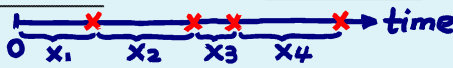
Bayes rule $\leftarrow P_{\Theta|\underline{X}}(\Theta = \theta_j|\underline{x})$ is maximized, i.e., **maximize** **joint pmf** \rightarrow **Bayes rule**

(reasonable?) **mode of $\Theta|\underline{X}$**

for the class that has maximum posterior probability.

Sum=1(?) **Ans. No** **LN, CH1~6** **p.21**

Example 10.11 (waiting times between emissions, 2nd Ed., TBp. 582)

- Data:** $\underline{X} = (X_1, X_2, \dots, X_n) = n$ waiting times between emissions of alpha particles from a radioactive substance. 
- On the basis of \underline{X} , a decision is to be made as to whether to classify the particles as coming from substance I or substance II.
- Statistical Modeling:** $\Omega = \{\theta_1, \theta_2\}$: **2 classes**

risk = misclassification probability

- X_1, X_2, \dots, X_n are i.i.d. $\sim E(\theta_1)$ if particles came from substance I
- X_1, X_2, \dots, X_n are i.i.d. $\sim E(\theta_2)$ if particles came from substance II
- prior:** $\pi_1 = \pi_2 = 1/2$
- loss function:** 0-1 loss

exponential \rightarrow **assumed known**

Q: What if we treat it as a testing problem?
e.g. $H_0: E(\theta_1)$ vs. $H_A: E(\theta_2)$
cf. more protection