SOME INTERESTING DECISION FUNCTIONS Michael F. Capobianco Polytechnic Institute of Brooklyn

Basic Considerations of Statistical Decision Theory

The problem is to decide which of q possible states of nature $\theta_1, \theta_2, \ldots, \theta_q$ is the true one by observing the outcome of some experiment which has n possible outcomes x_1, x_2, \ldots, x_n For each θ_i there is a probability distribution vector

$$\mathbf{P}_{j} = \begin{bmatrix} \mathbf{P}_{1j} \\ \mathbf{P}_{2j} \\ \vdots \\ \mathbf{P}_{nj} \end{bmatrix}$$

where $p_{ij} = P(x_i | \theta_j)$. We also form a loss vector for each θ_i

$$\mathbf{W}_{j} = \begin{bmatrix} \mathbf{w}_{1j} \\ \mathbf{w}_{2j} \\ \vdots \\ \mathbf{w}_{qj} \end{bmatrix}$$

where w_{ii} = the loss incurred in making decision d when $\hat{\theta}_{j}$ is the true state of nature; $w_{ij} \ge 0$ and = 0 if and only if i = j. In order to make decisions we need a mechanism for choosing a θ_{i} upon observing an x. Such a mechanism is called a decision function and can be represented by a matrix

$$A = \begin{bmatrix} a_{11} & a_{12} \cdots & a_{1n} \\ \vdots \\ a_{q1} & \cdots & a_{qn} \end{bmatrix}$$

where $a_{ij} = P(d_i | x_j)$.

The idea is to find a matrix A that in some sense minimizes the loss. We can compute the expected loss, called the risk, for any A under a given state of nature. This is denoted by $R(A, \theta_i)$ and

$$R(A, \theta_j) = W'_j A P_j$$

where W; is the transpose of W. In the absence of any further information one way of choosing a single decision function is by using the minimax criterion, i.e., choose the function with risk equal to min max $R(A, \theta)$. θ

Α

Another possibility is to use the maximum likelihood decision function, i.e., choose the function A such that

$$a_{ij} = \begin{cases} 1 & \text{if } p_{ji} \ge p_{j\ell} \text{ for all } \ell \\ 0 & \text{otherwise} \end{cases}$$

This is an example of an non-randomized decision function; each column of A has a single entry equal to 1 and all other entries equal to 0. In such a case we say that A is non-randomized. Clearly each column of any A must add up to 1.

If an a priori distribution is available, i.e., if one has a vector

$$\mathbf{P} = \begin{bmatrix} \mathbf{P}(\boldsymbol{\theta}_{1}) \\ \mathbf{P}(\boldsymbol{\theta}_{2}) \\ \vdots \\ \mathbf{P}(\boldsymbol{\theta}_{q}) \end{bmatrix}$$

where $P(\theta_i)$ = the probability that θ_i is the true state of nature, one can then find the expected risk (or Bayes risk)

$$\sum w'_j \text{ AP}_j P$$

where $\sum = |1 | \dots |$, a row vector of q l's. One now chooses the decision function with the minimum expected risk. This is called the Bayes decision function.

One property that a decision function should have is that of admissibility. To explain this term we introduce first the notion of dominance. If $R(A, \theta_j) \leq R(B, \theta_j)$ for all θ_j and strict inequality holds for at least one θ_i , then A is said to dominate B. A decision function is admissible if it is not dominated by any other one.

Proportional Likelihood Decision Function

We argue as follows: There seems to be a weakness in the maximum likelihood criteria in that it chooses that state of nature θ , which yields the observed x, with the highest probability, even though some other state of nature may yield x. with a probability almost as high. It seems reasonable that it would be better to give all states of nature a chance of being chosen which is proportional to their respective probabilities of yielding x. We, therefore, propose to form the matrix Å with

$$a_{ij} = \frac{P(x_j | \theta_i)}{\sum_{i=1}^{q} P(x_j | \theta_i)} = \frac{P_{ji}}{\sum_{i=1}^{q} P_{ji}}$$

The following example shows that such a decision function may be admissible.

$$\mathbf{P}_{1} = \begin{bmatrix} \frac{1}{3} \\ \frac{1}{3} \\ \frac{1}{3} \\ \frac{1}{3} \end{bmatrix} \qquad \mathbf{P}_{2} = \begin{bmatrix} \frac{1}{3} \\ \frac{1}{2} \\ \frac{1}{6} \end{bmatrix}$$

$$W_{1} = \begin{bmatrix} 0 \\ 1 \end{bmatrix} \qquad W_{2} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$
$$A_{1} = \begin{bmatrix} 1 & 1 & 1 \\ 0 & 0 & 0 \end{bmatrix} \qquad A_{2} = \begin{bmatrix} 0 & 1 & 1 \\ 1 & 0 & 0 \end{bmatrix}$$
$$A_{2} = \begin{bmatrix} 1 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \qquad A_{6} = \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 1 \end{bmatrix}$$
$$A_{3} = \begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix} \qquad A_{7} = \begin{bmatrix} 0 & 0 & 1 \\ 1 & 1 & 0 \end{bmatrix}$$
$$A_{4} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 1 \end{bmatrix} \qquad A_{8} = \begin{bmatrix} 0 & 0 & 0 \\ 1 & 1 & 1 \end{bmatrix}$$

Sample Calculation:

$$\mathbf{R}(\mathbf{A}_{4}, \boldsymbol{\theta}_{1}) = \mathbf{W}_{1}' \mathbf{A}_{4} \mathbf{P}_{1} = \begin{bmatrix} 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 1 \end{bmatrix} \begin{bmatrix} \frac{1}{3} \\ \frac{1}{3} \\ \frac{1}{3} \\ \frac{1}{3} \end{bmatrix} = \frac{2}{3}$$

Tabulation

	R(Α, θ ₁)	R(A, θ ₂)
A	0	1
A ₂	$\frac{1}{3}$	5.
A ₃	$\frac{1}{3}$	$\frac{1}{2}$
A ₄	$\frac{2}{3}$	$\frac{1}{3}$
A ₅	$\frac{1}{3}$	$\frac{2}{3}$
A.6	$\frac{2}{3}$	$\frac{1}{2}$
A ₇	$\frac{2}{3}$	$\frac{1}{6}$
A ₈	1	0

Note that A_2 , A_4 , A_5 and A_6 are inadmissible. Also A_3 and A_7 are both maximum likelihood functions. Now for the decision function proposed above.

$$A = \begin{bmatrix} \frac{1}{2} & \frac{2}{5} & \frac{2}{3} \\ \\ \frac{1}{2} & \frac{3}{5} & \frac{1}{3} \end{bmatrix}$$

Sample Calculation

а

$$r_{22} = \frac{\frac{1}{2}}{\frac{1}{3} + \frac{1}{2}} = \frac{\frac{1}{2}}{\frac{5}{6}} = \frac{3}{5}$$

Now

$$\mathbf{R}: (\mathbf{A}, \, \boldsymbol{\theta}_{1}) = \mathbf{W}_{1}' \, \mathbf{AP}_{1} = \begin{bmatrix} 0 & 1 \end{bmatrix} \begin{bmatrix} \frac{1}{2} & \frac{2}{5} & \frac{2}{3} \\ & & \\ \frac{1}{2} & \frac{3}{5} & \frac{1}{3} \end{bmatrix} \begin{bmatrix} \frac{1}{3} \\ \frac{1}{3} \\ \frac{1}{3} \\ \frac{1}{3} \end{bmatrix}$$

$$= \frac{1}{6} + \frac{1}{5} + \frac{1}{9} = \frac{43}{90},$$

nd
$$R(A, \theta_2) = W'_2 AP_2 = \begin{bmatrix} 1 & 0 \end{bmatrix} \begin{bmatrix} \frac{1}{2} & \frac{2}{5} & \frac{2}{3} \\ & & \\ \frac{1}{2} & \frac{3}{5} & \frac{1}{3} \end{bmatrix} \begin{bmatrix} \frac{1}{3} \\ \frac{1}{2} \\ \frac{1}{6} \end{bmatrix}$$

 $=\frac{1}{6}+\frac{1}{5}+\frac{1}{9}=\frac{43}{90}$

=

and

so that A is admissible. Note that all losses were taken as equal in order to make things most favorable for a maximum likelihood decision function. If we make this assumption in general i.e., suppose W is a $q \ge 1$ vector with 0 in the ith position and the quantity W in all other positions, then we can see that proportional likelihood decision functions become dominated by maximum likelihood decision functions, and in fact have a risk of W for each θ_i as $q \rightarrow \infty$. We say, therefore, that they are asymptotically inadmissible. The proof is as follows: Let A be the proportional likelihood function

$$R(A, \theta_{i}) = W_{i}' A P_{i} = W_{i}' \left[\sum_{k=1}^{n} a_{jk} P_{ki} \right]$$
$$= W \sum_{j \neq i}^{n} \sum_{k=1}^{n} a_{jk} P_{ki}$$

$$= W \sum_{j \neq i}^{n} \sum_{k=1}^{n} \frac{p_{kj} p_{ki}}{\sum_{\ell=1}^{q} p_{k\ell}} = W \sum_{k=1}^{n} p_{ki} \sum_{j \neq i} \frac{p_{kj}}{\sum_{\ell=1}^{q} p_{k\ell}}$$
$$= W \sum_{k=1}^{n} p_{ki} \frac{p_{kj} p_{ki}}{\sum_{\ell=1}^{q} p_{kj} p_{ki}} = W \sum_{k=1}^{n} p_{ki} (1 - a_{ik}),$$

which approaches

$$W\sum_{k=1}^{n} p_{ki} = W$$

as $q \rightarrow \infty$ because

$$a_{ik} = \frac{p_{ki}}{\sum_{\ell=1}^{q} p_{k\ell}} \leq \frac{1}{\sum_{\ell=1}^{q} p_{k\ell}} \longrightarrow 0$$

as $q \rightarrow \infty$.

Now for a maximum likelihood decision function A, n

$$R(A, \theta_i) = W'_1 AP_i = W \sum_{j \neq i} \sum_{k=1}^{\infty} a_{jk} P_{ki}$$

where

$$a_{jk} = \begin{cases} 1. & \text{if } p_{kj} \ge p_{kl} & \text{for all } l \\ 0 & \text{otherwise} \end{cases}$$

Therefore,

$$R(A, \theta_i) = W \sum_{k=1}^{n} p_{ki} \sum_{j \neq i}^{n} a_{jk} = W \sum_{k=1}^{n} p_{ki} (1 - a_{ik})$$
$$= W \sum_{k=1}^{n} p_{ki} ,$$

where the sum is taken over all k such that $p_{ki} < p_{kl}$ for some l. This is less than W for at least one θ_i .

Proportional Bayes Decision Functions

We will assume that in our above example we have on a priori distribution given by

$$\mathbf{P} = \begin{bmatrix} \frac{1}{3} \\ \frac{2}{3} \end{bmatrix}$$

and we will show how to find the Bayes decision function.

If we plot $R(A, \theta_1)$ against $R(A, \theta_2)$ for all the admissible A's and join these points by a broken line we have



It can easily be shown that any admissible randomized decision function can be obtained from the non-randomized admissible ones in the following way: Select two non-randomized admissible functions, say A_3 and A_7 , which are jointed by a straight line segment. Choose A_3 with probability a and A_7 with probability 1 - a. This yields a randomized function A such-that

$$R(A, \theta_{i}) = aR(A_{3}, \theta_{i}) + (1 - a) R(A_{7}, \theta_{i}), i = 1, 2$$

The point $(R(A, \theta_1), R(A, \theta_2))$ lies on the line segment jointing A_3 and A_7 . Hence, the points of the entire broken line are the risk pairs for all admissible functions. To find which of these is the Bayes function we form the equation

$$\frac{1}{3}x + \frac{2}{3}y = k$$

where, x is the risk under θ_1 , and y is the risk under θ_2 , and let k vary from zero up until this line first touches our broken line of admissible functions. As we can see from the diagram below, A_7 is the Bayes decision function.



A problem arises in this procedure if the

$$P(\theta_1) x + P(\theta_2) y = k$$
 (1)

has a slope equal to that of one of the line segments of admissible functions. In such a case there will be an infinite number of Bayes functions. In our example this would happen if $P(\theta_1) = P(\theta_2) = \frac{1}{2}$. Then (1) has a slope of -1, and so does the line segment joining A_3 and A_7 . We now must choose one of the functions along this segment. To do this we argue as follows:

$$R(A_3, \theta_1) < R(A_7, \theta_1)$$

while

line

$$R(A_7, \theta_2) < R(A_3, \theta_2)$$

Therefore if θ_1 were the true state of nature, A_3 would be better, while if θ_2 were true, A_7 would be better. Hence, we propose choosing A_3 with probability $P(\theta_1)$ and A_7 with probability $P(\theta_2)$. It seems that the resulting function is in some sense better than either A_3 or A_7 , but it is not clear how this can be expressed mathematically.

The procedures discussed can all be generalized to more than two states of nature. We used the above example in the interests of clarity of exposition.