**1.2 Decision Theory**. Usually in statistics, instead of just two possible probability distributions $P, Q$, as in the last section, there is an infinite family $\mathcal{P}$ of such distributions, defined on a *sample space*, which is a measurable space $(X, \mathcal{B})$, in other words a set $X$ together with a $\sigma$-algebra $\mathcal{B}$ of subsets of $X$. As noted previously, if $X$ is a subset of a Euclidean space, then $\mathcal{B}$ will usually be the $\sigma$-algebra of Borel subsets of $X$. If $X$ is a countable set, then $\mathcal{B}$ will usually be the $\sigma$-algebra of all subsets of $X$ (if also $X \subset \mathbb{R}^k$, then all its subsets are in fact Borel sets). A probability measure on $\mathcal{B}$ will be called a *law*. The family $\mathcal{P}$ of laws on $(X, \mathcal{B})$ is usually written as $\{P_\theta, \ \theta \in \Theta\}$, where $\Theta$ is called a *parameter space*. For example, if $\mathcal{P}$ is the set of all normal measures $N(\mu, \sigma^2)$ for $\mu \in \mathbb{R}$ and $\sigma > 0$, we can take $\theta = (\mu, \sigma)$ or $(\mu, \sigma^2)$ where in either case $\Theta$ is the open upper half-plane, that is, the set of all $(t, u) \in \mathbb{R}^2$ such that $u > 0$. We assume that the function $\theta \mapsto P_\theta$ from $\Theta$ to laws on $\mathcal{B}$ is one-to-one, in other words $P_\theta \neq P_\phi$ whenever $\theta \neq \phi$ in $\Theta$. So the sets $\mathcal{P}$ and $\Theta$ are in 1-1 correspondence and any structure on one can be taken over to the other. We also assume given a $\sigma$-algebra $\mathcal{T}$ of subsets of $\Theta$. Most often $\Theta$ will be a subset of some Euclidean space and $\mathcal{T}$ the family of Borel subsets of $\Theta$. The family $\{P_\theta, \ \theta \in \Theta\}$ will be called *measurable* on $(\Theta, \mathcal{T})$ if and only if for each $B \in \mathcal{B}$, the function $\theta \mapsto P_\theta(B)$ is measurable on $\Theta$. If $\Theta$ is finite or countable, then (as with sample spaces) $\mathcal{T}$ will usually be taken to be the collection of all its subsets. In that case the family $\{P_\theta, \ \theta \in \Theta\}$ is always measurable.

An *observation* will be a point $x$ of $X$. Given $x$, the statistician tries to make inferences about $\theta$, such as estimating $\theta$ by a function $\hat{\theta}(x)$. For example, if $X = \mathbb{R}^n$ and $P_\theta = N(\theta, 1)^n$, so $x = (X_1, \ldots, X_n)$ where the $X_i$ are i.i.d. with distribution $N(\theta, 1)$, then $\hat{\theta}(x) = \overline{X} := (X_1 + \cdots X_n)/n$ is the classical estimator of $\theta$.

In decision theory, there is also a measurable space $(D, \mathcal{S})$, called the *decision space*. A measurable function $d(\cdot)$ from $X$ into $D$ is called a *decision rule*. Such a rule says that if $x$ is observed, then action $d(x)$ should be taken.

One possible decision space $D$ would be the set of all $d_\theta$ for $\theta \in \Theta$, where $d_\theta$ is the decision (estimate) that $\theta$ is the true value of the parameter. Or, if we just have a set $\mathcal{P}$ of laws, then $d_P$ would be the decision that $P$ is the true law. Thus in the last section we had $\mathcal{P} = \{P, Q\}$ and for non-randomized tests, $D = \{d_P, d_Q\}$. There, a decision rule is equivalent to a measurable subset of $X$, which was taken to be the set where the decision will be $d_Q$. For randomized rules, still for $\mathcal{P} = \{P, Q\}$, the decision space $D$ can be taken as the interval $0 \leq d \leq 1$, where $d(x)$ is the probability that $Q$ will be chosen if $x$ is observed.

Another possible decision space is a collection of subsets of the parameter space. Suppose $P_\theta = N(\theta, \sigma^2)^n$ on $X = \mathbb{R}^n$ for $-\infty < \theta < \infty$ where $\sigma$ is fixed and known. Then $[\overline{X} - 1.96\sigma/n^{1/2}, \overline{X} + 1.96\sigma/n^{1/2}]$ is a "95% confidence interval for $\theta$," meaning that for all $\theta$, $P_\theta\{|\overline{X} - \theta| > 1.96\sigma/n^{1/2}\} \doteq 0.05$. Here the decision space $D$ could be taken as the set of all closed intervals $[a, b] \subset \mathbb{R}$. Giving a confidence interval (or a "confidence set" more generally) is one kind of decision rule.

In decision theory, we also assume given a *loss function* $L$, which is a measurable function: $\Theta \times D \to [0, \infty]$. Here $L(\theta, d)$ is the loss suffered when the decision $d$ is taken and $\theta$ is the true value of the parameter, sometimes called the "state of nature." The

following condition on a loss function will be noted for later reference, though not always assumed:

(1.2.1)   $L(P, d_P) = 0$  for every  $P \in \mathcal{P}$,  in other words  $L(P_\theta, d_\theta) = 0$  for all  $\theta \in \Theta$,

that is, a correct decision incurs no loss. If the decision rule is an estimator $T(x) = \hat{\theta}(x)$ for a real parameter $\theta$, then one frequently used loss function is squared-error loss, $L(\theta, t) = (\theta - t)^2$.

Many authors treat a *utility function $U(\theta, d)$* rather than a loss function. Here larger values of $U$ are more favorable. If $L$ is a loss function, then $U = -L$ gives a utility function, but not necessarily conversely: in accord with broad usage of the term among statisticians and economists, a utility function may take values positive, negative or zero, and the values $U(P, d_P)$ may not be 0 and may be different for different $P$. To give a mathematical definition, a *utility function* is a measurable, extended real valued function (i.e. its values are in $[-\infty, \infty]$) on $\mathcal{P} \times D$, or equivalently on $\Theta \times D$. If for $P$ and $Q$ in $\mathcal{P}$, $d_P \in D$ and $d_Q \in D$, then it is assumed that

(1.2.2) $$U(P, d_Q) \leq U(P, d_P),$$

that is, it's better to make the right decision than a wrong one. Values $U(P, d_Q) = -\infty$ are allowed, corresponding to the notion that a wrong decision could lead to "ruin" or "death" of the decision maker and possibly others.

If $D = \{d_P : P \in \mathcal{P}\}$ and $U(P, d_P) < \infty$ for all $P$, we can get a corresponding loss function satisfying (1.2.1) by setting

$$L(P, d_Q) = U(P, d_P) - U(P, d_Q).$$

As this suggests, a loss function measures how far off a decision was, *relative* to other possible decisions. Such an evaluation seems natural for statistics. A utility function, by contrast, measures outcomes on a more absolute scale, incorporating the possibility that some values of $\theta$ are more favorable than others, as is natural in economics. Loss functions and, especially, utility functions, reflect the preferences of the individual making the decision.

The *risk* of a decision rule $e(\cdot)$ at $\theta$ is defined by

$$r(\theta, e) := \int L(\theta, e(x)) dP_\theta(x),$$

that is, risk is expected loss. If $f$ and $g$ are two decision rules, let $f \preceq g$ mean that $r(\theta, f) \leq r(\theta, g)$ for all $\theta$. Also, $f$ *improves on* $g$, written $f \prec g$, will mean that $f \preceq g$ and for some $\theta$, $r(\theta, f) < r(\theta, g)$. A decision rule $g$ is called *inadmissible* if there is some rule $f$ with $f \prec g$. If there is no such $f$, then $g$ is called *admissible*. A class $H$ of decision rules will be called *complete* if for every decision rule $g$ not in $H$ there is an $h \in H$ with $h \prec g$. If only $h \preceq g$, then $H$ is called *essentially complete*.

As in the case of deciding between two laws in Sec. 1.1, there are randomized decision rules. For such rules, the decision space doesn't consist only of definite, specific decisions

such as $d_P$. Instead, $D$ contains probability measures $\nu$ on a space $A$ of specific actions. If $x$ is observed, and $d(x) \in D$ is a law $\nu_x$ on $A$, one will choose an action $a \in A$ at random according to $\nu_x$, then take the action $a$. If $A = \{d_P, d_Q\}$, a law $\nu$ on $A$ is given by a number $y$ with $0 \le y \le 1$, where $y = \nu(d_Q) = 1 - \nu(d_P)$.

In Sec. 1.1, we saw that randomized rules gave us admissible tests of all possible sizes, but also that for Bayes rules it was not necessary to randomize (Remark 1.1.9).

A decision rule $e(\cdot)$, which may be randomized, is called *minimax* if it minimizes the maximum risk, or mathematically

$$\sup_\theta r(\theta, e) \;=\; \inf_{d(\cdot)} \sup_\theta r(\theta, d),$$

where the infimum is over all possible decision rules, possibly randomized.

Minimax and randomized rules are both important in another subject closely related to decision theory, *game theory*. Here there are actions $a \in A$, parameters $\theta \in \Theta$, and a utility function $U(\theta, a)$. In game theory, at any rate in the basic form of it being treated here, there is no sample space $X$, and so no observation $x$ nor laws $P_\theta$. Also, an intelligent opponent can choose $\theta$, possibly even knowing one's (randomized) decision rule, although not one's specific action $a$. Thus, if there are minimax decision rules, then such rules should be used.

For example, in the game of "scissors-stone-paper," each of two players has available three possible actions, "scissors", "stone" and "paper," and the two take actions simultaneously. Here scissor beats paper, paper beats stone, and stone beats scissors. Suppose the winner of one round gains \$1 and the loser loses \$1. If both players take the same action, the outcome is a draw (no money changes hands). Then for repeated plays, one should use a randomized rule: if one always takes the same action, or any predictable sequence of actions, the opponent can always win. The randomized rule of choosing each action with probability 1/3 has average winnings 0 against any strategy of the opponent, and this strategy is minimax and is the unique minimax strategy: any other randomized rule, if known to the opponent, can be defeated and result in an average net loss (this is left as a problem).

Even when minimax rules are not available, it may be that for each specific action $a \in A$, there is a large loss $L(\theta, a)$ for some $\theta$. Then, possibly, any non-randomized decision rule choosing $a \in A$ has a large risk for some $\theta$. As in the insurance business, risks above a certain size may be unacceptable, so that one may prefer to choose a randomized rule $d(\cdot)$ to keep $\sup_\theta r(\theta, d)$ from being too large. In insurance, unlike simpler problems, "minimax" seems too extreme a requirement: policies would only be sold if the buyers could be persuaded to pay more in premiums than they could ever receive in benefits!

For a measurable space $(A, \mathcal{E})$ of specific actions, let $D_\mathcal{E}$ be the set of all probability measures on $(A, \mathcal{E})$. On $D_\mathcal{E}$, let $S_\mathcal{E}$ be the smallest $\sigma$-algebra for which all the evaluations $\nu \mapsto \nu(B)$ are measurable for $B \in \mathcal{E}$. For a loss function $L$ on $\Theta \times A$, and $\nu \in D_\mathcal{E}$, the loss at $\theta$ and $\nu$ will be

(1.2.3) $$L(\theta, \nu) \;:=\; \int_A L(\theta, a) d\nu(a).$$

The next fact extends the joint measurability of $L$ from simple to randomized strategies. Recall that for two measurable spaces $(U, \mathcal{U})$ and $(V, \mathcal{V})$, a function $F$ on $U \times V$ is called

*jointly measurable* if and only if it is measurable for the product $\sigma$-algebra $\mathcal{U} \otimes \mathcal{V}$, which is the smallest $\sigma$-algebra of subsets of $U \times V$ containing all sets $B \times C$ for $B \in \mathcal{U}$ and $C \in \mathcal{V}$.

**1.2.4 Proposition**. If $L$ is nonnegative and jointly measurable on $\Theta \times A$, then as defined on $\Theta \times D_{\mathcal{E}}$ by (1.2.3), it is jointly measurable into $[0, \infty]$, for the $\sigma$-algebra $\mathcal{S}_{\mathcal{E}}$ on $D_{\mathcal{E}}$.

If $f$ is a jointly measurable real-valued function on $\Theta \times A$ and $W_f$ is the set of all $(\theta, \nu)$ such that $f(\theta, \nu) := \int f(\theta, a) d\nu(a)$ is well-defined (not $\infty - \infty$), then $W_f$ is a measurable set for the product $\sigma$-algebra and $(\theta, \nu) \mapsto f(\theta, \nu)$ is jointly measurable on it for $\mathcal{T}$ and $\mathcal{S}_{\mathcal{E}}$ into $[-\infty, \infty]$.

**Proof.** Measurability holds if $L = 1_{B \times C}$, where $1_{B \times C}(\theta, a) \equiv 1_B(\theta) 1_C(a)$ for $B \in \mathcal{T}$ and $C \in \mathcal{E}$, since then $L(\theta, \mu) = 1_B(\theta) \mu(C)$, which is measurable for $\mathcal{T} \otimes \mathcal{S}_{\mathcal{E}}$. Next, any finite union $W$ of sets $B_i \times C_i$ for $B_i \in \mathcal{T}$ and $C_i \in \mathcal{E}$ can be written as a finite, disjoint union of such sets (RAP, Prop. 3.2.2). Adding their indicators, we get the result for $W$.

The set of all such $W$ forms an algebra (RAP, Prop. 3.2.3). If $L_n$ is a uniformly bounded sequence of measurable functions on $\Theta \times A$ for which the conclusion holds, and $L_n {\uparrow} L$ or $L_n {\downarrow} L$, then it also holds for $L$. Since the smallest monotone class including an algebra is a $\sigma$-algebra (RAP, Theorem 4.4.2), the result holds for $1_H$ for all $H$ in the product $\sigma$-algebra $\mathcal{T} \otimes \mathcal{E}$. Then it holds for $L$ simple (a finite linear combination of such $1_H$), then for $L$ nonnegative and measurable by monotone convergence (RAP, Prop. 4.1.5 and Theorem 4.3.2).

Then for any measurable $f$ on $\Theta \times A$, we write as usual $f = f^+ - f^-$ where $f^+ := \max(f, 0)$. Then $W_f$ is the set of $(\theta, \nu)$ such that $\int f^+(\theta, a) d\nu(a)$ and $\int f^-(\theta, a) d\nu(a)$ are not both $+\infty$, so it is a product measurable set. Applying the result for $L \geq 0$ to $f^+$ and $f_-$, we get that on $W_f$, $f(\theta, \nu)$ is a difference of two nonnegative measurable functions, not both $+\infty$, so the difference is a measurable function into $[-\infty, \infty]$. $\square$

**Remark**. Although loss functions are assumed nonnegative, utility functions can be positive or negative, so Proposition 1.2.4 can be applied when $f$ is a utility function.

A function $\nu_{.} : x \mapsto \nu_x$ from $X$ into $D_{\mathcal{E}}$ will be called a *randomized decision rule* if it is measurable from $(X, \mathcal{B})$ to $(D_{\mathcal{E}}, \mathcal{S}_{\mathcal{E}})$. The *risk* of the rule for a given $\theta$ is

$$r(\theta, \nu_{.}) := \int L(\theta, \nu_x) dP_\theta(x),$$

which is a measurable function of $\theta$. The definitions of admissible and inadmissible rules extend directly to $D_{\mathcal{E}}$ in place of $A$. Conversely, a randomized decision rule may be viewed as a non-randomized one with the space $D_{\mathcal{E}}$ in place of the action space $A$, and the risk $r(\theta, \nu)$ in place of the loss $L(\theta, a)$. So, let's just consider ordinary decision rules $a(\cdot)$.

Given a prior distribution $\pi$ on $(\Theta, \mathcal{T})$, the *risk* of a decision rule $a(\cdot)$ is defined as $r(a(\cdot)) := \int \int L(\theta, a(x)) dP_\theta(x) d\pi(\theta)$. A decision rule $a(\cdot)$ will be called a *Bayes* rule if it attains the minimum risk (for the given $\pi$) over all possible decision rules, and if this minimum risk is finite.

Even for non-Bayesian statisticians, who don't believe priors should be used in practice, at least not in all cases, priors can be useful technically in showing that a decision rule is admissible:

**1.2.5 Theorem**. Suppose $a(\cdot)$ is a decision rule and for some prior $\pi$ on $\Theta$, $a(\cdot)$ is Bayes for $\pi$ and is unique Bayes, meaning that for any other Bayes decision rule $b(\cdot)$ for $\pi$, and all $\theta$, $a(x) = b(x)$ for $P_\theta$-almost all $x$. Then $a(\cdot)$ is admissible.

**Proof.** If there were a decision rule $c(\cdot) \prec a(\cdot)$, then clearly $c(\cdot)$ is also Bayes, but for some $\theta$, $c(x) \neq a(x)$ with positive probability for $P_\theta$, a contradiction. $\qquad\square$

**1.2.6 Theorem**. Suppose the parameter space $\Theta$ is countable. If $a(\cdot)$ is a decision rule and there is a prior $\pi$ on $\Theta$, positive at every point, such that $a(\cdot)$ is Bayes for $\pi$, then $a(\cdot)$ is admissible.

**Proof.** If $b(\cdot)$ were another decision rule with $b \prec a$, then the risk of $b(\cdot)$ for $\pi$ would be smaller than that of $a(\cdot)$, a contradiction. $\qquad\square$

A set $C$ in a Euclidean space $\mathbb{R}^k$ is called *convex* if whenever $x, y \in C$ and $0 \le t \le 1$ we have $tx + (1-t)y \in C$.

**1.2.7 Proposition**. Let $(X, \mathcal{B})$ be the sample space and $(A, \mathcal{E})$ the action space where $A$ is a convex, Borel measurable subset of a Euclidean space $\mathbb{R}^k$ with Borel $\sigma$-algebra $\mathcal{E}$. Let $\|a\| := (a_1^2 + \cdots + a_k^2)^{1/2}$. Let $x \mapsto \nu_x : X \to D_\mathcal{E}$ be a randomized decision rule. Then $B_{\nu.} := \{x \in X : \int \|a\| d\nu_x(a) < \infty\} \in \mathcal{B}$. Let $a(x) := \int a \, d\nu_x(a)$ for $x \in B_{\nu.}$, where integration of vectors is coordinatewise. Then $x \mapsto a(x)$ is measurable from $B_{\nu.}$ into $A$.

**Proof.** By Proposition 1.2.4 (for a fixed $\theta$), $\nu \mapsto \int \|a\| d\nu(a)$ is measurable from $(D_\mathcal{E}, \mathcal{S}_\mathcal{E})$ into $[0, \infty]$. Since $x \mapsto \nu_x$ is measurable from $(X, \mathcal{B})$ to $(D_\mathcal{E}, \mathcal{S}_\mathcal{E})$, it follows that $B_{\nu.} \in \mathcal{B}$. For $x \in B_{\nu.}$, the integral $a(x)$ is well-defined and has no infinite coordinates. We have $a(x) \in A$ by the proof of Jensen's inequality (RAP, 10.2.6). By Proposition 1.2.4 (again for a single $\theta$), $x \mapsto a(x)$ is measurable. $\qquad\square$

## PROBLEMS

1. Show that whenever the set of all admissible decision rules is essentially complete, it is actually complete and is the smallest complete class.

2. In the situation of Sec. 1.1 where decision rules are randomized tests of $P$ vs. $Q$,
   (a) Show that there is a smallest complete class and describe it. (Hint: use the result of Problem 1.)
   (b) Show that for some $P$ and $Q$ there is an essentially complete class smaller than the class in (a). In terms of $P$ and $Q$, describe an essentially complete class which is as small as possible.

3. Show that the set of all admissible decision rules may not be essentially complete and in fact, may be empty. Hint: without any sample or parameter space, let the action space be the set of positive real numbers $a$ with loss function $L(a) = 1/a$. Describe what are the complete and essentially complete classes in this case.

4. Prove that if a player of scissors-stone-paper has a known randomized strategy anything other than playing each action with probability $1/3$, the opponent can win on the average for some suitable strategy.

5. Let $P_\theta = N(\theta, 1)^n$ on $\mathbb{R}^n$ for $-\infty < \theta < \infty$. Let the action space be $\mathbb{R}$ and the loss function $L(\theta, a) = (a - \theta)^2$. Show that a decision rule (estimator) which is a linear function of $\overline{X}$, $d(x) = b\overline{X} + c$, is inadmissible if $|b| > 1$ but admissible if $b = 0$. (It is admissible for $b = 1$ and $c = 0$, but this problem doesn't ask for a proof of that.)

## NOTES TO SEC. 1.2

A classic reference on decision theory is Ferguson (1967). A more recent work is Berger (1980, 1985). Decision theory was also the subject of the last chapter in each of two general texts by other leading statisticians: Bickel and Doksum (1977) and Cox and Hinkley (1974). In the second edition of Bickel and Doksum (2001), decision theory is more integrated into the book, beginning with the first chapter, as it is here.

## REFERENCES

Berger, James O. (1980, 2d ed. 1985). *Statistical Decision Theory*. Springer, New York.

Bickel, Peter J., and Kjell A. Doksum (1977, 2001). See Sec. 1.1.

Cox, David R., and David V. Hinkley (1974). *Theoretical Statistics*. Chapman and Hall, London.

Ferguson, Thomas S. (1967). See Sec. 1.1.