

CHAPTER 3

Probability theory and static models

3.1 INTRODUCTION

This chapter lays the foundation for a model of stochastic events and processes, events and processes which deterministic models cannot adequately describe. We want such a model to account explicitly for the sources of uncertainty described in Chapter 1. To do so, we will first describe random events probabilistically, and then in the next chapter we will add dynamics to the model through a study of random processes.

Probability theory [3–5, 8–10] is basically addressed to assigning probabilities to events of interest associated with the outcomes of some experiment. The fundamental context of this theory is then a general space of outcomes, called a sample space. However, it is more convenient to work in a Euclidean space, so we consider a function, or mapping, from the sample space to Euclidean space, called a random variable. The basic entity for describing probabilities and random variables becomes the probability distribution function, or its derivative, the probability density function, which will exist for the problems of interest to us. These concepts are discussed in Sections 3.2 and 3.3. Since we will eventually want to generate probability information about certain variables of interest, based upon measurements of related quantities, conditional probability densities and functions of random variables will be of primary importance and are described in Sections 3.4 and 3.5.

Having a mathematical model of probability of events, one can consider the concept of expected values of a random variable, the average value of that variable one would obtain over the entire set of possible outcomes of an experiment. Again looking ahead to estimation of variables based on measurement data, the idea of conditional expectation arises naturally. Expectations of certain functions yield moments of random variables, a set of parameters

which describe the shape of the distribution or density function; these moments are most easily generated by characteristic functions. These topics are the subjects of Sections 3.6–3.8.

Because Gaussian random variables will form a basis of our system model, they are described in detail in Section 3.9. The initial model will involve linear operations on Gaussian inputs, so the results of such operations are discussed in Section 3.10. Finally, Section 3.11 solves the optimal estimation problem for cases in which a static linear Gaussian system model is an adequate description.

3.2 PROBABILITY AND RANDOM VARIABLES

Probability theory can be developed in an intuitive manner by describing probabilities of events of interest in terms of the *relative frequency of occurrence*. Using this concept, the probability of some event A , denoted as $P(A)$, can be generated as follows: if the event A is observed to occur $N(A)$ times in a total of N trials, then $P(A)$ is defined by

$$P(A) \triangleq \lim_{N \rightarrow \infty} \frac{N(A)}{N} \quad (3-1)$$

provided that this limit in fact exists. In other words, we conduct a number of experimental trials and observe the ratio of the number of times the event of interest occurs to the total number of trials. As we make more and more trials, if this ratio converges to some value, we call that value the probability of the event of interest.

Although this is a conceptually appealing basis for probability theory, it does not allow precise treatment of many problems and issues of direct importance to us. Modern probability theory is more rigorously based on an axiomatic definition of probability. This axiomatic definition must still be a valid mathematical model of empirically observed frequencies of occurrence, but it is meant to extract the essence of the ideas involved and to deal with them in a precise, rather than heuristic, manner.

To describe an experiment in precise terms, let Ω be the fundamental *sample space* containing all possible outcomes of the experiment conducted. Each single *elementary outcome* of the experiment is denoted as an ω ; these ω 's then are the elements of Ω : $\omega \in \Omega$. In other words, the sample space is just the collection of possible outcomes of the experiment, each of which being thought of as a point in that space Ω . Now let A be defined as a specific *event* of interest, a specific set of outcomes of the experiment. Thus, each such event A is a subset of Ω : $A \subset \Omega$. An event A is said to occur if the observed outcome ω is an element of A , if $\omega \in A$.

EXAMPLE 3.1 Consider two consecutive tosses of a fair coin. The sample space Ω is composed of four elements ω : if HT represents heads on the first throw and tails on the second, and

so forth, then the four possible elementary outcomes are HH, HT, TH, and TT. This is depicted schematically in Fig. 3.1. Ω is just the collection of those four outcomes.

Let us say we are interested in three events:

A_1 = at least one tail was thrown

A_2 = exactly one tail was thrown

A_3 = exactly two tails were thrown

Then A_1 , A_2 , and A_3 are subsets of Ω ; each is a collection of points ω . These are also depicted in Fig. 3.1.

Now suppose we conduct one trial of the experiment, and we observe $HT = \omega_2$. Then, since this point is an element of sets A_1 and A_2 , but not of A_3 , we say that the events A_1 and A_2 occurred on that trial, but event A_3 did not occur. ■

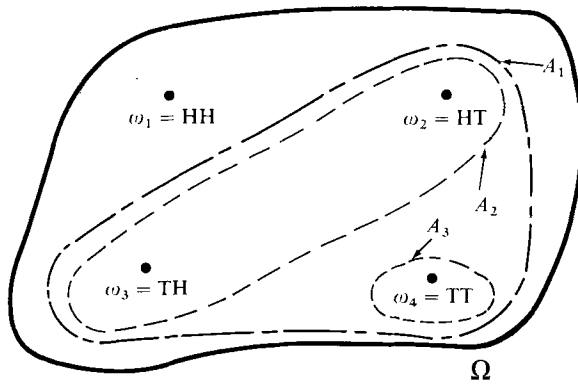


FIG. 3.1 Two tosses of a coin.

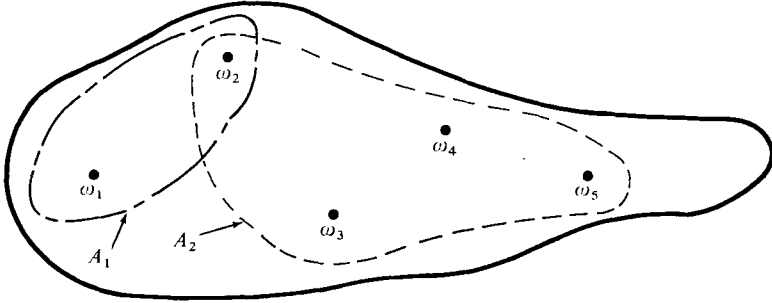
The sample space Ω can be discrete, with a finite or countably infinite number of elements, such as the space for the coin toss experiment described in the previous example. On the other hand, it could also be continuous, with an uncountable number of elements, such as the space appropriate to describe the continuous range of possible voltage values across a certain capacitor in a circuit.

So far, we have the structure of a sample space Ω , composed of elements ω_i , whose subsets are denoted as A_i . This is represented in Fig. 3.2.

Now we are going to restrict our attention to a certain class of sets A_i , a broad class called a σ -algebra, denoted as \mathcal{F} . In other words, the sets A_1 , A_2 , A_3, \dots , which will be admissible for consideration will be elements of the class $\mathcal{F}: A_i \in \mathcal{F}$.

A σ -algebra \mathcal{F} is a class of sets A_i , each of which is a subset of Ω ($A_i \subset \Omega$), such that if A_i is an element of \mathcal{F} (i.e., if $A_i \in \mathcal{F}$), then:

- (1) $A_i^* \in \mathcal{F}$, where A_i^* is the complement of A_i , $A_i^* = \Omega - A_i$,
- (2) $\Omega \in \mathcal{F}$ [and then the empty set $\emptyset \in \mathcal{F}$ also, due to the preceding (1)],

FIG. 3.2 Sample space Ω .

(3) if $A_1, A_2, \dots \in \mathcal{F}$, then their union and intersection are also in \mathcal{F} :

$$\bigcup_i A_i \in \mathcal{F} \quad \text{and} \quad \bigcap_i A_i \in \mathcal{F}$$

where all possible finite and countably infinite unions and intersections are included.

Whereas Ω is a collection of points (elementary outcomes, ω), \mathcal{F} is a collection of sets (events A_i), one of which is Ω itself.

For the purposes of our applications, we can let the sample space Ω be the set of points in n -dimensional Euclidean space R^n and let \mathcal{F} be the class of sets generated by sets of the form (each of which is a subset of Ω):

$$A = \{\omega: \omega \leq \mathbf{a}, \omega \in \Omega\} \quad (3-2)$$

and their complements, unions, and intersections. The notation in Eq. (3-2) requires some explanation. In other words, A is the set of ω 's that are elements of Ω (vectors in the n -dimensional Euclidean space, and thus, the boldfacing of ω to denote vector quantity; in the general case, ω will not be boldfaced); such that $(:)\omega \leq \mathbf{a}$, where ω and \mathbf{a} are n -dimensional vectors and \mathbf{a} is specified. Furthermore, $\omega \leq \mathbf{a}$ is to be interpreted componentwise: $\omega \leq \mathbf{a}$ means $\omega_1 \leq a_1$, $\omega_2 \leq a_2$, \dots , $\omega_n \leq a_n$ for the n components ω_i and a_i of ω and \mathbf{a} , respectively. This particular σ -algebra is of sufficient interest to have acquired its own name, and it is called a *Borel field*, denoted as \mathcal{F}_B . Taking complements, unions, and intersections of sets described by (3-2) leads to finite intervals (open, closed, half open) and point values along each of the n directions. Thus, a Borel field is composed of virtually all subsets of Euclidean n -space (R^n) that might be of interest in describing a probability problem associated with $\Omega = R^n$.

EXAMPLE 3.2 Consider generation of such sets of interest for $\Omega = R^1$, the real line. Let a_1 and a_2 be points on the real line, with $a_1 < a_2$. Then let

$$A_1 = \{\omega: \omega \leq a_1, \omega \in R^1\} = (-\infty, a_1]$$

$$A_2 = \{\omega: \omega \leq a_2, \omega \in R^1\} = (-\infty, a_2]$$

The complement of A_1 , which is also a member of \mathcal{F}_B by the definition of a σ -algebra, is

$$A_1^* = (-\infty, a_1]^* = (a_1, \infty)$$

Then the intersection of A_1^* and A_2 is

$$A_1^* \cap A_2 = (a_1, a_2]$$

Thus, we are able to generate any half-open interval, open on the left.

To generate points, we can look at a countably infinite intersection of half-open sets of the form

$$B_K = \{\omega : (b - \{1/K\}) < \omega \leq b\} = (b - \{1/K\}, b]$$

to generate

$$\bigcap_{K=1}^{\infty} B_K = \{\text{the point, } \omega = b\}$$

Note that we needed an infinite, not just finite, intersection to generate the point, and thus we needed to assure that such countable intersections yield sets in the σ -algebra when we first defined σ -algebra.

With a point b and a set $(b, c]$, we can generate the closed set $[b, c]$ by a simple union. Complementing and intersecting then yields open and half-open (open on the right) sets.

Thus, as claimed, the Borel field on the real line includes essentially all sets of possible interest. ■

Now define the *probability function* (or probability measure) $P(\cdot)$ to be a real scalar-valued function defined on \mathcal{F} that assigns a value, $P(A)$, to each A which is a member of \mathcal{F} ($A \in \mathcal{F}$) such that:

- (1) $P(A) \geq 0$ for all $A \in \mathcal{F}$,
- (2) $P(\Omega) = 1$,
- (3) if A_1, A_2, \dots are elements of \mathcal{F} and are disjoint, or mutually exclusive: i.e., if

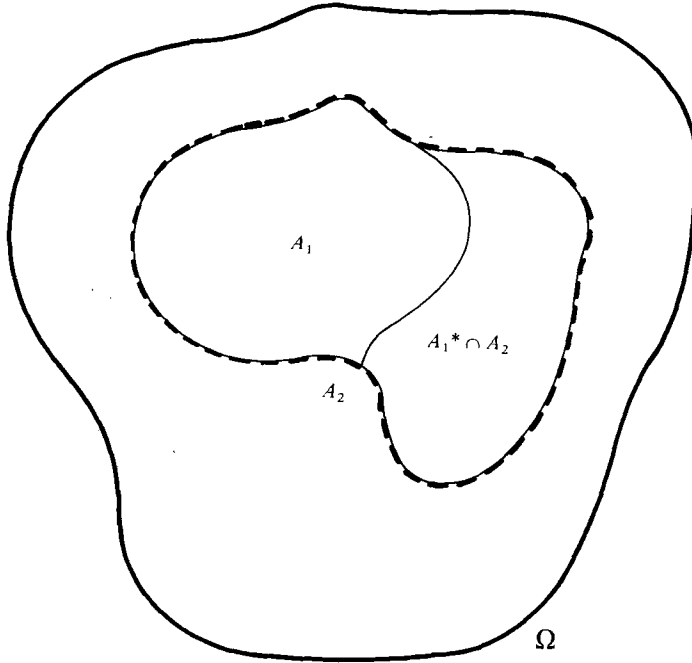
$$A_i \cap A_j = \emptyset \quad \text{for all } i \neq j$$

then

$$P\left(\bigcup_{i=1}^N A_i\right) = \sum_{i=1}^N P(A_i)$$

for all finite and countably infinite N .

Such a definition for the function to determine the probability of the events $A \in \mathcal{F}$ does correspond to one's intuition of probability gained through the concept of relative frequency of occurrence. Each set of interest (i.e., each $A \in \mathcal{F}$) is assigned a probability value between 0 and 1 (P is a mapping from \mathcal{F} into $[0, 1]$), and the probability of the sure event is one. Moreover, if A_1 is a subset of A_2 , then the probability of set A_2 is at least as great as the probability of A_1 , as expected. That this is a direct consequence of the axiomatic approach can be seen from Fig. 3.3. The set A_2 can be decomposed into two disjoint sets,

FIG. 3.3 $A_1 \subset A_2$ implies $P(A_1) \leq P(A_2)$.

A_1 and $A_1^* \cap A_2$. Then according to part (3) of the definition of a probability function, $P(A_2) = P(A_1) + P(A_1^* \cap A_2)$. From part (1), $P(A_1^* \cap A_2) \geq 0$, and so $P(A_2) \geq P(A_1)$ as desired.

Now we have what is called a *probability space*, defined by the triplet (Ω, \mathcal{F}, P) of the sample space, the underlying σ -algebra, and the probability function, all defined axiomatically as in the preceding. This entity serves as the basis of rigorously developed probability theory. Besides yielding results consistent with our intuitions about probability, this approach allows us to probe the essence of a problem and to determine whether or not it is posed properly. Subsequently, this rigor will also allow us to ensure that our definition of a random variable for a particular problem is in fact an appropriate choice for that application.

EXAMPLE 3.3 Let us consider the toss of a die to investigate a probability problem in the context of a rigorously defined probability space. Suppose that, for some reason, you are interested only in the occurrence of one of two events,

$$A_1 = \{\text{a 1 or a 2 was thrown}\} = \{1 \text{ or } 2\} \quad \text{and} \quad A_2 = \{\text{a 3 was thrown}\} = \{3\}$$

First of all, the sample space Ω is made up of the six possible outcomes: $\{1\}$, $\{2\}$, $\{3\}$, $\{4\}$, $\{5\}$, and $\{6\}$.

One possible means of generating a σ -algebra would be to let \mathcal{F}_1 be composed of \emptyset , $\Omega = \{1, 2, 3, 4, 5, \text{ or } 6\}$, the six elementary outcomes (ω 's), all possible unions of the outcomes two at a time, all possible unions three at a time, and so forth. A probability function would then have to assign a probability value to all such sets in \mathcal{F}_1 . However, if one is only interested in A_1 and A_2 just defined, there is no need to be able to assign probabilities to such sets as $\{3 \text{ or } 4 \text{ or } 5\}$.

Another means of generating an appropriate σ -algebra, a "minimal" σ -algebra for this particular example, would be to let \mathcal{F}_2 be composed of \emptyset , Ω , A_1 , A_2 , and all possible complements, unions, and intersections thereof. Thus, \mathcal{F}_2 is made up of \emptyset , Ω , $\{1 \text{ or } 2\}$, $\{3\}$, $\{1 \text{ or } 2\}^* = \{3 \text{ or } 4 \text{ or } 5 \text{ or } 6\}$, $\{3\}^* = \{1 \text{ or } 2 \text{ or } 4 \text{ or } 5 \text{ or } 6\}$, $\{1 \text{ or } 2\} \cup \{3\} = \{1 \text{ or } 2 \text{ or } 3\}$, $\{1 \text{ or } 2 \text{ or } 3\}^* = \{4 \text{ or } 5 \text{ or } 6\}$. Conceptually, these are the only sets for which a probability must be defined in order to determine solutions to any probability questions posed in terms of events A_1 and A_2 . Experiments could be conducted to assign probabilities to *only* these sets through the idea of relative frequency of occurrence, and the data generated would be complete.

Now let $P(A_1) = P(\{1 \text{ or } 2\}) = P_1$ and $P(A_2) = P(\{3\}) = P_2$. Then the probability function $P(\cdot)$ would assign probabilities to all of the elements of \mathcal{F} as follows:

$$\begin{aligned} P(\emptyset) &= 0 & P(A_1^*) &= P(\{3 \text{ or } 4 \text{ or } 5 \text{ or } 6\}) = 1 - P_1 \\ P(\Omega) &= 1 & P(A_2^*) &= P(\{1 \text{ or } 2 \text{ or } 4 \text{ or } 5 \text{ or } 6\}) = 1 - P_2 \\ P(A_1) &= P_1 & P(A_1 \cup A_2) &= P(\{1 \text{ or } 2 \text{ or } 3\}) = P_1 + P_2 \\ P(A_2) &= P_2 & P(\{A_1 \cup A_2\}^*) &= P(\{4 \text{ or } 5 \text{ or } 6\}) = 1 - P_1 - P_2 \end{aligned}$$

$P(A_1^*)$ is established by the fact that A_1 and A_1^* are disjoint sets whose union is Ω , so $P(A_1 \cup A_1^*) = 1 = P(A_1) + P(A_1^*) = P_1 + P(A_1^*)$; similarly for $P(A_2^*)$ and $P(\{A_1 \cup A_2\}^*)$. Since A_1 and A_2 are disjoint, $P(A_1 \cup A_2) = P(A_1) + P(A_2) = P_1 + P_2$. These are established by the axiomatic definitions, and would be verifiable by experimental observation (the theory is just abstractly modeling empirical results). ■

Once a probability space is properly defined for a given problem, the probability of all events of interest can be established, and theoretically we could be finished. The sample space Ω defines the possible outcomes of the experiment, \mathcal{F} is the collection of events (sets) of interest, and P assigns a probability to every one of these events. However, we can deal with numerical representations of sets in a space more readily than with the abstract subsets themselves. Consequently, for quantitative analysis, we need a mapping from the sample space Ω to the real numbers. It is for this reason that we introduce the concept of a random variable.

A scalar *random variable* $x(\cdot)$ is a real-valued point *function* which assigns a real scalar value to each point ω in Ω , denoted as $x(\omega) = x$, such that every set $A \subset \Omega$ of the form

$$A = \{\omega : x(\omega) \leq \xi\} \quad (3-3)$$

for ξ any value on the real line ($\xi \in R^1$), is an element of the σ -algebra \mathcal{F} (i.e., $A \in \mathcal{F}$). The name "random variable" is perhaps unfortunate in that it does not seem to imply the fact that we are talking about a function, as opposed to values the function can assume. In fact, $x(\cdot)$ is a function, or mapping, from Ω into R^1 .

The notation used warrants discussion. Random variables will be set in sans serif type, $x(\cdot)$ or simply x , to emphasize the fact that they are functions of

point values ω from the sample space Ω . The value that this function assumes for a particular ω , a *realization* of the random variable, will be the corresponding italicized letter. The corresponding Greek symbol will be used to denote a given vector or dummy variable (as, for integration), in the space of realizations. Thus, the notation $\{\omega: \mathbf{x}(\omega) \leq \xi\}$ is meant to read, "the set of ω in Ω such that the values assumed by the random variable function $\mathbf{x}(\cdot)$, for those ω as its argument, $\mathbf{x}(\omega) = \mathbf{x}$, are less than or equal to the given number ξ on the real line."

A *vector random variable* or *random vector* $\mathbf{x}(\cdot)$ is just the generalization of the random variable concept to the vector case: a real-valued point function which assigns a real vector value to each point ω in Ω , denoted as $\mathbf{x}(\omega)$, such that every set A of the form

$$A = \{\omega: \mathbf{x}(\omega) \leq \xi\} \quad (3-4)$$

for any $\xi \in R^n$, is an element of \mathcal{F} . Although these definitions might at first seem contorted, there is good reason for their form. Scalar random variables are specifically mappings from Ω into R^1 such that inverse images of half-open intervals of the form $(-\infty, \xi]$ in R^1 are events in Ω that belong to \mathcal{F} . That is to say, they are events in Ω for which probabilities have been defined through the probability function P . Vector random variables are simply extensions of the same idea—mappings from Ω into R^n such that inverse images of sets of the form $\{\mathbf{x}(\omega) \in R^n: -\infty < x_i(\omega) \leq \xi_i; i = 1, 2, \dots, n\}$ are events in Ω to which probabilities have been ascribed. (From a measure theoretic point of view, this just says that random variables are measurable functions.)

EXAMPLE 3.4 Perhaps the best way to understand the concept of a random variable is to consider a function that is *not* a random variable for a given problem. Let Ω be the interval $(0, 10]$ on the real line, and suppose we are interested in distinguishing whether ω takes on a value in the interval $I_1 = (0, 5]$ or in $I_2 = (5, 10]$. The minimal σ -algebra \mathcal{F} is made up of all possible complements, unions, and intersections of these two intervals, so that

$$\mathcal{F} = \{\emptyset, \Omega = (0, 10], I_1 = (0, 5], I_2 = (5, 10]\}$$

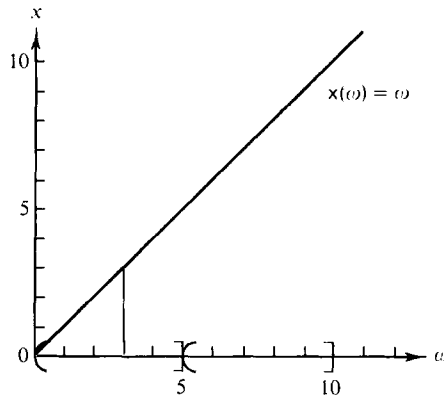


FIG. 3.4 Function that is not a random variable.

The value of ω can be anywhere along the line segment $(0, 10]$, and we just want to tell in which half of the segment it lies.

Now we want to establish an appropriate form for the function x to assume. Try defining $x(\cdot)$ such that $x(\omega) = \omega$, as in Fig. 3.4. This is not a suitable choice. Choose, for example, $\xi = 3$, as shown. Then the set A defined by

$$A = \{\omega: x(\omega) \leq 3\} = (0, 3]$$

is *not* an element of the class \mathcal{F} . By definition, a random variable x *must* be defined such that all sets of the form

$$A = \{\omega: x(\omega) \leq \xi\}$$

are in \mathcal{F} , for any choice of $\xi \in R^1$.

For this example, we must define x as assuming a constant value over $(0, 5]$ and a (different) constant over $(5, 10]$. One such random variable is shown in Fig. 3.5. Note for instance, that for this definition of x ,

$$A_1 = \{\omega: x(\omega) \leq 3\} = \emptyset; \quad A_2 = \{\omega: x(\omega) \leq 6\} = (0, 5]; \quad A_3 = \{\omega: x(\omega) \leq 20\} = (0, 10]$$

are all elements of \mathcal{F} . ■

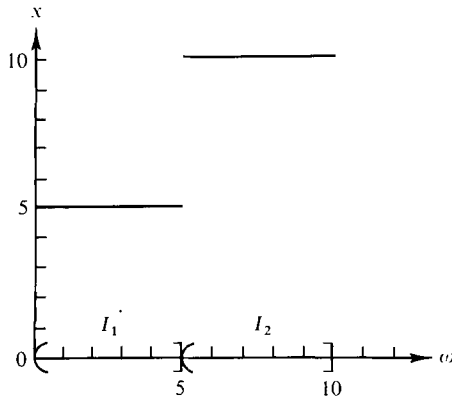


FIG. 3.5 Random variable definition. $x(\omega) = 5$ if $\omega \in I_1$ and $x(\omega) = 10$ if $\omega \in I_2$.

The sets I_1 and I_2 in the preceding example were irreducible elements of the σ -algebra \mathcal{F} , and a proper definition of a random variable required the function x to assume constant values over these sets. To generalize this concept, we call the set (event) $A \subset \Omega$ an “atom” of the σ -algebra \mathcal{F} if $A \in \mathcal{F}$ and no subset of A is an element of \mathcal{F} other than A itself and the null set \emptyset . A random variable can only assume a single value on an atom of the underlying σ -algebra \mathcal{F} .

Thus, we have the relationships depicted in Fig. 3.6. A random variable is a mapping from the fundamental sample space Ω into Euclidean n -space R^n . Each atom in Ω is mapped into a single vector in R^n . Conversely, the inverse image of sets in R^n of the form

$$A_i = \{x(\omega) \in R^n: -\infty < x_i(\omega) \leq \xi_i; i = 1, 2, \dots, n\}$$

are events in Ω ($A_i \subset \Omega$) for which probabilities have been defined ($A_i \in \mathcal{F}$).

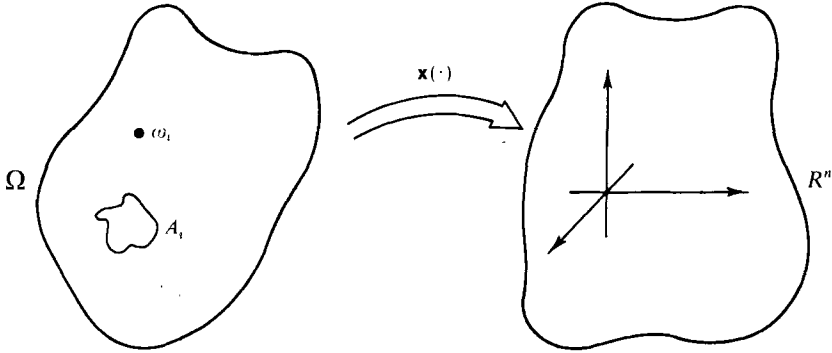


FIG. 3.6 The random variable mapping.

For the problems we will address in the sequel, the sample space Ω is R^n itself and the underlying σ -algebra is the Borel field \mathcal{F}_B generated by sets of the form $A_i = \{\omega: \omega \leq \mathbf{a}, \omega \in \Omega\}$. An appropriate random variable definition for this case is simply the identity mapping suggested in Example 3.4:

$$\mathbf{x}(\omega) = \omega \quad (3-5)$$

Note that an atom in $\Omega = R^n$ is just a single point in the space (a single vector), and that the random variable just mentioned does map each such atom into a single vector in R^n . Thus, each realization $\mathbf{x}(\omega)$ is an n -dimensional vector whose components can take on any value in $(-\infty, \infty)$.

By the definition of a random variable \mathbf{x} , all sets of the form

$$A = \{\omega: \mathbf{x}(\omega) \leq \xi\} = \{\omega: x_1(\omega) \leq \xi_1, x_2(\omega) \leq \xi_2, \dots, x_n(\omega) \leq \xi_n\}$$

have probabilities: probabilities are defined for them because $A \subset \Omega$ and $A \in \mathcal{F}$, and $P(\cdot)$ assigns probabilities to all such sets A . Therefore, the *probability distribution function* $F_{\mathbf{x}}(\cdot)$, a real scalar-valued function defined by

$$F_{\mathbf{x}}(\xi) = P(\{\omega: \mathbf{x}(\omega) \leq \xi\}) \quad (3-6a)$$

$$= "P(\mathbf{x} \leq \xi)" \quad (3-6b)$$

$$= "P(x_1 \leq \xi_1, x_2 \leq \xi_2, \dots, x_n \leq \xi_n)" \quad (3-6c)$$

always exists. We have defined the various sets and functions to this point so as to assure that such a function exists. The quotation marks in (3-6) are meant to emphasize that such notation, very typical in probability theory literature, should be interpreted in terms of the probability of a set of ω 's in the original sample space Ω . Sometimes the notation $F(\xi)$ is used rather than $F_{\mathbf{x}}(\xi)$, if the random variable concerned is obvious from context. Moreover, since

$$F_{\mathbf{x}}(\xi) = F_{x_1, x_2, \dots, x_n}(\xi_1, \xi_2, \dots, \xi_n) \quad (3-7)$$

this is sometimes called the joint probability distribution function of x_1, x_2, \dots , and x_n .

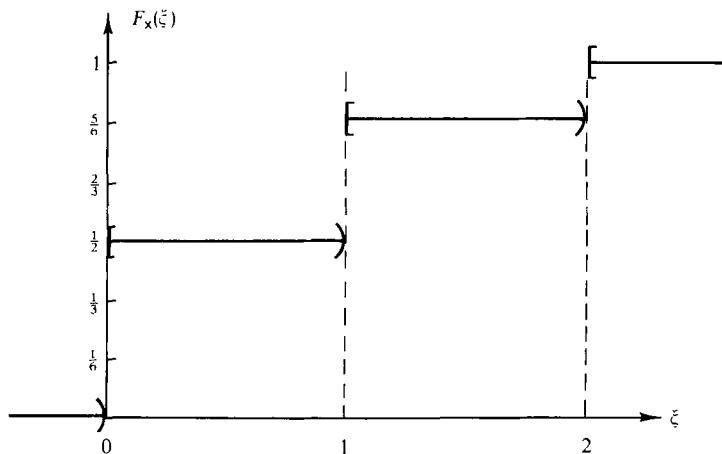


FIG. 3.7 Probability distribution function.

EXAMPLE 3.5 Consider the die-toss experiment introduced in Example 3.3, in which we were interested only in the sets $A_1 = \{1 \text{ or } 2\}$ and $A_2 = \{3\}$. Define a random variable x through

$$x(\omega) = \begin{cases} 0 & \text{if } \omega \notin A_1 \text{ or } A_2 \\ 1 & \text{if } \omega \in A_1 \\ 2 & \text{if } \omega \in A_2 \end{cases}$$

As in Example 3.3, let $P(A_1) = P_1$ and $P(A_2) = P_2$ ($\frac{1}{3}$ and $\frac{1}{6}$, respectively, for a fair die).

Now we will establish the probability distribution function

$$F(\xi) = P(\{\omega: x(\omega) \leq \xi\})$$

For $\xi < 0$, $P(\{\omega: x(\omega) \leq \xi\}) = P(\emptyset) = 0$.

For $0 \leq \xi < 1$, $P(\{\omega: x(\omega) \leq \xi\}) = P(\{A_1 \cup A_2\}^c) = 1 - P_1 - P_2 = \frac{1}{2}$.

For $1 \leq \xi < 2$, $P(\{\omega: x(\omega) \leq \xi\}) = P(A_2^c) = 1 - P_2 = \frac{5}{6}$.

For $2 \leq \xi < \infty$, $P(\{\omega: x(\omega) \leq \xi\}) = P(\Omega) = 1$.

Plotting $F_x(\xi)$ versus ξ yields the graphical depiction of the probability distribution function in Fig. 3.7. ■

Figure 3.8 summarizes the concepts that have been discussed. We started with an abstract sample space Ω , composed of elements (points) ω that were the elementary outcomes of an experiment. There were also certain subsets A of Ω ($A \subset \Omega$) of interest, called events, and specifically these sets were from a class \mathcal{F} ($A \in \mathcal{F}$) called a σ -algebra. For any $A \in \mathcal{F}$, we could evaluate the *set function* $P(\cdot)$, a mapping from \mathcal{F} into $[0, 1]$, to generate probabilities as $P(A)$. The triplet (Ω, \mathcal{F}, P) then defined what was termed a probability space.

We also defined a *point function* $x(\cdot)$ called a random variable, a mapping from Ω into R^n , which could be evaluated for each $\omega \in \Omega$ to yield realizations $x(\omega)$. The probabilities established as $P(A)$ and the realizations $x(\omega)$ of the random variable x are then related by the probability distribution function $F_x(\cdot)$, a mapping from R^n into $[0, 1]$, that yields $F_x(\xi)$ as the probability of the set of $\omega \in \Omega$ such that $x(\omega) \leq \xi$.

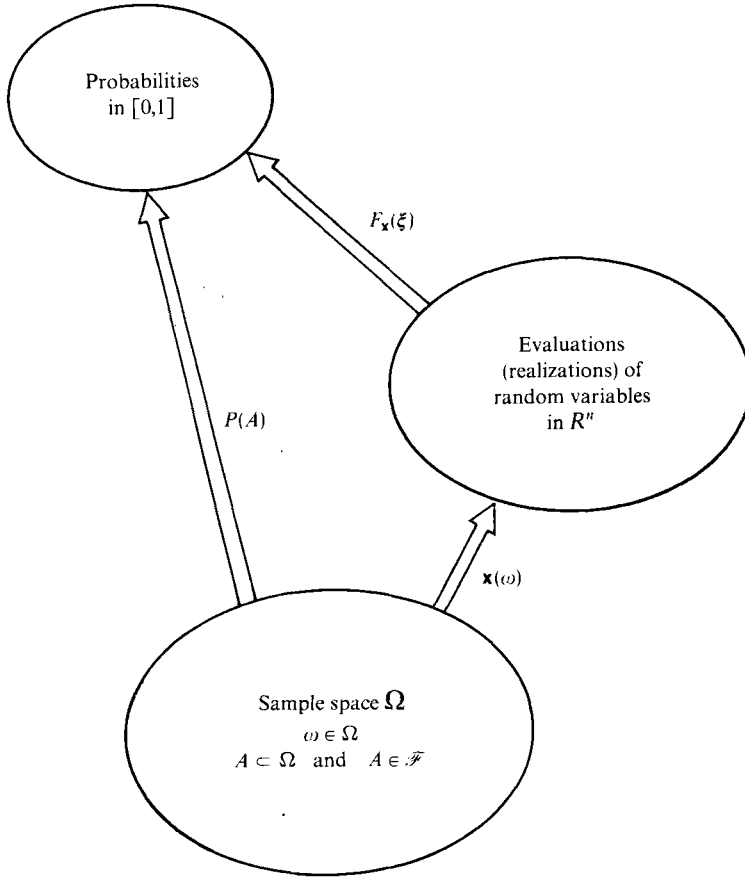


FIG. 3.8 Probability and random variables.

3.3 PROBABILITY DISTRIBUTIONS AND DENSITIES

As discussed in the previous section, the probability distribution function is a basic entity associated with any random variable that allows us to generate probabilities of sets of interest. We are assured of its existence. On the other hand, we are not assured of the existence of its derivative everywhere, but if it does exist, it is often easier to use and more revealing in terms of graphical interpretations.

Given a vector random variable \mathbf{x} ,

$$\mathbf{x} \triangleq \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \quad (3-8)$$

the probability distribution function $F_{\mathbf{x}}$ can be evaluated as a scalar function of the dummy vector $\xi = [\xi_1, \xi_2, \dots, \xi_n]^T$:

$$F_{\mathbf{x}}(\xi) \triangleq F_{x_1, x_2, \dots, x_n}(\xi_1, \xi_2, \dots, \xi_n) \quad (3-9a)$$

$$\triangleq P(\{\omega: x_1(\omega) \leq \xi_1, x_2(\omega) \leq \xi_2, \dots, x_n(\omega) \leq \xi_n\}) \quad (3-9b)$$

Note again that we specifically avoid the notation $F_{\mathbf{x}}(\mathbf{x})$, as is often used, to prevent giving the impression that $F_{\mathbf{x}}$ is in any way a function of particular realizations \mathbf{x} of \mathbf{x} . As can be seen from Eq. (3-9), $F_{\mathbf{x}}$ is a monotonic nondecreasing function of any component ξ_i of the vector ξ : for instance, the probability of the set of ω such that $x_i(\omega) \leq 2$ must be at least as large as the probability of the set of ω such that $x_i(\omega) \leq 1$.

Other properties of this function that become apparent from its definition include:

$$F_{\mathbf{x}}(\infty, \infty, \dots, \infty) = P(\{\omega: x_1(\omega) \leq \infty, \dots, x_n(\omega) \leq \infty\}) = 1 \quad (3-10)$$

$$F_{\mathbf{x}}(\xi_1, \dots, -\infty, \dots, \xi_n) = P(\{\omega: \dots, x_i(\omega) \leq -\infty, \dots\}) = 0 \quad (3-11)$$

If all of its arguments are ∞ , the value $F_{\mathbf{x}}$ assumes is one; if any single argument is $-\infty$, it takes on the value zero (these statements are more properly expressed in the limit as certain arguments tend to $+\infty$ or $-\infty$). If we are interested only in probabilities concerning the first k of the n random variables, and x_{k+1}, \dots, x_n can take on any values, then:

$$\begin{aligned} F_{x_1, \dots, x_k}(\xi_1, \dots, \xi_k) &= P(\{\omega: x_1(\omega) \leq \xi_1, \dots, x_k(\omega) \leq \xi_k\}) \\ &= P(\{\omega: x_1(\omega) \leq \xi_1, \dots, x_k(\omega) \leq \xi_k, \\ &\quad x_{k+1}(\omega) \leq \infty, \dots, x_n(\omega) \leq \infty\}) \\ &= F_{x_1, \dots, x_n}(\xi_1, \dots, \xi_k, \infty, \dots, \infty) \end{aligned} \quad (3-12)$$

Equation (3-12) embodies the concept of a “*marginal*” probability distribution of x_1, x_2, \dots, x_k . Note that the first k components were chosen only so Eq. (3-12) could be written readily; any argument of ∞ in $F_{\mathbf{x}}(\xi)$ yields the corresponding result [the variables can be reordered so there is no loss of generality in (3-12)].

The probability distribution function can be used to generate probabilities of other sets of interest as well. For instance, in the scalar case, the probability of sets of ω such that x assumes a value in a given half-open interval $(\xi_1, \xi_2]$ can be readily established. The set $\{\omega: x(\omega) \leq \xi_2\}$ can be decomposed into the union of two disjoint sets:

$$\{\omega: x(\omega) \leq \xi_2\} = \{\omega: x(\omega) \leq \xi_1\} \cup \{\omega: x(\omega) \in (\xi_1, \xi_2]\}$$

Because the sets on the right are disjoint, we have

$$P(\{\omega: x(\omega) \leq \xi_2\}) = P(\{\omega: x(\omega) \leq \xi_1\}) + P(\{\omega: x(\omega) \in (\xi_1, \xi_2]\})$$

so we can write

$$P(\{\omega: \mathbf{x}(\omega) \in (\xi_1, \xi_2]\}) = F_{\mathbf{x}}(\xi_2) - F_{\mathbf{x}}(\xi_1) \quad (3-13)$$

To generate probabilities of open or closed sets, we need to evaluate probabilities that \mathbf{x} assumes a single value. If the distribution function is discontinuous at some ξ_0 , then there is a finite probability that $\mathbf{x}(\cdot)$ assumes that value. In (3-13), let $\xi_2 = \xi_0$ and $\xi_1 = \xi_0 - \varepsilon$, and take the limit as $\varepsilon \rightarrow 0$ to yield

$$P(\{\omega: \mathbf{x}(\omega) = \xi_0\}) = F_{\mathbf{x}}(\xi_0) - F_{\mathbf{x}}(\xi_0^-) \quad (3-14)$$

i.e., the probability is equal to the magnitude of the jump discontinuity. We have observed such discontinuities in Example 3.5. Thus, for instance, since we have disjoint sets,

$$\begin{aligned} P(\{\omega: \mathbf{x}(\omega) \in [\xi_1, \xi_2]\}) &= P(\{\omega: \mathbf{x}(\omega) = \xi_1\}) + P(\{\omega: \mathbf{x}(\omega) \in (\xi_1, \xi_2]\}) \\ &= [F_{\mathbf{x}}(\xi_1) - F_{\mathbf{x}}(\xi_1^-)] + [F_{\mathbf{x}}(\xi_2) - F_{\mathbf{x}}(\xi_1)] \\ &= F_{\mathbf{x}}(\xi_2) - F_{\mathbf{x}}(\xi_1^-) \end{aligned} \quad (3-15)$$

We will discuss the generation of probabilities for general sets of interest after we introduce the concept of a probability density function.

If a scalar-valued function $f_{\mathbf{x}}(\cdot)$ exists such that

$$F_{\mathbf{x}}(\xi_1, \xi_2, \dots, \xi_n) = \int_{-\infty}^{\xi_1} \int_{-\infty}^{\xi_2} \dots \int_{-\infty}^{\xi_n} f_{\mathbf{x}}(\rho_1, \rho_2, \dots, \rho_n) d\rho_1 d\rho_2 \dots d\rho_n \quad (3-16a)$$

or, in a simpler notation to represent the same expression,

$$F_{\mathbf{x}}(\xi) = \int_{-\infty}^{\xi} f_{\mathbf{x}}(\rho) d\rho \quad (3-16b)$$

holds for all values of $\xi = [\xi_1, \xi_2, \dots, \xi_n]^T$, then this function $f_{\mathbf{x}}$ is the *probability density function* of \mathbf{x} . Unlike the probability distribution function, we are not always assured of the existence of $f_{\mathbf{x}}$. If $F_{\mathbf{x}}$ is absolutely continuous, then the density function does exist (absolute continuity can be defined rigorously through measure theory, but basically $F_{\mathbf{x}}$ is absolutely continuous if the number of points where it is not differentiable is countable). If such a density function exists, then \mathbf{x} is termed a *continuous random variable*.

By the fundamental theorem of calculus, we can use (3-16) to deduce

$$f_{\mathbf{x}}(\xi) = \frac{\partial^n}{\partial \xi_1 \partial \xi_2 \dots \partial \xi_n} F_{\mathbf{x}}(\xi) \quad (3-17)$$

This relationship and (3-16), combined with properties of $F_{\mathbf{x}}$, yield some properties of $f_{\mathbf{x}}$. Since $F_{\mathbf{x}}$ is monotonic nondecreasing,

$$f_{\mathbf{x}}(\xi) \geq 0 \quad \text{for all } \xi \quad (3-18)$$

In view of (3-10), it is a property of a density function that

$$\int_{-\infty}^{\infty} f_{\mathbf{x}}(\xi) d\xi = 1 \quad (3-19)$$

If we are interested only in the first k of the n components of \mathbf{x} , and $\mathbf{x}_{k+1}, \dots, \mathbf{x}_n$ can take on any values, then we can establish the *marginal density function* by integrating out the dependence upon the last $(k - n)$ components:

$$f_{\mathbf{x}_1, \dots, \mathbf{x}_k}(\xi_1, \dots, \xi_k) = \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} f_{\mathbf{x}_1, \dots, \mathbf{x}_n}(\xi_1, \dots, \xi_k, \dots, \xi_n) d\xi_{k+1} \cdots d\xi_n \quad (3-20)$$

as can be seen from (3-12). With $F_{\mathbf{x}}$ continuous, (3-14) yields

$$P(\{\omega: \mathbf{x}(\omega) = \xi_0\}) = 0 \quad (3-21)$$

so that, using (3-13) and (3-15),

$$P(\{\omega: \mathbf{x}(\omega) \in [\xi_1, \xi_2] \text{ or } [\xi_1, \xi_2]\}) = F_{\mathbf{x}}(\xi_2) - F_{\mathbf{x}}(\xi_1) = \int_{\xi_1}^{\xi_2} f_{\mathbf{x}}(\rho) d\rho \quad (3-22)$$

with extensions to the vector case.

EXAMPLE 3.6 Two forms of random variables useful for modeling empirically observed phenomena are the uniformly and Gaussian (or normally) distributed random variables. In the scalar case, their probability distribution and density functions can be plotted as in Fig. 3.9.

The uniformly distributed random variable models a situation in which a quantity of interest can take on any value in a specified range (limited by physical considerations as gimbal stops on a servo motor, by definition of units as angular orientation being described in the range $[0, 2\pi)$, etc.) and in which there is no reason to believe certain ranges of values to be more probable than others. The Gaussian, or normal, random variable serves as a good model for many observed phenomena and will be discussed at length in Section 3.9.

Note that, analogous to a mass density function, the probability density function indicates where the probability (mass) is concentrated. It is partly this graphic portrayal of probable ranges of values that makes the density function more attractive to use than the distribution function. ■

Now we want to obtain the probability of general sets of interest associated with vector random variables. In the scalar case, we can see from (3-22) that the probability of the set of ω such that $\mathbf{x}(\omega)$ lies in the infinitesimal interval from ξ_1 to $(\xi_1 + d\xi_1)$ is just

$$P(\{\omega: \mathbf{x}(\omega) \in [\xi_1, \xi_1 + d\xi_1]\}) = f_{\mathbf{x}}(\xi_1) d\xi_1 \quad (3-23a)$$

This generalizes to the n -dimensional case as

$$P(\{\omega: \mathbf{x}_i(\omega) \in [\xi_i, \xi_i + d\xi_i]; i = 1, 2, \dots, n\}) = f_{\mathbf{x}}(\xi_1, \dots, \xi_n) d\xi_1 \cdots d\xi_n \quad (3-23b)$$

Thus, the probability that the random variable takes on a value within an infinitesimal hypercube is just the *probability (mass) density* evaluated at the location of the hypercube (it is constant over the *infinitesimal* volume) multiplied by the *volume* of the hypercube, $[d\xi_1 \cdots d\xi_n]$. Then any set A in R^n can be generated from such hypercubes, and so the probability that $\mathbf{x}(\omega) = \mathbf{x}$ lies in

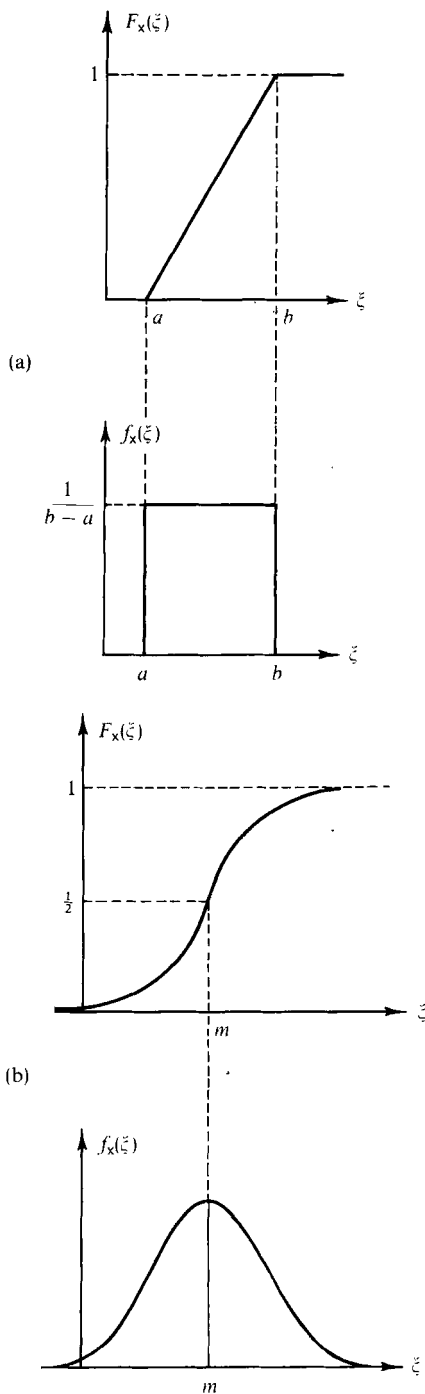


FIG. 3.9 (a) Uniform and (b) Gaussian (normal) distributions and densities.

the set A is

$$P(\{\omega: \mathbf{x}(\omega) \in A\}) = \int_A f_{\mathbf{x}}(\xi) d\xi \quad (3-24a)$$

$$= \int \cdots \int_A f_{\mathbf{x}}(\xi_1, \dots, \xi_n) d\xi_1 \cdots d\xi_n \quad (3-24b)$$

In this manner, the probability associated with sets of interest can be generated through ordinary (Riemann) integration with the probability density function.

To think of this geometrically, consider Fig. 3.10. In case (a), \mathbf{x} is scalar and $f_{\mathbf{x}}(\xi)$ is a simple curve plotted against ξ . Sets A of interest are intervals along the abscissa (possibly disjoint), and $P(\{\omega: \mathbf{x}(\omega) \in A\})$ can be determined as the area under the curve delimited by A . If \mathbf{x} is two dimensional, as in case (b), $f_{\mathbf{x}}(\xi_1, \xi_2)$ is a surface over the ξ_1 - ξ_2 plane; sets A are areas in the ξ_1 - ξ_2 plane, and $P(\{\omega: \mathbf{x}(\omega) \in A\})$ can be calculated as the volume under the surface $f_{\mathbf{x}}(\xi_1, \xi_2)$ with A as a cross section.

If we want to calculate the probability of general sets without use of the density function (as, for cases in which $F_{\mathbf{x}}$ has discontinuities so $f_{\mathbf{x}}$ cannot be

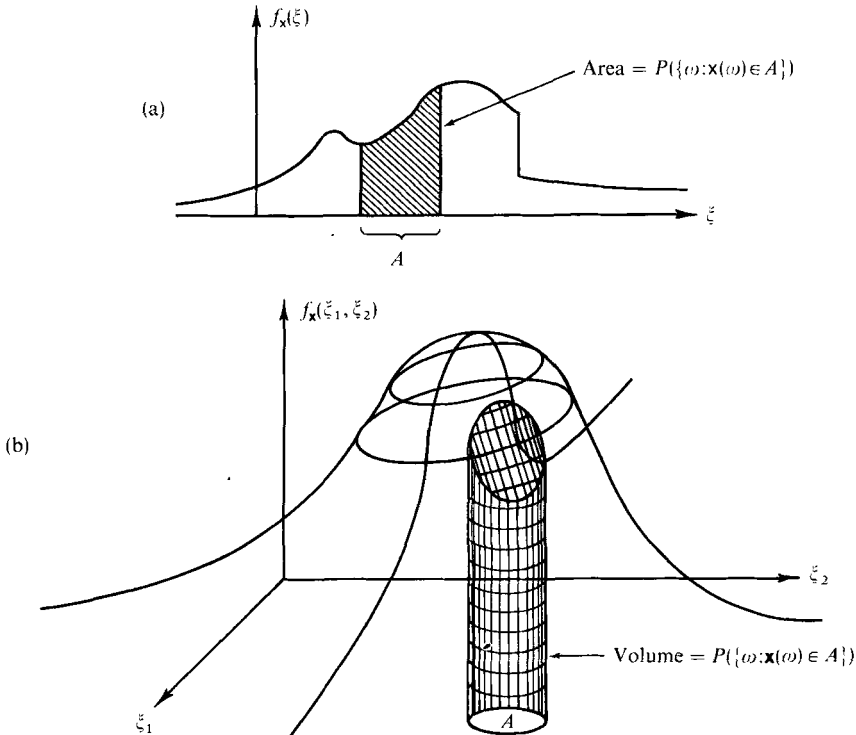


FIG. 3.10 Probability of sets. (a) Scalar-valued \mathbf{x} . (b) Two-dimensional vector-valued \mathbf{x} .

defined everywhere), the preceding ideas can be extended through

$$P(\{\omega: \mathbf{x}(\omega) \in A\}) = \int_A dF_{\mathbf{x}}(\xi) \quad (3-25)$$

Meaning is given to this expression through measure theory: we need to define this Lebesgue–Stieltjes integral of $F_{\mathbf{x}}$ over A . If $f_{\mathbf{x}}$ exists, then (3-24) and (3-25) yield the same result, with (3-24) being more attractive because it is in terms of ordinary Riemann integration. For our applications, we will be able to assume the existence of $f_{\mathbf{x}}$. We will therefore be able to avoid measure theory considerations, but such an extension *can* be made.

Let us reflect on what has been accomplished through the random variable mapping \mathbf{x} . Recall Fig. 3.6: \mathbf{x} maps Ω into R^n such that each atom (each irreducible set) in Ω maps into a single vector in R^n . Thus, the sets of interest in R^n will be elements of the Borel field \mathcal{F}_B associated with R^n . For all sets $A \subset R^n$ and $A \in \mathcal{F}_B$, we can define probabilities through

$$P_{\mathbf{x}}(A) = \int_A dF_{\mathbf{x}}(\xi) \quad (3-26)$$

where $P_{\mathbf{x}}(\cdot)$ is the probability function (Borel measure) associated with R^n , or, if the probability density function exists, as

$$P_{\mathbf{x}}'(A) = \int_A f_{\mathbf{x}}(\xi) d\xi = P_{\mathbf{x}}(A) \quad (3-27)$$

Equation (3-24) relates that $P_{\mathbf{x}}(A \subset R^n)$ defined in (3-26) is equal to $P(\{\omega: \mathbf{x}(\omega) \in A\} \subset \Omega)$ for all sets A of interest. We now have a *new probability space*, $(R^n, \mathcal{F}_B, P_{\mathbf{x}})$, generated by the mapping \mathbf{x} from the original probability space:

$$(\Omega, \mathcal{F}, P) \xrightarrow{\mathbf{x}(\cdot)} (R^n, \mathcal{F}_B, P_{\mathbf{x}}) \quad (3-28)$$

Quite often, one can describe a problem conveniently in terms of the probability space $(R^n, \mathcal{F}_B, P_{\mathbf{x}})$, neglecting the original probability space. For our applications, $\mathbf{x}(\cdot)$ is the identity mapping, so the issue is rather academic. However, there are many problem areas in which recollection of the fundamental probability space and associated sets of interest yields clear insights into subtle and troublesome considerations.

3.4 CONDITIONAL PROBABILITY AND DENSITIES

Suppose we have two random variables \mathbf{x} and \mathbf{y} mapping from a sample space Ω into R^n and R^m , respectively. Further suppose that \mathbf{x} can assume only discrete values \mathbf{x}_i and similarly \mathbf{y} can take only discrete values \mathbf{y}_j , with i and j integers (a finite or countably infinite number). If we knew that \mathbf{y} has assumed a particular realization \mathbf{y}_j , that knowledge would, in general, affect the determination of the probability that $\mathbf{x}(\omega) = \mathbf{x}_i$ for a given value \mathbf{x}_i . Figure 3.11

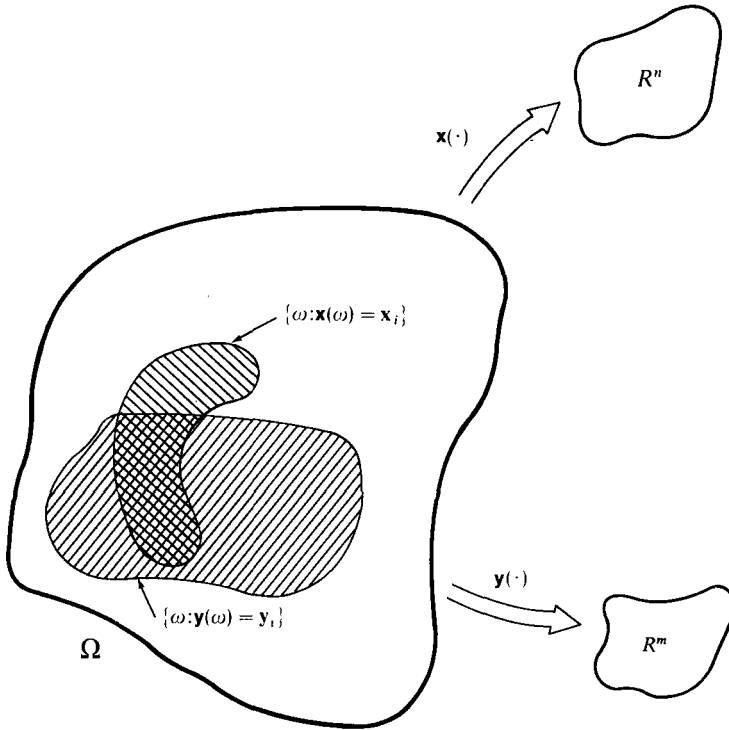


FIG. 3.11 Conditional probability via sets in Ω . The crosshatching is $\{\omega: \mathbf{x}(\omega) = \mathbf{x}_i \text{ and } \mathbf{y}(\omega) = \mathbf{y}_j\}$.

depicts the various sets of interest in the sample space Ω , namely $\{\omega: \mathbf{x}(\omega) = \mathbf{x}_i\}$ and $\{\omega: \mathbf{y}(\omega) = \mathbf{y}_j\}$. If we did not know the value \mathbf{y} assumes, then we would simply evaluate $P(\{\omega: \mathbf{x}(\omega) = \mathbf{x}_i\})$ using the probability measure defined for sets A , $A \subset \Omega$ and $A \in \mathcal{F}$. However, if we *know* that $\mathbf{y}(\omega) = \mathbf{y}_j$, we can restrict our attention to $\{\omega: \mathbf{y}(\omega) = \mathbf{y}_j\} \subset \Omega$ instead of considering all Ω . *Within that set* of ω , we want to know the probability of the set of ω such that $\mathbf{x}(\omega) = \mathbf{x}_i$ as well: the probability of $\{\omega: \mathbf{x}(\omega) = \mathbf{x}_i \text{ and } \mathbf{y}(\omega) = \mathbf{y}_j\}$, the cross hatched set in Fig. 3.11, relative to the set $\{\omega: \mathbf{y}(\omega) = \mathbf{y}_j\}$. Thus, the conditional probability that $\mathbf{x}(\omega) = \mathbf{x}_i$, conditioned on the fact that $\mathbf{y}(\omega) = \mathbf{y}_j$, can be defined as

$$P(\mathbf{x}(\omega) = \mathbf{x}_i | \mathbf{y}(\omega) = \mathbf{y}_j) = \frac{P(\mathbf{x}(\omega) = \mathbf{x}_i \text{ and } \mathbf{y}(\omega) = \mathbf{y}_j)}{P(\mathbf{y}(\omega) = \mathbf{y}_j)} \quad (3-29a)$$

$$= \frac{P(\{\omega: \mathbf{x}(\omega) = \mathbf{x}_i \text{ and } \mathbf{y}(\omega) = \mathbf{y}_j\})}{P(\{\omega: \mathbf{y}(\omega) = \mathbf{y}_j\})} \quad (3-29b)$$

Note that this conditional probability need only be defined for sets $A \subset \{\omega: \mathbf{y}(\omega) = \mathbf{y}_j\} \subset \Omega$, $A \in \mathcal{F}' \subset \mathcal{F}$ (\mathcal{F}' can be a "coarser" σ -algebra than \mathcal{F} , consisting of fewer elements).

To evaluate such a conditional probability numerically, first one could calculate the numerator in (3-29) as the ratio of the number of trials in which both events occur to the total number of trials. Then the denominator could be generated as the ratio of trials in which $\mathbf{y}(\omega) = \mathbf{y}_j$ to the total number of trials. Note also that if we summed over all \mathbf{x}_i 's, we would obtain a probability of one for any given \mathbf{y}_j .

However, this definition of conditional probability is valid only if $P(\mathbf{y}(\omega) = \mathbf{y}_j) > 0$. In our applications, we will be considering continuous random variables \mathbf{y} , for which $P(\mathbf{y}(\omega) = \mathbf{y}_j) = 0$, so the previous definition breaks down. Measure theory can be used to develop the concept of conditional probabilities rigorously and in more generality than we will require (we will assume the existence of appropriate densities). More will be said about this measure theoretic approach in Section 3.7, once the idea of expectation has been introduced. For now, we will develop the concept of a conditional density function, which will be of basic importance for estimation problems addressed in the sequel.

Let us first provide an interpretation of $f_{\mathbf{x}|\mathbf{y}}(\xi|\mathbf{y}_0)$, the conditional density of \mathbf{x} as a function of ξ , conditioned on knowledge that the random variable \mathbf{y} has assumed the realization $\mathbf{y}_0: \mathbf{y}(\omega) = \mathbf{y}_0$. Let \mathbf{x} map Ω into R^n and \mathbf{y} map Ω into R^m , and let $A \subset R^n$ and $B \subset R^m$ be point sets of interest in the corresponding spaces, as in Fig. 3.12a. The conditional probability that $\mathbf{x}(\omega)$ lies in A , conditioned on the fact that $\mathbf{y}(\omega) \in B$, is

$$P(\mathbf{x}(\omega) \in A | \mathbf{y}(\omega) \in B) = \frac{P(\mathbf{x}(\omega) \in A \text{ and } \mathbf{y}(\omega) \in B)}{P(\mathbf{y}(\omega) \in B)} \quad (3-30)$$

provided that $P(\mathbf{y}(\omega) \in B)$ is nonzero. The probabilities on the right hand side of (3-30) can be evaluated using the probability function P associated with (Ω, \mathcal{F}, P) , or with the functions P'_{xy} and P'_y associated with $(R^{nm}, \mathcal{F}_B, P'_{xy})$ and $(R^m, \mathcal{F}_B, P'_y)$ as described in (3-27). Thus, letting $f_{\mathbf{x},\mathbf{y}}(\cdot, \cdot)$ denote the joint probability density of \mathbf{x} and \mathbf{y} ,

$$P(\mathbf{x}(\omega) \in A | \mathbf{y}(\omega) \in B) = \frac{\int_A \left[\int_B f_{\mathbf{x},\mathbf{y}}(\xi, \gamma) d\gamma \right] d\xi}{\int_B f_{\mathbf{y}}(\rho) d\rho} \quad (3-31a)$$

$$= \int_A \frac{\int_B f_{\mathbf{x},\mathbf{y}}(\xi, \gamma) d\gamma}{\int_B f_{\mathbf{y}}(\rho) d\rho} d\xi \quad (3-31b)$$

From (3-31b) it can be seen that the conditional density for \mathbf{x} , conditioned on the fact that $\mathbf{y}(\omega) \in B$, would be

$$f_{\mathbf{x}}(\xi | \mathbf{y}(\omega) \in B) = \frac{\int_B f_{\mathbf{x},\mathbf{y}}(\xi, \gamma) d\gamma}{\int_B f_{\mathbf{y}}(\rho) d\rho} \quad (3-32)$$

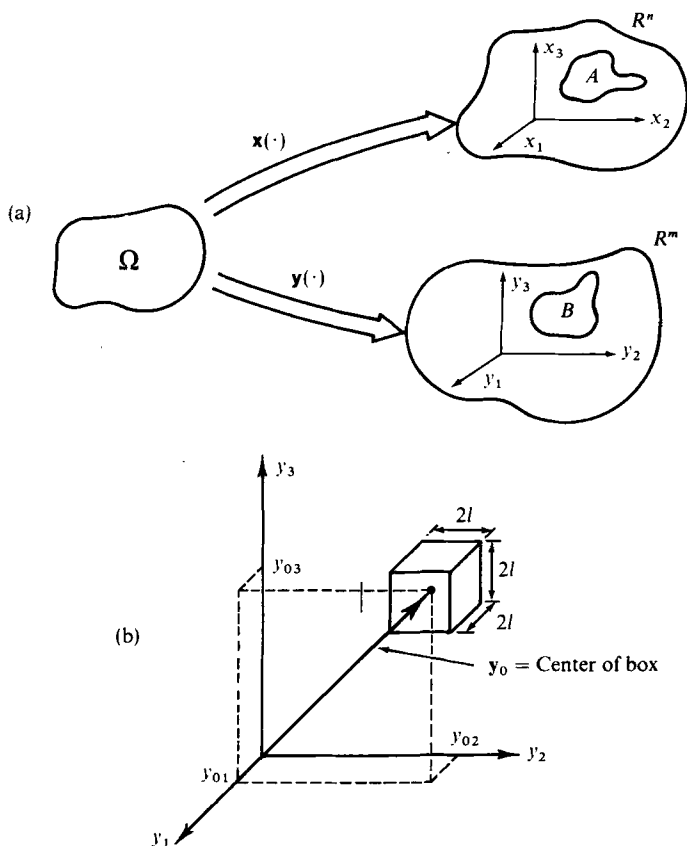


FIG. 3.12 (a) Sets of interest in R^n and R^m . (b) A particular set $B \subset R^m$. \mathbf{y}_0 is center of box.

Now consider a particular set $B \subset R^m$, namely a hypercube centered at \mathbf{y}_0 of dimension $2l$ on each side, as shown in Fig. 3.12b:

$$B = \{\mathbf{y} \in R^m : |y_1 - y_{01}| \leq l, |y_2 - y_{02}| \leq l, \dots, |y_m - y_{0m}| \leq l\} \quad (3-33)$$

We can write the density functions $f_{\mathbf{x},\mathbf{y}}(\xi, \gamma)$ and $f_{\mathbf{y}}(\rho)$ in terms of their evaluations at the given value of \mathbf{y}_0 as

$$f_{\mathbf{x},\mathbf{y}}(\xi, \gamma) = f_{\mathbf{x},\mathbf{y}}(\xi, \mathbf{y}_0) + \delta f_{\mathbf{x},\mathbf{y}}(\xi, \gamma - \mathbf{y}_0) \quad (3-34a)$$

$$f_{\mathbf{y}}(\rho) = f_{\mathbf{y}}(\mathbf{y}_0) + \delta f_{\mathbf{y}}(\rho - \mathbf{y}_0) \quad (3-34b)$$

Thus, (3-32) can be written as

$$f_{\mathbf{x}}(\xi | \mathbf{y}(\omega) \in B) = \frac{\int_B [f_{\mathbf{x},\mathbf{y}}(\xi, \mathbf{y}_0) + \delta f_{\mathbf{x},\mathbf{y}}(\xi, \gamma - \mathbf{y}_0)] d\gamma}{\int_B [f_{\mathbf{y}}(\mathbf{y}_0) + \delta f_{\mathbf{y}}(\rho - \mathbf{y}_0)] d\rho} \quad (3-35)$$

Now let V_B be the "volume" of the hypercube in m dimensions, $(2l)^m$, to write

$$\begin{aligned} f_{\mathbf{x}}(\xi | \mathbf{y}(\omega) \in B) &= \frac{V_B f_{\mathbf{x},\mathbf{y}}(\xi, \mathbf{y}_0) + \int_B \delta f_{\mathbf{x},\mathbf{y}}(\xi, \boldsymbol{\gamma} - \mathbf{y}_0) d\boldsymbol{\gamma}}{V_B f_{\mathbf{y}}(\mathbf{y}_0) + \int_B \delta f_{\mathbf{y}}(\boldsymbol{\rho} - \mathbf{y}_0) d\boldsymbol{\rho}} \\ &= \frac{f_{\mathbf{x},\mathbf{y}}(\xi, \mathbf{y}_0) + (1/V_B) \int_B \delta f_{\mathbf{x},\mathbf{y}}(\xi, \boldsymbol{\gamma} - \mathbf{y}_0) d\boldsymbol{\gamma}}{f_{\mathbf{y}}(\mathbf{y}_0) + (1/V_B) \int_B \delta f_{\mathbf{y}}(\boldsymbol{\rho} - \mathbf{y}_0) d\boldsymbol{\rho}} \end{aligned} \quad (3-36)$$

We assume that $f_{\mathbf{x},\mathbf{y}}$ and $f_{\mathbf{y}}$ are continuous, so that the mean value theorem can be used to write

$$\int_B \delta f_{\mathbf{x},\mathbf{y}} d\boldsymbol{\gamma} = V_B \delta f_{\mathbf{x},\mathbf{y}}(\xi, \mathbf{b}_1) \quad (3-37a)$$

$$\int_B \delta f_{\mathbf{y}} d\boldsymbol{\rho} = V_B \delta f_{\mathbf{y}}(\mathbf{b}_2) \quad (3-37b)$$

for \mathbf{b}_1 and \mathbf{b}_2 vectors somewhere in the hypercube B . Thus,

$$f_{\mathbf{x}}(\xi | \mathbf{y}(\omega) \in B) = \frac{f_{\mathbf{x},\mathbf{y}}(\xi, \mathbf{y}_0) + \delta f_{\mathbf{x},\mathbf{y}}(\xi, \mathbf{b}_1)}{f_{\mathbf{y}}(\mathbf{y}_0) + \delta f_{\mathbf{y}}(\mathbf{b}_2)} \quad (3-38)$$

Now consider reducing the hypercube down to the point \mathbf{y}_0 by letting $l \rightarrow 0$. This causes $\mathbf{b}_1 \rightarrow \mathbf{y}_0$, $\mathbf{b}_2 \rightarrow \mathbf{y}_0$, $\delta f_{\mathbf{x},\mathbf{y}}(\xi, \mathbf{b}_1) \rightarrow 0$, and $\delta f_{\mathbf{y}}(\mathbf{b}_2) \rightarrow 0$, so that

$$\lim_{l \rightarrow 0} f_{\mathbf{x}}(\xi | \mathbf{y}(\omega) \in B) = f_{\mathbf{x},\mathbf{y}}(\xi, \mathbf{y}_0) / f_{\mathbf{y}}(\mathbf{y}_0) \quad (3-39)$$

defines what is meant by $f_{\mathbf{x}|\mathbf{y}}(\xi | \mathbf{y}_0)$, sometimes denoted as $f_{\mathbf{x}}(\xi | \mathbf{y} = \mathbf{y}_0)$. This development is not the most general possible, but is sufficient for the continuous random variable problems to be considered later.

A fundamental result, which some use as the basic definition of a *conditional probability density*, is:

$$f_{\mathbf{x}|\mathbf{y}}(\xi | \boldsymbol{\rho}) = \frac{f_{\mathbf{x},\mathbf{y}}(\xi, \boldsymbol{\rho})}{f_{\mathbf{y}}(\boldsymbol{\rho})} \quad (3-40)$$

The denominator in (3-40) can be interpreted as a term to normalize the expression, so that the total "area under the density" is unity:

$$\begin{aligned} \int_{-\infty}^{\infty} f_{\mathbf{x}|\mathbf{y}}(\xi | \boldsymbol{\rho}) d\xi &= \int_{-\infty}^{\infty} \frac{f_{\mathbf{x},\mathbf{y}}(\xi, \boldsymbol{\rho})}{f_{\mathbf{y}}(\boldsymbol{\rho})} d\xi = \frac{\int_{-\infty}^{\infty} f_{\mathbf{x},\mathbf{y}}(\xi, \boldsymbol{\rho}) d\xi}{f_{\mathbf{y}}(\boldsymbol{\rho})} \\ &= \frac{f_{\mathbf{y}}(\boldsymbol{\rho})}{f_{\mathbf{y}}(\boldsymbol{\rho})} = 1 \end{aligned} \quad (3-41)$$

A graphical representation of (3-40) is useful for insight. Figure 3.13 portrays the joint density function $f_{\mathbf{x},\mathbf{y}}(\xi, \boldsymbol{\rho})$ as a surface above the ξ - $\boldsymbol{\rho}$ plane. To generate $f_{\mathbf{x}|\mathbf{y}}(\xi | \mathbf{y}_0)$ from this surface, a plane is passed through the surface at $\boldsymbol{\rho} = \mathbf{y}_0$, orthogonal to the $\boldsymbol{\rho}$ axis, resulting in the shaded region being $f_{\mathbf{x},\mathbf{y}}(\xi, \mathbf{y}_0)$. If that

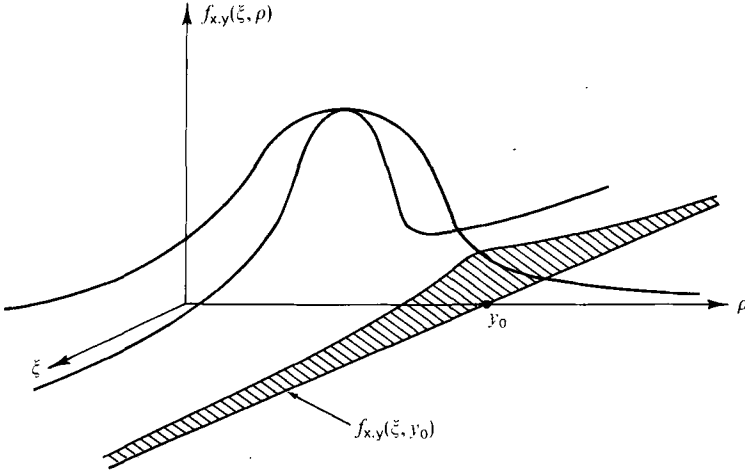


FIG. 3.13 Generation of conditional density.

function is divided by the number $f_y(y_0)$, the resulting function is normalized: its height is adjusted so that the shaded region is of area one. Note that $f_y(y_0)$ can be obtained from integrating $f_{x,y}(\xi, \rho)$ over all ξ , and evaluating the resulting function at $\rho = y_0$.

Eventually, we will want to consider the problem of estimating the value of a state vector \mathbf{x} based upon a set of measurements $\mathbf{z}_1 = \mathbf{z}_1, \mathbf{z}_2 = \mathbf{z}_2, \dots, \mathbf{z}_N = \mathbf{z}_N$. To accomplish this objective, we will propagate the conditional density $f_{\mathbf{x}|\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_N}(\xi|\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_N)$ of the state \mathbf{x} (modeled as a random variable), conditioned on knowledge of the entire set of measurements. This density embodies all of the information needed for estimation purposes.

Equation (3-40) is one form of *Bayes' rule*. Another useful form of this rule is

$$f_{\mathbf{x}|\mathbf{y}}(\xi|\rho) = \frac{f_{\mathbf{x},\mathbf{y}}(\xi, \rho)}{f_{\mathbf{y}}(\rho)} = \frac{f_{\mathbf{y}|\mathbf{x}}(\rho|\xi)f_{\mathbf{x}}(\xi)}{f_{\mathbf{y}}(\rho)} \quad (3-42)$$

Through this expression, one can readily generate the conditional density $f_{\mathbf{x}|\mathbf{y}}(\xi|\rho)$ if it is possible to write $f_{\mathbf{y}|\mathbf{x}}(\rho|\xi)$ and the unconditional densities for \mathbf{x} and \mathbf{y} . This will be exploited to derive $f_{\mathbf{x}|\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_N}(\xi|\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_N)$ for the estimation problem as just described. The denominator in Eq. (3-42) can be expanded to yield yet another useful form of Bayes' rule:

$$f_{\mathbf{x}|\mathbf{y}}(\xi|\rho) = \frac{f_{\mathbf{y}|\mathbf{x}}(\rho|\xi)f_{\mathbf{x}}(\xi)}{\int_{-\infty}^{\infty} f_{\mathbf{x},\mathbf{y}}(\xi, \rho) d\xi} = \frac{f_{\mathbf{y}|\mathbf{x}}(\rho|\xi)f_{\mathbf{x}}(\xi)}{\int_{-\infty}^{\infty} f_{\mathbf{y}|\mathbf{x}}(\rho|\xi)f_{\mathbf{x}}(\xi) d\xi} \quad (3-43)$$

Note that the integrand in the denominator is of the same form as the numerator.

EXAMPLE 3.7 Consider two scalar random variables \mathbf{x} and \mathbf{y} , described through the joint density function

$$f_{\mathbf{x},\mathbf{y}}(\xi, \rho) = (1/\pi) \exp\{-2[\xi - 1]^2 - 2[\rho - 2]^2 + 2\sqrt{3}[\xi - 1][\rho - 2]\}$$

The conditional probability density $f_{\mathbf{x}|\mathbf{y}}(\xi|\rho)$ can be generated from (3-40) once $f_{\mathbf{y}}(\rho)$ is established through the concept of marginal densities:

$$f_{\mathbf{y}}(\rho) = \int_{-\infty}^{\infty} f_{\mathbf{x},\mathbf{y}}(\xi, \rho) d\xi = \frac{1}{\sqrt{2\pi}} \exp\{-\frac{1}{2}[\rho - 2]^2\}$$

Then, (3-40) yields, after some algebraic reduction,

$$f_{\mathbf{x}|\mathbf{y}}(\xi|\rho) = (1/\pi) \exp\{-2[\xi - (\sqrt{3}/2)\rho + \sqrt{3} - 1]^2\}$$

This is the density function for \mathbf{x} as a function of ξ , given that \mathbf{y} has assumed some given realization ρ : given a particular ρ , the density is completely specified. ■

Through conditional probabilities and densities, we are specifying inter-relationships among random variables. The two extremes of such relationships are independence and functional dependence, to be discussed next.

Consider two random variables, \mathbf{x} mapping Ω into R^n and \mathbf{y} mapping Ω into R^m , and two admissible events $A \subset R^n$ and $B \subset R^m$. Then \mathbf{x} and \mathbf{y} are independent if

$$P(\{\omega: \mathbf{x}(\omega) \in A \text{ and } \mathbf{y}(\omega) \in B\}) = P(\{\omega: \mathbf{x}(\omega) \in A\})P(\{\omega: \mathbf{y}(\omega) \in B\}) \quad (3-44)$$

for all A and B . This is a fundamental definition, in terms of sets in the sample space Ω . To relate it to distribution functions, let A and B be chosen as sets of a particular form:

$$A = \{\mathbf{x}: \mathbf{x} = \mathbf{x}(\omega) \leq \xi\} = \{x_1 \leq \xi_1, x_2 \leq \xi_2, \dots, x_n \leq \xi_n\} \quad (3-45a)$$

$$B = \{\mathbf{y}: \mathbf{y} = \mathbf{y}(\omega) \leq \rho\} \quad (3-45b)$$

Then, by the definition of the appropriate distribution functions,

$$P(\{\omega: \mathbf{x}(\omega) \in A \text{ and } \mathbf{y}(\omega) \in B\}) = F_{\mathbf{x},\mathbf{y}}(\xi, \rho) \quad (3-46a)$$

$$P(\{\omega: \mathbf{x}(\omega) \in A\}) = F_{\mathbf{x}}(\xi) \quad (3-46b)$$

$$P(\{\omega: \mathbf{y}(\omega) \in B\}) = F_{\mathbf{y}}(\rho) \quad (3-46c)$$

for this particular choice of A and B . Thus, if \mathbf{x} and \mathbf{y} are independent, then

$$F_{\mathbf{x},\mathbf{y}}(\xi, \rho) = F_{\mathbf{x}}(\xi)F_{\mathbf{y}}(\rho) \quad (3-47)$$

for all ξ and ρ . If the distribution functions in (3-47) all have well-defined derivatives, one can conclude that, if \mathbf{x} and \mathbf{y} are independent, then

$$f_{\mathbf{x},\mathbf{y}}(\xi, \rho) = f_{\mathbf{x}}(\xi)f_{\mathbf{y}}(\rho) \quad (3-48)$$

for all ξ and ρ .

Another, more restrictive, approach is to define random vectors \mathbf{x} and \mathbf{y} to be independent if their joint density $f_{\mathbf{x},\mathbf{y}}(\xi, \rho)$ can be equated to the product of the separate marginal densities $f_{\mathbf{x}}(\xi)$ and $f_{\mathbf{y}}(\rho)$, as in (3-48). However, such a "definition" is valid only if the densities involved exist, whereas the concept of independence does not inherently require such existence.

When using the fundamental sample space and its subsets to think of independence, confusion sometimes arises between independent events and mutually exclusive events. If the occurrence of event A implies that B did not occur and vice versa, i.e., if $A \cap B = \emptyset$, then A and B are mutually exclusive. Said another way, if $P(A) = 1$ implies $P(B) = 0$ and $P(B) = 1$ implies $P(A) = 0$, then A and B are mutually exclusive. However, if knowledge of $P(A)$ gives you no information about $P(B)$ and vice versa, i.e., if $P(A \text{ and } B) = P(A)P(B)$, then A and B are independent events.

If \mathbf{x} and \mathbf{y} are independent, then (3-48) and Bayes' rule together yield $f_{\mathbf{x}|\mathbf{y}}(\xi|\rho)$ as

$$f_{\mathbf{x}|\mathbf{y}}(\xi|\rho) = f_{\mathbf{x},\mathbf{y}}(\xi, \rho)/f(\rho) = f_{\mathbf{x}}(\xi)f_{\mathbf{y}}(\rho)/f_{\mathbf{y}}(\rho) = f_{\mathbf{x}}(\xi) \quad (3-49)$$

i.e., the conditional density for \mathbf{x} , conditioned on knowledge that \mathbf{y} has assumed a realization ρ , is equal to the unconditional density for \mathbf{x} . This makes sense conceptually: if \mathbf{x} and \mathbf{y} are to be independent, then knowledge of the value of $\mathbf{y}(\omega)$ should give no information about the value of $\mathbf{x}(\omega)$.

In practice, physical arguments are often presented to establish the independence of two random variables. The validity of such arguments can be established through empirical testing as well. For example, if \mathbf{x} describes the outcome of one toss of a coin, and \mathbf{y} models the outcome of another toss of a coin, intuition dictates that there is no causal relationship between \mathbf{x} and \mathbf{y} —that the outcome of one toss does not affect the other toss. Experimental testing can substantiate (or contradict) such intuition. In other cases, uncertainty in the value of two quantities can be ascribed to physically unrelated sources. For instance, consider a sampled signal from an aircraft inertial system, modeled as a true position indication corrupted by a noise random variable \mathbf{n}_{INS} , and a sample of ground-based radar data, modeled similarly as a true position indication corrupted by noise $\mathbf{n}_{\text{radar}}$. The sources of uncertainty modeled through \mathbf{n}_{INS} (accelerometer bias, gyro drift, aircraft bending, etc.) are physically unrelated to the effects modeled by $\mathbf{n}_{\text{radar}}$ (electronic noise, atmospheric effects, etc.), and so these random variables are assumed to be independent.

The other extreme case of random variable interrelationship is *functional dependence*. If \mathbf{x} is a deterministic function of \mathbf{y} , $\mathbf{x} = \phi(\mathbf{y})$, then the conditional probability density function $f_{\mathbf{x}|\mathbf{y}}(\xi|\rho)$ is an impulse:

$$f_{\mathbf{x}|\mathbf{y}}(\xi|\rho) = \delta[\xi - \phi(\rho)] = \delta[\phi(\rho) - \xi] \quad (3-50)$$

where the delta function $\delta(\cdot)$ is defined as the function which satisfies the conditions:

$$\int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} \delta(\xi) d\xi_1 \cdots d\xi_n = 1; \quad \delta(\xi) = 0 \quad \text{for all } \xi \neq \mathbf{0} \quad (3-51)$$

It assumes a value of zero everywhere in R^n except where its argument is $\mathbf{0}$, and its integrated value over all R^n is unity. Equation (3-50) asserts that the conditional density function collapses down to an impulse function along the graph $\xi = \phi(\rho)$. If $\mathbf{x} = \phi(\mathbf{y})$ and $\mathbf{y}(\omega) = \rho$, then $\mathbf{x}(\omega) = \phi[\mathbf{y}(\omega)] = \phi(\rho)$ with no uncertainty: all of the probability is concentrated at $\xi = \phi(\rho)$, rather than being spread over a range of ξ values. In general, we will want to avoid impulse density functions, employing discontinuous distribution functions in such cases instead. However, the geometric insight of a density function collapsing down along certain loci of ξ values will be of practical use in certain applications, such as a system model outputting a number of "perfect" measurements.

3.5 FUNCTIONS OF RANDOM VARIABLES

The preceding section introduced the concept of functions of random variables, and this warrants some further attention. Let \mathbf{x} be a vector random variable that maps the sample space Ω into n -dimensional Euclidean space R^n . Now consider a continuous mapping $\theta(\cdot)$ from R^n into R^m , thus generating a vector $\mathbf{y} \in R^m$ from a vector $\mathbf{x} \in R^n$, as depicted in Fig. 3.14. Actually, $\theta(\cdot)$ can be out of a larger class of functions than the continuous functions, called Baire functions (Borel measurable functions), composed of continuous functions and limits of continuous functions, but this generality will not be needed in our applications.

Now define the m -vector-valued function \mathbf{y} as the composite mapping $\theta[\mathbf{x}(\cdot)]$. Then \mathbf{y} is itself a random variable, denoted as \mathbf{y} :

$$\mathbf{y}(\cdot) = \theta[\mathbf{x}(\cdot)] \quad (3-52a)$$

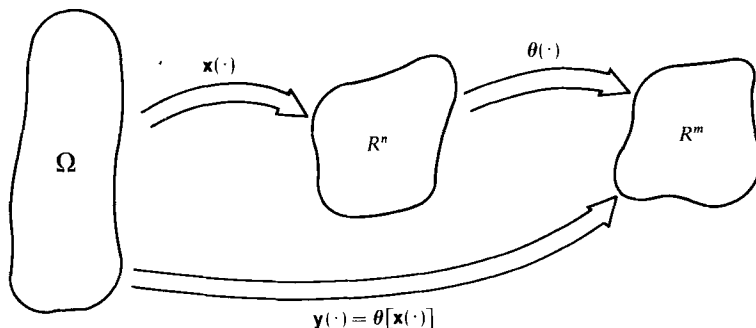


FIG. 3.14 Function of a random variable.

or

$$\begin{aligned} y_1(\cdot) &= \theta_1[x_1(\cdot), x_2(\cdot), \dots, x_n(\cdot)] \\ &\vdots \\ y_m(\cdot) &= \theta_m[x_1(\cdot), x_2(\cdot), \dots, x_n(\cdot)] \end{aligned} \quad (3-52b)$$

Stated simply, every Baire function of a random variable is a random variable.

Recall Eq. (3-28): $\mathbf{x}(\cdot)$ generates a new probability space $(R^n, \mathcal{F}_B, P_x)$ from the original probability space (Ω, \mathcal{F}, P) , with P_x defined in (3-26) or (3-27). If $\theta(\cdot)$ is a Baire function (Borel measurable) on R^n , then for every set of interest B in the range space R^m , the inverse image in R^n , $\{x \in R^n: \theta(x) \in B\}$, is an event for which probability has been defined through P_x . If we were to view $(R^n, \mathcal{F}_B, P_x)$ as the underlying probability space, then this just defines $\theta(\cdot)$ itself as a random variable mapping from the sample space R^n into the space R^m .

Analogous to the discussion concerning (3-28), we would then expect to generate a new probability space, $(R^m, \mathcal{F}_B, P_y)$. The sets of interest in R^m will be the elements of the Borel field \mathcal{F}_B associated with R^m , and for all sets $B \subset R^m$ and $B \in \mathcal{F}_B$, we can define probabilities through an appropriate probability function (measure), $P_y(\cdot)$, to be described shortly. Thus,

$$(\Omega, \mathcal{F}, P) \xrightarrow{\mathbf{x}(\cdot)} (R^n, \mathcal{F}_B, P_x) \xrightarrow{\theta(\cdot)} (R^m, \mathcal{F}_B, P_y) \quad (3-53)$$

Then $\mathbf{y}(\cdot) = \theta[\mathbf{x}(\cdot)]$ is a random variable that directly maps into this new probability space:

$$(\Omega, \mathcal{F}, P) \xrightarrow{\mathbf{y}(\cdot) = \theta[\mathbf{x}(\cdot)]} (R^m, \mathcal{F}_B, P_y) \quad (3-54)$$

The random variable \mathbf{y} has a *distribution induced by the distribution of \mathbf{x}* :

$$\begin{aligned} F_y(\rho) &= P(\{\omega: \mathbf{y}(\omega) \leq \rho\}) = P(\{\omega: \theta[\mathbf{x}(\omega)] \leq \rho\}) \\ &= P_x(\{\mathbf{x}: \theta(\mathbf{x}) \leq \rho\}) \end{aligned} \quad (3-55)$$

The *induced probability density function*, if it exists, is given by

$$f_y(\rho) = \frac{\hat{\partial}^m}{\hat{\partial} \rho_1 \hat{\partial} \rho_2 \cdots \hat{\partial} \rho_m} F_y(\rho) \quad (3-56)$$

EXAMPLE 3.8 Consider the die-toss experiment discussed previously in Examples 3.3 and 3.5, in which we were interested only in the sets $A_1 = \{1 \text{ or } 2\}$ and $A_2 = \{3\}$. We defined a random variable $\mathbf{x}(\cdot)$ through

$$\mathbf{x}(\omega) = \begin{cases} 0 & \text{if } \omega \notin A_1 \text{ or } A_2 \\ 1 & \text{if } \omega \in A_1 \\ 2 & \text{if } \omega \in A_2 \end{cases}$$

and derived its probability distribution function.

Now let us say that you will receive a payoff according to what you roll on the die. Let the payoff function $\theta(\cdot)$ be defined as

$$\theta(x) = x^2 + 1$$

so that $\theta(x)$ is a random variable y defined as

$$y = \theta(x) = x^2 + 1$$

Now we want to establish the induced distribution of y , $F_y(\rho)$, to describe our potential payoff. This can be generated using

$$F_y(\rho) = P(\{\omega: y(\omega) \leq \rho\}) = P(\{\omega: \theta[x(\omega)] \leq \rho\}) = P(\{\omega: x^2(\omega) + 1 \leq \rho\})$$

For $\rho < 1$, $P(\{\omega: x^2(\omega) + 1 \leq \rho\}) = P(\emptyset) = 0$.

For $1 \leq \rho < 2$, $P(\{\omega: x^2(\omega) + 1 \leq \rho\}) = P(\{A_1 \cup A_2\}^*) = \frac{1}{2}$.

For $2 \leq \rho < 5$, $P(\{\omega: x^2(\omega) + 1 \leq \rho\}) = P(A_2^*) = \frac{5}{6}$.

For $5 \leq \rho < \infty$, $P(\{\omega: x^2(\omega) + 1 \leq \rho\}) = P(\Omega) = 1$.

Thus we obtain the distribution function for y induced by the distribution of x as plotted in Fig. 3.15. ■

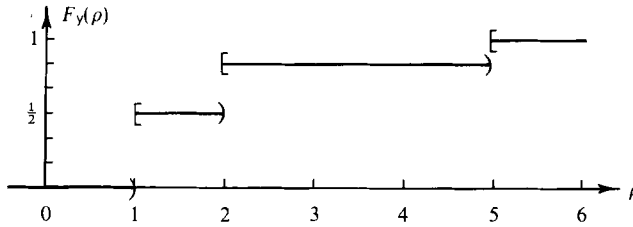


FIG. 3.15 Induced probability distribution function.

EXAMPLE 3.9 The scalar Gaussian random variable x is defined on $\Omega = R^1$ and is described through the density function

$$f_x(\xi) = (1/\sqrt{2\pi P}) \exp\{-(1/2P)(\xi - m)^2\}$$

Let the random variable y be defined through

$$y = \theta(x) = x^3$$

Now we want to generate the density function to describe y .

First, the distribution function is

$$\begin{aligned} F_y(\rho) &= P(\{\omega: y(\omega) \leq \rho\}) = P(\{\omega: x^3(\omega) \leq \rho\}) \\ &= P_x(\{x: x^3 \leq \rho\}) = P_x(\{x: x \leq \rho^{1/3}\}) = F_x(\rho^{1/3}) = \int_{-\infty}^{\rho^{1/3}} f_x(\xi) d\xi \end{aligned}$$

Then the desired density function is the derivative of $F_y(\rho)$ with respect to ρ . Using Leibnitz' rule, this yields

$$\begin{aligned} f_y(\rho) &= \frac{dF_y(\rho)}{d\rho} = \frac{d}{d\rho} \left\{ \int_{-\infty}^{\rho^{1/3}} f_x(\xi) d\xi \right\} \\ &= \frac{d}{d\rho} \{\rho^{1/3}\} \cdot f_x(\rho^{1/3}) = \frac{1}{3} \rho^{-2/3} \cdot \frac{1}{\sqrt{2\pi P}} \exp\left\{-\frac{1}{2P}(\rho^{1/3} - m)^2\right\} \quad \blacksquare \end{aligned}$$

Thus, given a random variable x and its distribution (or density if it exists) and the functional relationship $y = \theta(x)$, the induced distribution (density) for y can be determined. If densities do exist, then a useful result can be sum-

marized in the following manner. Let \mathbf{x} and \mathbf{y} be n -dimensional vector random variables, with $\mathbf{y} = \boldsymbol{\theta}(\mathbf{x})$. Suppose $\boldsymbol{\theta}^{-1}$ exists and both $\boldsymbol{\theta}$ and $\boldsymbol{\theta}^{-1}$ are continuously differentiable. Then

$$f_{\mathbf{y}}(\boldsymbol{\rho}) = f_{\mathbf{x}}[\boldsymbol{\theta}^{-1}(\boldsymbol{\rho})] \|\partial \boldsymbol{\theta}^{-1}(\boldsymbol{\rho}) / \partial \boldsymbol{\rho}\| \quad (3-57)$$

where $\|\partial \boldsymbol{\theta}^{-1}(\boldsymbol{\rho}) / \partial \boldsymbol{\rho}\| > 0$ is the absolute value of the Jacobian determinant arising naturally from integrations with change of variables. A proof of this theorem is outlined in Problem 3.8 at the end of this chapter.

EXAMPLE 3.10 Let us apply (3-57) to Example 3.9. Since $\mathbf{y} = \boldsymbol{\theta}(\mathbf{x}) = x^3$, the inverse function $\boldsymbol{\theta}^{-1}$ exists and can be expressed as (the real root)

$$\boldsymbol{\theta}^{-1}(\rho) = \rho^{1/3}$$

and thus its derivative is

$$\partial \boldsymbol{\theta}^{-1}(\rho) / \partial \rho = \frac{1}{3} \rho^{-2/3}$$

From (3-57),

$$\begin{aligned} f_{\mathbf{y}}(\rho) &= f_{\mathbf{x}}(\rho^{1/3}) \cdot \frac{1}{3} \rho^{-2/3} \\ &= \frac{1}{3} \rho^{-2/3} (1/\sqrt{2\pi P}) \exp\{-(1/2P)(\rho^{1/3} - m)^2\} \end{aligned}$$

as found previously. ■

Now consider the set function $P_{\mathbf{y}}$ introduced earlier. As in (3-26), for all sets $B \subset R^m$ and $B \in \mathcal{F}_{\mathbf{B}}$, we can define probabilities as

$$P_{\mathbf{y}}(B) = \int_B dF_{\mathbf{y}}(\boldsymbol{\rho}) \quad (3-58)$$

where meaning is given to the right hand side through measure theory. If density functions exist,

$$P_{\mathbf{y}}'(B) = \int_B f_{\mathbf{y}}(\boldsymbol{\rho}) d\boldsymbol{\rho} = P_{\mathbf{y}}(B) \quad (3-59)$$

The idea of induced densities and distributions is then embodied in the following result. If \mathbf{x} is a random variable mapping Ω into R^n , and $\mathbf{y} = \boldsymbol{\theta}(\mathbf{x})$, where $\boldsymbol{\theta}$ maps R^n into R^m , then

$$P_{\mathbf{y}}(B) = P_{\mathbf{x}}(\{\mathbf{x}: \boldsymbol{\theta}(\mathbf{x}) \in B\}) = P(\{\omega: \boldsymbol{\theta}[\mathbf{x}(\omega)] \in B\}) \quad (3-60)$$

The discussion in this section can help prevent the confusion that often arises in estimation. Consider having measurements available with which you want to estimate some quantities of interest, denoted as the n -dimensional vector $\boldsymbol{\theta}$. Suppose that the measurements are modeled through an m -dimensional vector of random variables $\mathbf{z}(\cdot)$, so that the numbers coming from the measuring devices are the realizations of the random variables for a particular outcome $\omega: \mathbf{z}(\omega) = \mathbf{z}$. Now you want to generate a mapping $\hat{\boldsymbol{\theta}}(\cdot)$ from R^m into R^n , called an estimator, that will map the given realizations into a “best” estimate of the value of $\boldsymbol{\theta}$. The composite mapping $\hat{\boldsymbol{\theta}}[\mathbf{z}(\cdot)]$ can then be considered a random

variable, often denoted as $\hat{\theta}$, called a randomized estimate or estimator, or just an estimate. Finally, estimation algorithms generate vectors of numbers, $\hat{\theta}(\mathbf{z}) \in R^n$, which are also called estimates. Whether one is concerned with a functional mapping, a random variable, or a vector of numbers is a fundamental distinction to be made, but one that can be misinterpreted unless one takes care to be aware of this aspect. The notation adopted herein specifically attempts to clarify this issue.

3.6 EXPECTATION AND MOMENTS OF RANDOM VARIABLES

The distribution or density function for a random variable is the entity of fundamental interest in Bayesian estimation, embodying all information known about the variable. Once it is generated, an "optimal" estimate can be defined using some chosen criterion. Similarly, it can be used to compute the expected value of some function of the random variable, where this "expected value" is just the average value one would obtain over the ensemble of outcomes of an "experiment." The expected value of particular functions will generate moments of a random variable, which are parameters (statistics) that characterize the distribution or density function. Although one would like to portray these functions completely through estimation, it is generally more feasible to evaluate expressions for a finite number of moments instead, thereby generating a partial description of the functions. In the case of Gaussian random variables, it will turn out that specification of only the first two moments will *completely* describe the distribution or density function.

Let \mathbf{x} be an n -dimensional random variable vector described through a density function $f_{\mathbf{x}}(\xi)$, and let \mathbf{y} be an m -dimensional vector function of \mathbf{x} :

$$\mathbf{y}(\cdot) = \theta[\mathbf{x}(\cdot)] \quad (3-61)$$

where $\theta(\cdot)$ is continuous. Thus \mathbf{y} is also a random variable with induced density $f_{\mathbf{y}}(\rho)$. Then the *expectation* of \mathbf{y} is

$$E[\mathbf{y}] = \int_{-\infty}^{\infty} \theta(\xi) f_{\mathbf{x}}(\xi) d\xi = \int_{-\infty}^{\infty} \rho f_{\mathbf{y}}(\rho) d\rho \quad (3-62)$$

If the density function $f_{\mathbf{x}}(\xi)$ does not exist, then $E[\mathbf{y}]$ can still be defined as

$$E[\mathbf{y}] = \int_{\Omega} \theta[\mathbf{x}(\omega)] dP(\omega) = \int_{R^n} \theta(\xi) dF_{\mathbf{x}}(\xi) = \int_{R^m} \rho dF_{\mathbf{y}}(\rho) \quad (3-63)$$

using measure theory to give meaning to the indicated integrals [12, 13, 15]. Note that the succeeding integrals in (3-63) are carried out over Ω , R^n , and R^m , respectively, using the appropriate probability functions. The integration performed over the original sample space Ω naturally provides the most basic

definition of expectation. However, for our applications, we will assume $f_{\mathbf{x}}(\xi)$ exists, so (3-62) will suffice as a definition.

EXAMPLE 3.11 Let us calculate the expected payoff for the die-toss experiment described in Example 3.8. A distribution function as in Fig. 3.15 indicates that the random variable assumes only discrete point values, with probability equal to the magnitude of each associated discontinuity [see Eq. (3-14)]. In such a case, the integrals indicated in (3-63) become summations, as

$$\begin{aligned} E[y] &= \int \rho dF_y(\rho) = \sum_i \rho_i \Delta F_y(\rho_i) = 1 \cdot (1 - P_1 - P_2) + 2 \cdot P_1 + 5 \cdot P_2 \\ &= 1 \cdot \frac{1}{2} + 2 \cdot \frac{1}{3} + 5 \cdot \frac{1}{6} = 2 \end{aligned}$$

A distribution function that varies continuously except for a finite number of jump discontinuities can be decomposed into the sum of a continuous function and a function composed only of the jump discontinuities as in Fig. 3.15. This allows expectations to be evaluated through addition of results obtained from the separate functions. ■

Since expectation is by definition an integration, it is a linear operation. In other words, for c equal to a scalar constant,

$$E[c\mathbf{y}] = cE[\mathbf{y}] \quad (3-64a)$$

$$E[\mathbf{y}_1 + \mathbf{y}_2] = E[\mathbf{y}_1] + E[\mathbf{y}_2] \quad (3-64b)$$

Combining (3-64a) and (3-64b) yields the useful result that, for \mathbf{A} a known matrix,

$$E[\mathbf{A}\mathbf{y}] = \mathbf{A}E[\mathbf{y}] \quad (3-65)$$

Now let us consider some specific functions $\theta(\cdot)$. First let $\theta(\mathbf{x}) = \mathbf{x}$ to generate the *first moment* of \mathbf{x} or the *mean* of \mathbf{x} . Define an n -dimensional vector \mathbf{m} , whose components are the mean values $m_i \triangleq E[x_i]$:

$$\mathbf{m} \triangleq \begin{bmatrix} m_1 \\ \vdots \\ m_n \end{bmatrix} \triangleq \begin{bmatrix} E[x_1] \\ \vdots \\ E[x_n] \end{bmatrix} = \begin{bmatrix} \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} \xi_1 f_{\mathbf{x}}(\xi) d\xi_1 \cdots d\xi_n \\ \vdots \\ \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} \xi_n f_{\mathbf{x}}(\xi) d\xi_1 \cdots d\xi_n \end{bmatrix} \quad (3-66)$$

This vector is the mean or first moment of \mathbf{x} , and notationally (3-66) can be written equivalently as

$$\mathbf{m} \triangleq E[\mathbf{x}] = \int_{-\infty}^{\infty} \xi f_{\mathbf{x}}(\xi) d\xi \quad (3-67)$$

Next consider Eq. (3-62), letting $\theta(\mathbf{x}) = \mathbf{x}\mathbf{x}^T$:

$$\mathbf{y} = \theta(\mathbf{x}) = \mathbf{x}\mathbf{x}^T = \begin{bmatrix} x_1^2 & x_1x_2 & \cdots & x_1x_n \\ \vdots & \vdots & & \vdots \\ x_nx_1 & x_nx_2 & \cdots & x_n^2 \end{bmatrix} \quad (3-68)$$

Define a matrix, denoted as Ψ , as the n -by- n matrix whose i - j component is the *correlation* of x_i and x_j (and thus the diagonal terms are autocorrelations, or mean squared values, the square roots of which are termed root mean squared, or *RMS*, values):

$$\Psi_{ij} \triangleq E[x_i x_j] = \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} \xi_i \xi_j f_{\mathbf{x}}(\xi) d\xi_1 \cdots d\xi_n \quad (3-69)$$

This matrix is the *second (noncentral) moment* of \mathbf{x} or the *autocorrelation matrix* of \mathbf{x} , and can be written as

$$\Psi \triangleq E[\mathbf{x}\mathbf{x}^T] = \int_{-\infty}^{\infty} \xi \xi^T f_{\mathbf{x}}(\xi) d\xi \quad (3-70)$$

where again this simply is a compact notation to be interpreted in the light of Eq. (3-69).

Let us consider yet another function, $\theta(\mathbf{x}) = [(\mathbf{x} - \mathbf{m})(\mathbf{x} - \mathbf{m})^T]$. This allows us to define an n -by- n matrix \mathbf{P} whose i - j component is the *covariance* of x_i and x_j :

$$\begin{aligned} P_{ij} &\triangleq E[(x_i - m_i)(x_j - m_j)] \\ &= \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} (\xi_i - m_i)(\xi_j - m_j) f_{\mathbf{x}}(\xi) d\xi_1 \cdots d\xi_n \end{aligned} \quad (3-71)$$

Note specifically that the mean values $m_i = E[x_i]$ and $m_j = E[x_j]$ in (3-70) are *not* random variables, but are statistics: deterministic numbers. The matrix \mathbf{P} is the *second central moment* of \mathbf{x} or the *covariance matrix* of \mathbf{x} , and can be written as

$$\mathbf{P} \triangleq E[(\mathbf{x} - \mathbf{m})(\mathbf{x} - \mathbf{m})^T] = \int_{-\infty}^{\infty} (\xi - \mathbf{m})(\xi - \mathbf{m})^T f_{\mathbf{x}}(\xi) d\xi \quad (3-72)$$

In cases where there might be ambiguity as to what correlation or covariance matrix is being discussed, subscripts will be employed in the notation, such as \mathbf{P}_{xx} or Ψ_{yy} .

Because the covariance will be significant in our work, it will be characterized further. The matrix \mathbf{P} is a symmetric, positive semidefinite matrix (its eigenvalues are nonnegative). The variances of the separate components of \mathbf{x} are along the diagonal:

$$P_{ii} \triangleq E[(x_i - m_i)^2] \quad (3-73)$$

The square root of a variance P_{ii} is termed the *standard deviation* of x_i , denoted as σ_i . Thus, the diagonal terms can be expressed as

$$P_{ii} \triangleq \sigma_i^2 \quad (3-74)$$

The *correlation coefficient* of x_i and x_j , denoted as r_{ij} , is defined as the ratio

$$r_{ij} \triangleq \frac{E[(x_i - m_i)(x_j - m_j)]}{(E[(x_i - m_i)^2])^{1/2}(E[(x_j - m_j)^2])^{1/2}} \triangleq \frac{P_{ij}}{\sigma_i \sigma_j} \quad (3-75)$$

Using (3-74) and (3-75), the covariance matrix \mathbf{P} can be written as

$$\mathbf{P} = \begin{bmatrix} \sigma_1^2 & r_{12}\sigma_1\sigma_2 & \cdots & r_{1n}\sigma_1\sigma_n \\ r_{12}\sigma_1\sigma_2 & \sigma_2^2 & \cdots & r_{2n}\sigma_2\sigma_n \\ \vdots & \vdots & \ddots & \vdots \\ r_{1n}\sigma_1\sigma_n & r_{2n}\sigma_2\sigma_n & \cdots & \sigma_n^2 \end{bmatrix} \quad (3-76)$$

If the correlation coefficient r_{ij} is zero, then the components x_i and x_j are said to be *uncorrelated*. Consequently, if \mathbf{P} is diagonal, i.e., if $r_{ij} = 0$ for all i and j with $i \neq j$, then \mathbf{x} is said to be composed of uncorrelated components.

Another expression for the covariance matrix can be derived in the following manner, using the facts that $E[\cdot]$ is linear and the mean vector \mathbf{m} is not random:

$$\begin{aligned} \mathbf{P} &= E[(\mathbf{x} - \mathbf{m})(\mathbf{x} - \mathbf{m})^T] = E[\mathbf{xx}^T - \mathbf{xm}^T - \mathbf{mx}^T + \mathbf{mm}^T] \\ &= E[\mathbf{xx}^T] - E[\mathbf{xm}^T] - E[\mathbf{mx}^T] + E[\mathbf{mm}^T] \\ &= E[\mathbf{xx}^T] - E[\mathbf{x}]\mathbf{m}^T - \mathbf{m}E[\mathbf{x}^T] + \mathbf{mm}^T \\ &= E[\mathbf{xx}^T] - \mathbf{mm}^T - \mathbf{mm}^T + \mathbf{mm}^T \\ \mathbf{P} &= E[\mathbf{xx}^T] - \mathbf{mm}^T \end{aligned} \quad (3-77)$$

This equation then directly relates the central and noncentral second moments. In the scalar case, it reduces to

$$P = E[x^2] - (E[x])^2 \quad (3-78)$$

The special cases of expectation that have been considered are just the first two moments of a random variable. (See Problems 3.24–3.27 for means of establishing best estimates of these moments from a finite set of empirical data.) Of course, there are higher ordered moments that can be used to characterize a probability density (or distribution) function. The mean relates where the density is centered, and the covariance gives an indication of the spread of the density about that mean value. In general, an endless number of moments would be required to specify a density function completely. In the particular case of a Gaussian random variable, the mean and covariance *completely* specify the density. By knowing the mean and covariance of a Gaussian random variable, you know *all* of the probability information contained in the associated density function, not just two parameters that partially describe its shape. This will be exploited to a great extent in linear estimation problems.

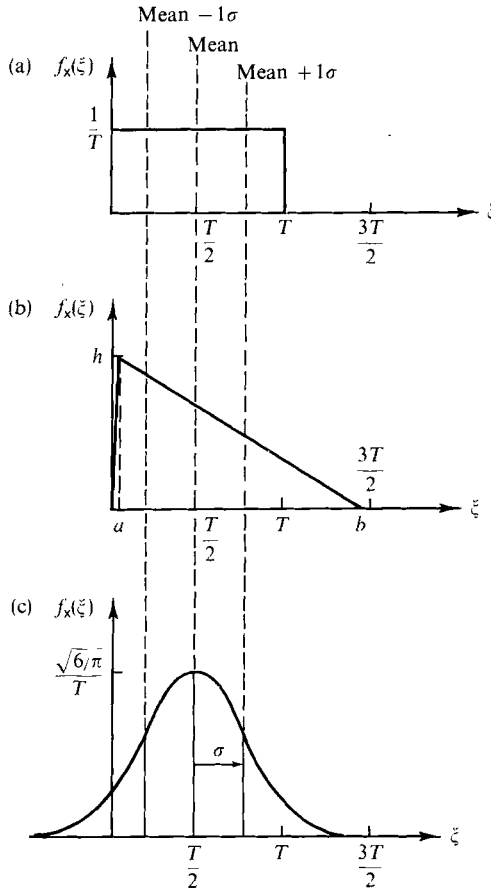


FIG. 3.16 Different random variables with equivalent first two moments. (a) Uniform. (b) Triangular. $a = (T/4)(3 - 5\sqrt{1/3}) \cong 0.03 T$, $b = (T/4)(3 + 5\sqrt{1/3}) \cong 1.47 T$, $h = (2/b) \cong 1.36 T$. (c) Gaussian. $\sigma = T/\sqrt{12} \cong 0.29 T$.

EXAMPLE 3.12 Consider the random variable x with uniform density between 0 and T , $f_x(\xi) = \{1/T \text{ for } \xi \in [0, T], 0 \text{ elsewhere}\}$, as depicted in Fig. 3.16a. The mean of x is

$$E[x] = \int_{-\infty}^{\infty} \xi f_x(\xi) d\xi = \int_0^T \xi \frac{1}{T} d\xi = \frac{T}{2}$$

and the variance of x is

$$E\left[\left(x - \frac{T}{2}\right)^2\right] = \int_0^T \left(\xi - \frac{T}{2}\right)^2 \frac{1}{T} d\xi = \frac{T^2}{12}$$

If y is defined as $\sin x$, then $E[y]$ is

$$E[y] = E[\sin x] = \int_0^T \sin \xi \frac{1}{T} d\xi = \frac{1}{T}(1 - \cos T)$$

Note that specification of just the first two moments of a random variable does not completely describe the associated distribution or density function. The triangular-shaped density function in Fig. 3.16b yields the same first two moments, despite its significantly different shape. Figure 3.16c depicts a Gaussian density yielding the same first two moments as in (a) and (b),

$$f_x(\xi) = \frac{1}{[2\pi(T^2/12)]^{1/2}} \exp\left\{-\frac{1}{2(T^2/12)}\left(\xi - \frac{T}{2}\right)^2\right\}$$

Furthermore, note that if one knew $f_x(\xi)$ were uniform or Gaussian, then knowledge of the first two moments would specify $f_x(\xi)$ completely. However, if $f_x(\xi)$ were triangular, these two parameters would not specify the density shape totally. An infinite number of moments is required to specify the shape of a general density function. ■

It will be useful to generalize the concept of the second moment of a single random variable \mathbf{x} to the second moment relationship between two random variables \mathbf{x} and \mathbf{y} . Such a concept is inherently involved in Eqs. (3-70) and (3-72), since the i - j component of Ψ or \mathbf{P} is the cross-correlation or covariance, respectively, of the scalar random variables x_i and x_j ; now we will generalize from the scalar to vector case. Let \mathbf{x} be an n -dimensional random vector and \mathbf{y} an m -dimensional random vector. Then the *cross-correlation matrix* of \mathbf{x} and \mathbf{y} is the n -by- m matrix whose i - j component is

$$E[x_i y_j] = \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} \xi_i \rho_j f_{\mathbf{x}, \mathbf{y}}(\xi, \rho) d\xi_1 \cdots d\xi_n d\rho_1 \cdots d\rho_m \quad (3-79)$$

This matrix is then expressed notationally as

$$\Psi_{xy} \triangleq E[\mathbf{xy}^T] = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \xi \rho^T f_{\mathbf{x}, \mathbf{y}}(\xi, \rho) d\xi d\rho \quad (3-80)$$

Similarly, the second central moment generalizes to the *cross-covariance matrix* \mathbf{x} and \mathbf{y} :

$$\mathbf{P}_{xy} \triangleq E[(\mathbf{x} - \mathbf{m}_x)(\mathbf{y} - \mathbf{m}_y)^T] = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (\xi - \mathbf{m}_x)(\rho - \mathbf{m}_y)^T f_{\mathbf{x}, \mathbf{y}}(\xi, \rho) d\xi d\rho \quad (3-81)$$

Two random vectors \mathbf{x} and \mathbf{y} are termed *uncorrelated* if their correlation matrix is equal to the outer product of their first order moments, i.e., if

$$E[\mathbf{xy}^T] = E[\mathbf{x}]E[\mathbf{y}^T] = \mathbf{m}_x \mathbf{m}_y^T \quad (3-82a)$$

or

$$E[x_i y_j] = E[x_i]E[y_j] \quad \text{for all } i \text{ and } j \quad (3-82b)$$

which is equivalent to the condition that $E\{[x_i - m_{x_i}][y_j - m_{y_j}]\} = 0$ for all i and j .

EXAMPLE 3.13 We want to show by a simple example that the preceding definition of uncorrelatedness corresponds to the previous definition of uncorrelated scalar random variables which involved the correlation coefficient described by Eq. (3-75). Consider two scalar random

variables, z_1 and z_2 . By (3-82), they are uncorrelated if $E[z_1 z_2] = E[z_1]E[z_2]$. Now let \mathbf{z} be the vector random variable made up of components z_1 and z_2 . The covariance of \mathbf{z} is then

$$\mathbf{P}_{zz} = \begin{bmatrix} E[z_1^2] - E[z_1]^2 & E[z_1 z_2] - E[z_1]E[z_2] \\ E[z_1 z_2] - E[z_1]E[z_2] & E[z_2^2] - E[z_2]^2 \end{bmatrix}$$

But, if z_1 and z_2 are uncorrelated, then $E[z_1 z_2] - E[z_1]E[z_2] = 0$, and the off-diagonal terms are zero. This is just the condition of the correlation coefficient of z_1 and z_2 being zero, as described earlier. ■

Whereas uncorrelatedness is a condition under which generalized second moments can be expressed as products of first order moments, independence is a condition under which the entire joint distribution or density function can be expressed as a product of marginal functions. As might be expected then, if \mathbf{x} and \mathbf{y} are independent, then they are uncorrelated, but not necessarily vice versa. This implication can be expressed simply as

$$\mathbf{x} \text{ and } \mathbf{y} \text{ independent} \rightarrow \mathbf{x} \text{ and } \mathbf{y} \text{ uncorrelated} \quad (3-83)$$

This can be demonstrated readily: by definition, we can write

$$E[\mathbf{xy}^T] = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \xi \rho^T f_{\mathbf{x},\mathbf{y}}(\xi, \rho) d\xi d\rho$$

If \mathbf{x} and \mathbf{y} are independent, then this becomes:

$$E[\mathbf{xy}^T] = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \xi \rho^T f_{\mathbf{x}}(\xi) f_{\mathbf{y}}(\rho) d\xi d\rho$$

Separating the integration yields the desired result:

$$E[\mathbf{xy}^T] = \int_{-\infty}^{\infty} \xi f_{\mathbf{x}}(\xi) d\xi \int_{-\infty}^{\infty} \rho^T f_{\mathbf{y}}(\rho) d\rho = E[\mathbf{x}]E[\mathbf{y}^T]$$

If \mathbf{x} and \mathbf{y} are uncorrelated, they are *not* necessarily independent. A counter-example to such an implication is given in the following example.

EXAMPLE 3.14 (modified from [3]) Let z be uniformly distributed between 0 and 1:

$$f_z(\zeta) = \begin{cases} 1 & \zeta \in [0, 1] \\ 0 & \text{otherwise} \end{cases}$$

Now define \mathbf{x} and \mathbf{y} as $\mathbf{x} = \sin(2\pi z)$ and $\mathbf{y} = \cos(2\pi z)$. It will now be shown that \mathbf{x} and \mathbf{y} are uncorrelated, but not independent. They are uncorrelated since

$$E[\mathbf{x}] = E[\mathbf{y}] = E[\mathbf{xy}] = E[\mathbf{x}]E[\mathbf{y}] = 0$$

However, consider higher order moments, as the fourth generalized moment:

$$E[\mathbf{x}^2 \mathbf{y}^2] = \frac{1}{8}$$

But,

$$E[\mathbf{x}^2] = E[\mathbf{y}^2] = \frac{1}{2}$$

so that

$$E[\mathbf{x}^2]E[\mathbf{y}^2] = \frac{1}{4}$$

which is not equal to $E[\mathbf{x}^2 \mathbf{y}^2]$. If \mathbf{x} and \mathbf{y} were independent, these would be equal. ■

Another related concept is that of orthogonality. Two random vectors \mathbf{x} and \mathbf{y} are termed *orthogonal* if their correlation matrix is the zero matrix: if $E[\mathbf{xy}^T] = \mathbf{0}$. Obviously, this concept is interrelated with \mathbf{x} and \mathbf{y} being uncorrelated, and this relation is as follows. If either \mathbf{x} or \mathbf{y} (or both) is zero-mean, then orthogonality and uncorrelatedness of \mathbf{x} and \mathbf{y} imply each other. However, if neither is zero-mean, then \mathbf{x} and \mathbf{y} may be uncorrelated or orthogonal or neither, but they cannot be both orthogonal and uncorrelated. Orthogonality provides one means of defining an optimal estimate: if we generate an estimate $\hat{\mathbf{x}}$ of \mathbf{x} based on measurement data \mathbf{z} , then that estimate can be termed optimal if the error $(\mathbf{x} - \hat{\mathbf{x}})$ is orthogonal to the data. This geometrical concept is instrumental in deriving optimal estimators by means of "orthogonal projections," the original means of derivation of the Kalman filter. We, however, will employ a Bayesian approach to estimation in the sequel.

3.7 CONDITIONAL EXPECTATIONS

The concept of expectation of some function of random variables answers the question, if we were to conduct a large (endless) number of experiments, what average value (over the entire ensemble of experimental outcomes, $\omega \in \Omega$) of that function would we achieve? Conditional expectations provide the same information, but incorporate insights into occurrence of events in Ω gained through observations of realizations of related random variables. In this section, we will first define conditional expectation under the assumption that the appropriate conditional density function exists. After investigating its properties and applications, the definition will be generalized to allow consideration of this concept without such an assumption.

Let \mathbf{x} and \mathbf{y} be random variables mapping Ω into R^n and R^m , respectively, and let \mathbf{z} be a continuous (Baire) function of \mathbf{x} ,

$$\mathbf{z}(\cdot) = \theta[\mathbf{x}(\cdot)] \quad (3-84)$$

so that \mathbf{z} is itself a random variable mapping Ω into R^r . Then the *conditional expected value*, or *conditional mean*, of \mathbf{z} , conditioned on the fact that \mathbf{y} has assumed the realization $\mathbf{y} \in R^m$, i.e., $\mathbf{y}(\omega) = \mathbf{y}$, is

$$E_{\mathbf{x}}[\mathbf{z}|\mathbf{y} = \mathbf{y}] = \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} \theta(\xi) f_{\mathbf{x}|\mathbf{y}}(\xi|\mathbf{y}) d\xi_1 \cdots d\xi_n \quad (3-85a)$$

$$= \int_{-\infty}^{\infty} \theta(\xi) f_{\mathbf{x}|\mathbf{y}}(\xi|\mathbf{y}) d\xi \quad (3-85b)$$

The subscript \mathbf{x} on $E_{\mathbf{x}}[\mathbf{z}|\mathbf{y} = \mathbf{y}]$ denotes that the expectation operation (integration) is performed over the possible values of \mathbf{x} , and sometimes this subscript is not included in the notation. For a given value $\mathbf{y} \in R^m$, $E_{\mathbf{x}}[\mathbf{z}|\mathbf{y} = \mathbf{y}]$ is a vector in R^r . Thus, $E_{\mathbf{x}}[\mathbf{z}|\mathbf{y} = \cdot]$ is a mapping from R^m into R^r , a function of the values $\mathbf{y} \in R^m$. Recall Section 3.5, Functions of Random Variables. If these \mathbf{y} values are realizations of the random variable \mathbf{y} , then the conditional

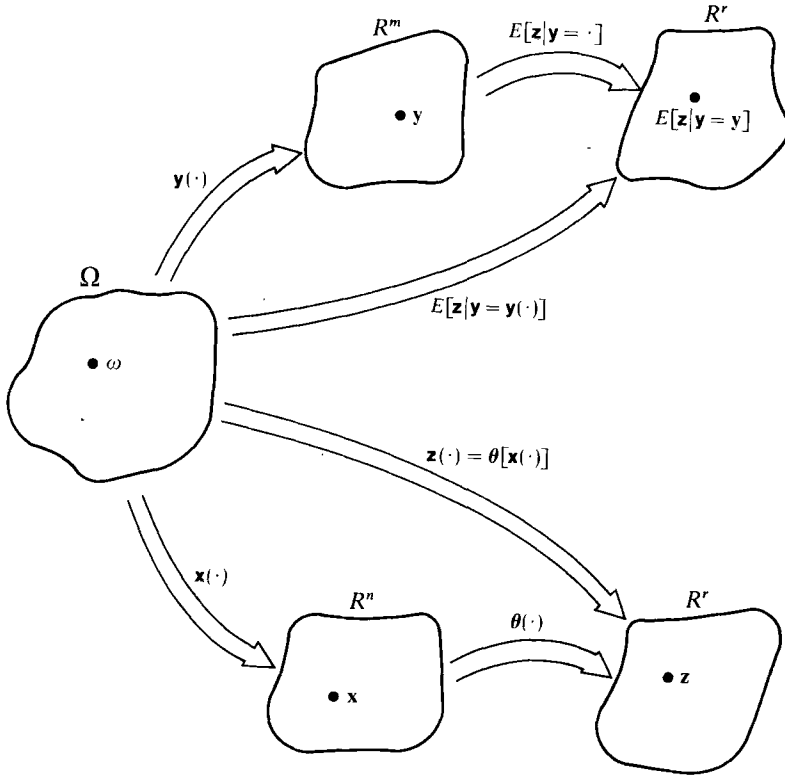


FIG. 3.17 Conditional expectation functional relationships.

expectation can be viewed as a random variable, i.e., the composite mapping $E_{\mathbf{x}}[\mathbf{z}|\mathbf{y} = \mathbf{y}(\cdot)]$ mapping Ω into R^r . These interrelationships are depicted in Fig. 3.17.

Moreover, the random variable $E_{\mathbf{x}}[\mathbf{z}|\mathbf{y} = \mathbf{y}(\cdot)]$ is unique and has the property that

$$E_{\mathbf{y}}\{E_{\mathbf{x}}[\mathbf{z}|\mathbf{y} = \mathbf{y}(\cdot)]\} = E_{\mathbf{x}}[\mathbf{z}] \quad (3-86)$$

Conceptually, this is reasonable. If we take the conditional expectation of \mathbf{z} , conditioned on a realized value of \mathbf{y} , and look at its expected value over all possible realizations of \mathbf{y} , then the result is the unconditional expectation of \mathbf{z} . Let us demonstrate the validity of (3-86) mathematically as well. By the definition of expectation, we can write

$$E_{\mathbf{x}}[\mathbf{z}] = \int_{-\infty}^{\infty} \theta(\xi) f_{\mathbf{x}}(\xi) d\xi$$

Now $f_{\mathbf{x}}(\xi)$ can be written as the marginal density derived from $f_{\mathbf{x},\mathbf{y}}(\xi, \rho)$ to yield

$$E_{\mathbf{x}}[\mathbf{z}] = \int_{-\infty}^{\infty} \theta(\xi) \left[\int_{-\infty}^{\infty} f_{\mathbf{x},\mathbf{y}}(\xi, \rho) d\rho \right] d\xi$$

Bayes' rule can be applied to derive an equivalent expression as

$$E_{\mathbf{x}}[\mathbf{z}] = \int_{-\infty}^{\infty} \theta(\xi) \left[\int_{-\infty}^{\infty} f_{\mathbf{x}|\mathbf{y}}(\xi|\rho) f_{\mathbf{y}}(\rho) d\rho \right] d\xi$$

Now we assume convergence in the definition of the integrals taken in different orders, so that we can interchange the order of integration (to be more precise, we invoke the Fubini theorem [7, 12, 13] from functional analysis) to yield

$$E_{\mathbf{x}}[\mathbf{z}] = \int_{-\infty}^{\infty} \left[\int_{-\infty}^{\infty} \theta(\xi) f_{\mathbf{x}|\mathbf{y}}(\xi|\rho) d\xi \right] f_{\mathbf{y}}(\rho) d\rho$$

The bracketed term is a function of ρ alone, where ρ is a dummy variable corresponding to realized values of \mathbf{y} (ρ is used as the dummy variable to distinguish it from a single realization \mathbf{y} of \mathbf{y}). Thus, the preceding expression is just the expected value (over all possible \mathbf{y} realizations, ρ) of the bracketed term, which is itself the conditioned expectation of \mathbf{z} : this directly yields

$$E_{\mathbf{x}}[\mathbf{z}] = E_{\mathbf{y}}\{E_{\mathbf{x}}[\mathbf{z}|\mathbf{y} = \mathbf{y}(\cdot)]\}$$

as desired.

The conditional expectation can also be viewed as a function $E_{\mathbf{x}}[\cdot|\mathbf{y} = \mathbf{y}]$ that maps a random variable \mathbf{z} into a vector $E_{\mathbf{x}}[\mathbf{z}|\mathbf{y} = \mathbf{y}] \in R^r$. As in the case of unconditional expectations, such an operation is defined through an integration and is *linear*. Thus, if \mathbf{A} is a known matrix,

$$E_{\mathbf{x}}[\mathbf{A}\mathbf{x}|\mathbf{y} = \mathbf{y}] = \mathbf{A}E_{\mathbf{x}}[\mathbf{x}|\mathbf{y} = \mathbf{y}] \quad (3-87)$$

$$E_{\mathbf{x}\mathbf{y}}[\mathbf{x} + \mathbf{y}|\mathbf{z} = \mathbf{z}] = E_{\mathbf{x}}[\mathbf{x}|\mathbf{z} = \mathbf{z}] + E_{\mathbf{y}}[\mathbf{y}|\mathbf{z} = \mathbf{z}] \quad (3-88)$$

Two special cases of the conditional mean of \mathbf{z} defined in (3-85) are of particular interest to our applications: the conditional mean and covariance of \mathbf{x} . The *conditional mean* of \mathbf{x} , given that \mathbf{y} has assumed the value \mathbf{y} , is generated by letting $\theta(\mathbf{x}) = \mathbf{x}$:

$$E_{\mathbf{x}}[\mathbf{x}|\mathbf{y} = \mathbf{y}] = \int_{-\infty}^{\infty} \xi f_{\mathbf{x}|\mathbf{y}}(\xi|\mathbf{y}) d\xi \quad (3-89)$$

The *conditional covariance* of \mathbf{x} , given that $\mathbf{y}(\omega) = \mathbf{y}$, is then defined as

$$\mathbf{P}_{\mathbf{x}|\mathbf{y}} = E_{\mathbf{x}} \left[(\mathbf{x} - E_{\mathbf{x}}[\mathbf{x}|\mathbf{y} = \mathbf{y}])(\mathbf{x} - E_{\mathbf{x}}[\mathbf{x}|\mathbf{y} = \mathbf{y}])^T | \mathbf{y} = \mathbf{y} \right] \quad (3-90a)$$

$$= \int_{-\infty}^{\infty} (\xi - E_{\mathbf{x}}[\mathbf{x}|\mathbf{y} = \mathbf{y}])(\xi - E_{\mathbf{x}}[\mathbf{x}|\mathbf{y} = \mathbf{y}])^T f_{\mathbf{x}|\mathbf{y}}(\xi|\mathbf{y}) d\xi \quad (3-90b)$$

If we want to generate an estimate of \mathbf{x} using measurement data $\mathbf{y}(\omega) = \mathbf{y}$, one possible estimator that is optimal with respect to many criteria is the random variable $E_{\mathbf{x}}[\mathbf{x}|\mathbf{y} = \mathbf{y}(\cdot)]$. Then $(\mathbf{x} - E_{\mathbf{x}}[\mathbf{x}|\mathbf{y} = \mathbf{y}(\cdot)])$ can be interpreted as the random variable to model the error in the estimate: the difference between \mathbf{x} and our estimate of \mathbf{x} . The conditional mean of this error vector would be zero. Consequently, $\mathbf{P}_{\mathbf{x}|\mathbf{y}}$ would be not only the conditional covariance of \mathbf{x} ,

but also the conditional covariance of the error in our estimate of the value of \mathbf{x} .

EXAMPLE 3.15 Let $f_{\mathbf{x}|y}(\xi|y) = (1/\sqrt{2\pi})\exp\{-\frac{1}{2}(\xi - y)^2\}$. Then the conditional mean and variance of \mathbf{x} are

$$E_{\mathbf{x}}[\mathbf{x}|y = y] = y, \quad E_{\mathbf{x}}[(\mathbf{x} - E_{\mathbf{x}}[\mathbf{x}|y = y])^2|y = y] = 1$$

Thus, for different realizations y of \mathbf{y} , the conditional mean is altered but the conditional variance is unchanged for this particular density. ■

To this point, we have assumed the existence of conditional probability density functions. Although this is not a restrictive assumption for our applications, conditional expectations and probabilities can be defined without assuming such existence. Consider a probability space composed of a sample space Ω , σ -algebra \mathcal{F} , and probability function P ; let \mathbf{x} be a proper random variable (the set $\{\omega: \mathbf{x}(\omega) \leq \xi\}$ is in \mathcal{F}). Now let \mathcal{F}' be a σ -algebra that is a subset of \mathcal{F} (\mathcal{F}' is a "coarser" σ -algebra, with fewer or the same number of elements as \mathcal{F}). The conditional expectation of \mathbf{x} relative to \mathcal{F}' , $E[\mathbf{x}|\mathcal{F}']$, is any ω function that is a proper random variable relative to \mathcal{F}' (the set $\{\omega: E[\mathbf{x}|\mathcal{F}'] \leq \xi\}$ is in \mathcal{F}' ; i.e., measurable relative to \mathcal{F}') satisfying

$$\int_{\Lambda} E\{\mathbf{x}|\mathcal{F}'\} dP(\omega) = \int_{\Lambda} \mathbf{x}(\omega) dP(\omega) \quad (3-91)$$

for any $\Lambda \subset \Omega$ and $\Lambda \in \mathcal{F}'$. Integration over sets in Ω is defined through measure theory. Letting Λ be Ω itself yields the fact that (3-86) is satisfied by the basic definition of a conditional expectation. The existence and uniqueness of such a random variable is guaranteed by the *Radon-Nikodym theorem* [7, 12, 13, 15], whether or not a density function exists at all.

Now let \mathbf{y} be a vector-valued random variable and let \mathcal{F}' be the minimal σ -algebra with respect to which \mathbf{y} is a proper random variable (\mathcal{F}' is generated by complementation and countable intersection and union of sets of the form $\{\omega: \mathbf{y}(\omega) \leq \rho\}$). Let $B \in \mathcal{F}'$ be the set $\{\omega: \mathbf{y}(\omega) = \mathbf{y}\}$. Then $E\{\mathbf{x}|\mathbf{y} = \mathbf{y}\}$ is defined [7, 15] to be

$$E\{\mathbf{x}|\mathbf{y} = \mathbf{y}\} = E\{\mathbf{x}|\mathcal{F}'\} \Big|_{\omega \in B} \quad (3-92)$$

i.e., it is the random variable function of ω , $E\{\mathbf{x}|\mathcal{F}'\}$, evaluated with a particular ω chosen from the set B .

EXAMPLE 3.16 Let us return to the die-toss experiment described in Examples 3.3, 3.5, 3.8, and 3.11. The payoff random variable $y(\cdot)$ can be defined through

$$y(\omega) = \begin{cases} 1 & \text{if } \omega \notin A_1 \text{ or } A_2 \\ 2 & \text{if } \omega \in A_1 \\ 5 & \text{if } \omega \in A_2 \end{cases}$$

with distribution function defined in Fig. 3-15.

Now let z be an indicator of the outcome of the die toss being one of the lower three or upper three numbers:

$$z(\omega) = \begin{cases} 0 & \text{if } \omega \in \{1 \text{ or } 2 \text{ or } 3 \text{ thrown}\} = A_1 \cup A_2 \\ 1 & \text{if } \omega \in \{4 \text{ or } 5 \text{ or } 6 \text{ thrown}\} = (A_1 \cup A_2)^* \end{cases}$$

The smallest σ -algebra \mathcal{F}' associated with z would be

$$\mathcal{F}' = \{\emptyset, \Omega, A_1 \cup A_2, (A_1 \cup A_2)^*\}$$

which is composed of fewer elements than \mathcal{F} (see Example 3.3). Now let us use (3-91) and (3-63) to generate the conditional expectation, $E\{y|\mathcal{F}'\}$:

$$\int_{\Lambda} E\{y|\mathcal{F}'\} dP(\omega) = \int_{\Lambda} x(\omega) dP(\omega) = \int \rho dF_y(\rho) = \sum_i \rho_i \Delta F_y(\rho_i)$$

Let $\Lambda = (A_1 \cup A_2) \in \mathcal{F}'$ to obtain

$$\int_{(A_1 \cup A_2)} E\{y|\mathcal{F}'\} dP(\omega) = \sum_i \rho_i \Delta F_y(\rho_i) = 2 \cdot \frac{1}{3} + 5 \cdot \frac{1}{6} = \frac{3}{2}$$

But $(A_1 \cup A_2)$ is an atom of \mathcal{F}' , so to be a proper random variable on \mathcal{F}' , $E\{y|\mathcal{F}'\}$ must be constant over $(A_1 \cup A_2)$, so this becomes

$$\begin{aligned} \frac{3}{2} &= E\{y|\mathcal{F}'\} \Big|_{\omega \in (A_1 \cup A_2)} \int_{(A_1 \cup A_2)} dP(\omega) \\ &= E\{y|\mathcal{F}'\} \Big|_{\omega \in (A_1 \cup A_2)} \left[\frac{1}{3} + \frac{1}{6} \right] \end{aligned}$$

so $E\{y|\mathcal{F}'\} \Big|_{\omega \in (A_1 \cup A_2)} = 3$. But $z(\omega) = 0$ if $\omega \in (A_1 \cup A_2)$, so (3-92) yields

$$E\{y|z=0\} = E\{y|\mathcal{F}'\} \Big|_{\omega \in (A_1 \cup A_2)} = 3.$$

Similar reasoning then yields

$$E\{y|z=1\} = E\{y|\mathcal{F}'\} \Big|_{\omega \in (A_1 \cup A_2)^*} = 1 \quad \blacksquare$$

If conditional density functions exist, the definitions in (3-92) and (3-89) are equivalent. We will assume such existence for our applications and will exploit the conditional density function conceptualization.

3.8 CHARACTERISTIC FUNCTIONS

If \mathbf{x} is an n -vector-valued random variable, its *characteristic function* $\phi_{\mathbf{x}}(\cdot)$ is defined as a scalar function of the dummy vector $\boldsymbol{\mu}$ as

$$\phi_{\mathbf{x}}(\boldsymbol{\mu}) \triangleq E_{\mathbf{x}}[e^{j\boldsymbol{\mu}^T \mathbf{x}}] \quad (3-93a)$$

$$= \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} e^{j\boldsymbol{\mu}^T \boldsymbol{\xi}} f_{\mathbf{x}}(\boldsymbol{\xi}) d\xi_1 \cdots d\xi_n \quad (3-93b)$$

where $j = \sqrt{-1}$. Fourier transform theory can be used to describe the characteristics of $\phi_{\mathbf{x}}$ in terms of the corresponding $f_{\mathbf{x}}$.

One fundamental reason for considering characteristic functions is that *moments* of a random variable can be generated readily through them. Consider

taking the partial derivative of $\phi_{\mathbf{x}}(\boldsymbol{\mu})$ with respect to the k th component of $\boldsymbol{\mu}$, μ_k :

$$\frac{\partial \phi_{\mathbf{x}}(\boldsymbol{\mu})}{\partial \mu_k} = j \int_{-\infty}^{\infty} \xi_k e^{j\boldsymbol{\mu}^T \boldsymbol{\xi}} f_{\mathbf{x}}(\boldsymbol{\xi}) d\boldsymbol{\xi}$$

Now divide by j and evaluate the result at $\boldsymbol{\mu} = \mathbf{0}$:

$$\frac{1}{j} \left[\frac{\partial \phi_{\mathbf{x}}(\boldsymbol{\mu})}{\partial \mu_k} \right] \bigg|_{\boldsymbol{\mu}=\mathbf{0}} = \int_{-\infty}^{\infty} \xi_k f_{\mathbf{x}}(\boldsymbol{\xi}) d\boldsymbol{\xi} = E[x_k] \quad (3-94)$$

Thus, to obtain the mean of the k th component of \mathbf{x} , for $k = 1, 2, \dots, n$, one can evaluate the partial derivative of $\phi_{\mathbf{x}}(\boldsymbol{\mu})$ with respect to the corresponding component of $\boldsymbol{\mu}$, divide by j , and evaluate the result at $\boldsymbol{\mu} = \mathbf{0}$.

The second moments can be generated through the second partial derivatives. Since

$$\frac{\partial^2 \phi_{\mathbf{x}}(\boldsymbol{\mu})}{\partial \mu_k \partial \mu_l} = j^2 \int_{-\infty}^{\infty} \xi_k \xi_l e^{j\boldsymbol{\mu}^T \boldsymbol{\xi}} f_{\mathbf{x}}(\boldsymbol{\xi}) d\boldsymbol{\xi}$$

the second noncentral moment $E[x_k x_l]$ can be evaluated as

$$E[x_k x_l] = \frac{1}{j^2} \left[\frac{\partial^2 \phi_{\mathbf{x}}(\boldsymbol{\mu})}{\partial \mu_k \partial \mu_l} \right] \bigg|_{\boldsymbol{\mu}=\mathbf{0}} \quad (3-95)$$

In general, an N th noncentral moment of \mathbf{x} can be computed through

$$E[\underbrace{x_k x_l \cdots}_{N \text{ terms}}] = \frac{1}{j^N} \left[\frac{\partial^N \phi_{\mathbf{x}}(\boldsymbol{\mu})}{\partial \mu_k \partial \mu_l \cdots} \right] \bigg|_{\boldsymbol{\mu}=\mathbf{0}} \quad (3-96)$$

Another application of characteristic functions is the description of the *sum of two independent random variables*. Let \mathbf{x} and \mathbf{y} be two independent n -vector valued random variables, and define \mathbf{z} as their sum,

$$\mathbf{z} = \mathbf{x} + \mathbf{y} \quad (3-97)$$

If we know $f_{\mathbf{x}}(\boldsymbol{\xi})$ and $f_{\mathbf{y}}(\boldsymbol{\rho})$, how can we explicitly generate $f_{\mathbf{z}}(\boldsymbol{\zeta})$? Such a question will arise naturally in estimation when we describe the measurements available to us as true variable values corrupted by additive independent noise. To answer this question, first consider the conditional probability of lying in the infinitesimal hypercube in R^n with one corner at $\boldsymbol{\zeta}$ and of dimension $d\zeta_i = \varepsilon$, $i = 1, 2, \dots, n$, on each side. Starting with the definition of the conditional density $f_{\mathbf{z}|\mathbf{x}}(\boldsymbol{\zeta}|\boldsymbol{\xi})$, we can write [see Eq. (3-23)]:

$$\begin{aligned} f_{\mathbf{z}|\mathbf{x}}(\boldsymbol{\zeta}|\boldsymbol{\xi}) d\boldsymbol{\zeta} &= P(\{\omega: \boldsymbol{\zeta} < \mathbf{z}(\omega) \leq \boldsymbol{\zeta} + d\boldsymbol{\zeta}\}, \text{ given that } \mathbf{x}(\omega) = \boldsymbol{\xi}) \\ &= P(\{\omega: \boldsymbol{\zeta} < \mathbf{x}(\omega) + \mathbf{y}(\omega) \leq \boldsymbol{\zeta} + d\boldsymbol{\zeta}\} | \mathbf{x}(\omega) = \boldsymbol{\xi}) \\ &= P(\{\omega: \boldsymbol{\zeta} - \boldsymbol{\xi} < \mathbf{x}(\omega) + \mathbf{y}(\omega) - \boldsymbol{\xi} \leq \boldsymbol{\zeta} + d\boldsymbol{\zeta} - \boldsymbol{\xi}\} | \mathbf{x}(\omega) = \boldsymbol{\xi}) \\ &= P(\{\omega: \boldsymbol{\zeta} - \boldsymbol{\xi} < \mathbf{y}(\omega) \leq \boldsymbol{\zeta} - \boldsymbol{\xi} + d\boldsymbol{\zeta}\} | \mathbf{x}(\omega) = \boldsymbol{\xi}) \end{aligned}$$

But, since \mathbf{x} and \mathbf{y} are independent, this equals the unconditional probability that \mathbf{y} assumes values between the same limits:

$$f_{\mathbf{z}|\mathbf{x}}(\boldsymbol{\zeta}|\boldsymbol{\xi})d\boldsymbol{\zeta} = P(\{\omega: \boldsymbol{\zeta} - \boldsymbol{\xi} < \mathbf{y}(\omega) \leq \boldsymbol{\zeta} - \boldsymbol{\xi} + d\boldsymbol{\zeta}\}) = f_{\mathbf{y}}(\boldsymbol{\zeta} - \boldsymbol{\xi})d\boldsymbol{\zeta}$$

Thus, we have shown that

$$f_{\mathbf{z}|\mathbf{x}}(\boldsymbol{\zeta}|\boldsymbol{\xi}) = f_{\mathbf{y}}(\boldsymbol{\zeta} - \boldsymbol{\xi}) \quad (3-98)$$

By combining the concepts of marginal densities and Bayes' rule, this result can be used to write $f_{\mathbf{z}}(\boldsymbol{\zeta})$ as

$$\begin{aligned} f_{\mathbf{z}}(\boldsymbol{\zeta}) &= \int_{-\infty}^{\infty} f_{\mathbf{z},\mathbf{x}}(\boldsymbol{\zeta},\boldsymbol{\xi})d\boldsymbol{\xi} = \int_{-\infty}^{\infty} f_{\mathbf{z}|\mathbf{x}}(\boldsymbol{\zeta}|\boldsymbol{\xi})f_{\mathbf{x}}(\boldsymbol{\xi})d\boldsymbol{\xi} \\ &= \int_{-\infty}^{\infty} f_{\mathbf{y}}(\boldsymbol{\zeta} - \boldsymbol{\xi})f_{\mathbf{x}}(\boldsymbol{\xi})d\boldsymbol{\xi} \end{aligned} \quad (3-99)$$

This is a *convolution integral*, which is, in general, difficult to evaluate. The corresponding characteristic function is a simple *product*, as expected from the Fourier transform of a convolution.

$$\phi_{\mathbf{z}}(\boldsymbol{\mu}) = E[e^{j\boldsymbol{\mu}^T \mathbf{z}}] = \int_{-\infty}^{\infty} e^{j\boldsymbol{\mu}^T \boldsymbol{\zeta}} f_{\mathbf{z}}(\boldsymbol{\zeta})d\boldsymbol{\zeta}$$

Now substitute (3-99) into this result, letting $\boldsymbol{\zeta} = \boldsymbol{\xi} + \boldsymbol{\rho}$ in the exponential (since $\mathbf{z} = \mathbf{x} + \mathbf{y}$), and assuming we can interchange the order of integration,

$$\begin{aligned} \phi_{\mathbf{z}}(\boldsymbol{\mu}) &= \int_{-\infty}^{\infty} \left[\int_{-\infty}^{\infty} e^{j\boldsymbol{\mu}^T(\boldsymbol{\xi} + \boldsymbol{\rho})} f_{\mathbf{y}}(\boldsymbol{\zeta} - \boldsymbol{\xi}) f_{\mathbf{x}}(\boldsymbol{\xi}) d\boldsymbol{\xi} \right] d\boldsymbol{\zeta} \\ &= \int_{-\infty}^{\infty} \left[\int_{-\infty}^{\infty} e^{j\boldsymbol{\mu}^T(\boldsymbol{\xi} + \boldsymbol{\rho})} f_{\mathbf{y}}(\boldsymbol{\zeta} - \boldsymbol{\xi}) f_{\mathbf{x}}(\boldsymbol{\xi}) d\boldsymbol{\xi} \right] d\boldsymbol{\zeta} \\ &= \int_{-\infty}^{\infty} \left[\int_{-\infty}^{\infty} e^{j\boldsymbol{\mu}^T(\boldsymbol{\xi} + \boldsymbol{\rho})} f_{\mathbf{y}}(\boldsymbol{\rho}) f_{\mathbf{x}}(\boldsymbol{\xi}) d\boldsymbol{\rho} \right] d\boldsymbol{\xi} \\ &= \left[\int_{-\infty}^{\infty} e^{j\boldsymbol{\mu}^T \boldsymbol{\xi}} f_{\mathbf{x}}(\boldsymbol{\xi}) d\boldsymbol{\xi} \right] \left[\int_{-\infty}^{\infty} e^{j\boldsymbol{\mu}^T \boldsymbol{\rho}} f_{\mathbf{y}}(\boldsymbol{\rho}) d\boldsymbol{\rho} \right] \\ &= \phi_{\mathbf{x}}(\boldsymbol{\mu}) \phi_{\mathbf{y}}(\boldsymbol{\mu}) \end{aligned} \quad (3-100)$$

3.9 GAUSSIAN RANDOM VECTORS

A particular random variable of significance to our work is the Gaussian, or normal, vector-valued random variable. First, it provides an adequate model of the random behavior exhibited by many phenomena observed in nature. Second, Gaussian random variables yield tractable mathematical models upon which to base estimators and controllers.

The random n -dimensional vector \mathbf{x} is said to be a *Gaussian (normal) random vector*, or a normally distributed vector-valued random variable, if it can be described through a probability density function of the form

$$f_{\mathbf{x}}(\boldsymbol{\xi}) = \frac{1}{(2\pi)^{n/2} |\mathbf{P}|^{1/2}} \exp \left\{ -\frac{1}{2} [\boldsymbol{\xi} - \mathbf{m}]^T \mathbf{P}^{-1} [\boldsymbol{\xi} - \mathbf{m}] \right\} \quad (3-101)$$

where \mathbf{P} is a positive definite ($n \times n$) matrix, $|\cdot|$ denotes the determinant of a matrix, and $\exp\{\cdot\}$ denotes exponential. The matrix \mathbf{P} must be assumed positive definite to be assured of the existence of \mathbf{P}^{-1} . Actually, a more general definition of a Gaussian random vector, allowing positive semidefinite \mathbf{P} , can be achieved through the characteristic function. We emphasize the somewhat more restrictive characterization in (3-101) because the density function will provide more physical insight in estimation and control.

Note that the density function in (3-101) is completely defined by the two parameters \mathbf{m} and \mathbf{P} . We now claim, and will show later, that these parameters are in fact the mean vector and covariance matrix, respectively. Thus, unlike most other density functions, higher order moments are not required to generate a complete description of the density function.

Figure 3.18 depicts the density function for a scalar Gaussian random variable:

$$f_x(\xi) = \frac{1}{\sqrt{2\pi P}} \exp \left\{ -\frac{1}{2P} (\xi - m)^2 \right\} \quad (3-102)$$

Because the density is symmetric and unimodal (having one peak), m is both the mean and the mode, the value where the density assumes its peak value. The variance P determines the spread of the density about m as in the figure. If σ is the standard deviation, $\sigma = \sqrt{P}$, then 68.3% of the area under the curve lies in the interval between $(m - \sigma)$ and $(m + \sigma)$. Stated another way, the probability that x assumes a value in the interval $[m - \sigma, m + \sigma]$ is 0.683. Similarly, 95.4% of the probability weight lies between $(m - 2\sigma)$ and $(m + 2\sigma)$, and 99.7% lies between $(m - 3\sigma)$ and $(m + 3\sigma)$. For this reason, peak error specifications are often converted to 3σ values in practice if Gaussian models are employed.

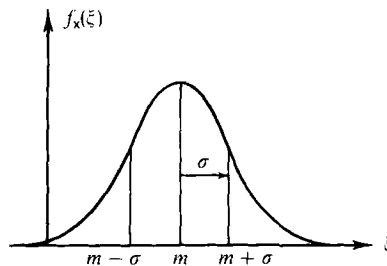


FIG. 3.18 Density function of a scalar Gaussian random variable.

A two-dimensional Gaussian random vector would be characterized by the density function

$$\begin{aligned}
 f_{\mathbf{x}}(\boldsymbol{\xi}) &= (2\pi)^{-1} \left| \begin{bmatrix} \sigma_1^2 & r_{12}\sigma_1\sigma_2 \\ r_{12}\sigma_1\sigma_2 & \sigma_2^2 \end{bmatrix} \right|^{-1/2} \\
 &\quad \times \exp -\frac{1}{2} \left\{ \begin{bmatrix} \xi_1 - m_1 \\ \xi_2 - m_2 \end{bmatrix}^T \begin{bmatrix} \sigma_1^2 & r_{12}\sigma_1\sigma_2 \\ r_{12}\sigma_1\sigma_2 & \sigma_2^2 \end{bmatrix}^{-1} \begin{bmatrix} \xi_1 - m_1 \\ \xi_2 - m_2 \end{bmatrix} \right\} \\
 &= \frac{1}{2\pi\sigma_1\sigma_2(1-r_{12}^2)^{1/2}} \exp \left\{ -\frac{1}{2(1-r_{12}^2)^2} \left[\frac{(\xi_1 - m_1)^2}{\sigma_1^2} \right. \right. \\
 &\quad \left. \left. + \frac{(\xi_2 - m_2)^2}{\sigma_2^2} - \frac{2r_{12}(\xi_1 - m_1)(\xi_2 - m_2)}{\sigma_1\sigma_2} \right] \right\} \quad (3-103)
 \end{aligned}$$

This is presented graphically in Fig. 3.19. The mean vector \mathbf{m} in the ξ_1 - ξ_2 plane locates the peak of the density function. Loci of constant density function

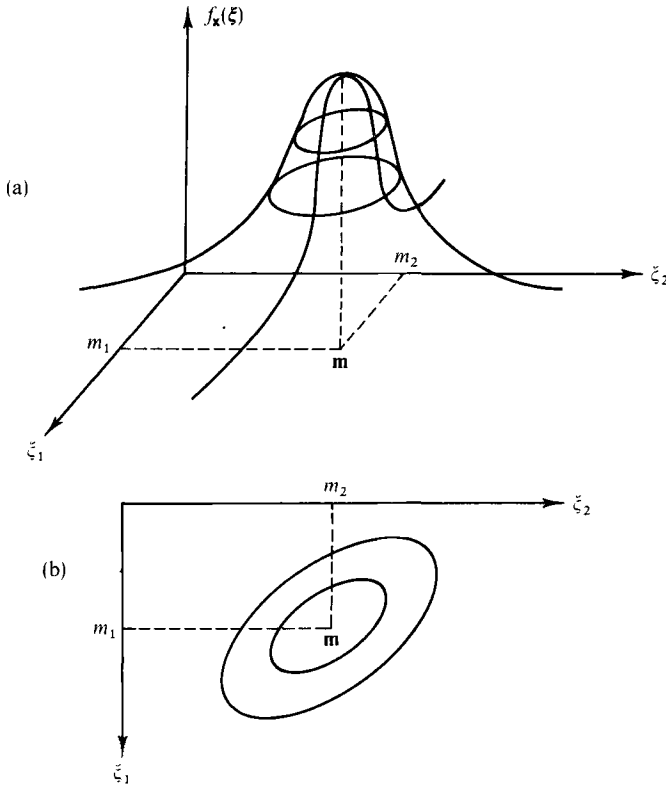


FIG. 3.19 Density function for a two-dimensional Gaussian random vector. (a) Three-dimensional depiction. (b) View from above. The ellipses are the loci of constant probability density value.

values, called surfaces of constant likelihood, are generated by passing planes parallel to the ξ_1 - ξ_2 plane through the density function surface, and are ellipses parallel to the ξ_1 - ξ_2 plane as shown in the diagram. This can also be seen by setting (3-103) equal to some constant, or equivalently,

$$\frac{(\xi_1 - m_1)^2}{\sigma_1^2} + \frac{(\xi_2 - m_2)^2}{\sigma_2^2} - \frac{2r_{12}(\xi_1 - m_1)(\xi_2 - m_2)}{\sigma_1\sigma_2} = k \quad (3-104)$$

which is the general equation of an ellipse in the ξ_1 - ξ_2 plane. Thus the covariance matrix determines the size and angular orientation of ellipses of constant likelihood. If the correlation coefficient r_{12} is zero and thus \mathbf{P} is diagonal, then the principal axes of the ellipses are parallel to the ξ_1 and ξ_2 axes, and $\sigma_1 = \sqrt{P_{11}}$ and $\sigma_2 = \sqrt{P_{22}}$ are the magnitudes of the semimajor and semiminor axes dimensions for the one-sigma ellipse. [This is readily apparent from (3-104).] In general, the eigenvalues of \mathbf{P} provide these magnitudes. If \mathbf{P} is singular and of rank one, then the density function surface collapses down to zero except over a single line (the limit of the ellipses) in the ξ_1 - ξ_2 plane: there is no uncertainty in the direction orthogonal to the line since you *know* \mathbf{x} assumes a value somewhere on the line.

These ideas generalize to higher-dimensional cases, with surfaces of constant likelihood becoming n -dimensional ellipsoids. For example, a probabilistic description of position in three dimensions would be expressible in terms of the size, shape, and orientation of the three-dimensional ellipsoid corresponding to a given probability that the true position lies within the ellipsoid (see Problem 3.11).

The *characteristic function* for a Gaussian random variable \mathbf{x} with density function as in (3-101) is

$$\phi_{\mathbf{x}}(\boldsymbol{\mu}) = \exp\{j\boldsymbol{\mu}^T \mathbf{m} - \frac{1}{2}\boldsymbol{\mu}^T \mathbf{P} \boldsymbol{\mu}\} \quad (3-105)$$

To show this, we will look at a density function in the general form of (3-101), translate our coordinate system origin to the mean location, and then rotate the coordinates to align them with the principal axes of the ellipsoids of constant likelihood. After working with this simpler description of the density, we will rotate and translate back to the original coordinate system to express the result. The geometric insights gained and linear algebra involved warrant presentation of the details of the proof.

DERIVATION OF CHARACTERISTIC FUNCTION By the definition of characteristic functions, we can write

$$\begin{aligned} \phi_{\mathbf{x}}(\boldsymbol{\mu}) &= \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} \exp\{j\boldsymbol{\mu}^T \boldsymbol{\xi}\} f_{\mathbf{x}}(\boldsymbol{\xi}) d\xi_1 \cdots d\xi_n \\ &= \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} \frac{1}{(2\pi)^{n/2} |\mathbf{P}|^{1/2}} \exp\{j\boldsymbol{\mu}^T \boldsymbol{\xi} - \frac{1}{2}(\boldsymbol{\xi} - \mathbf{m})^T \mathbf{P}^{-1} (\boldsymbol{\xi} - \mathbf{m})\} d\xi_1 \cdots d\xi_n \end{aligned}$$

Now translate the coordinate system by making a change of variable $\gamma = \xi - \mathbf{m}$. Note that $d\gamma_i = d\xi_i$ for $i = 1, 2, \dots, n$. For convenience, define the scalar a as

$$a = \frac{\exp\{j\boldsymbol{\mu}^T \mathbf{m}\}}{(2\pi)^{n/2} |\mathbf{P}|^{1/2}}$$

Thus, $\phi_{\mathbf{x}}(\boldsymbol{\mu})$ becomes

$$\phi_{\mathbf{x}}(\boldsymbol{\mu}) = a \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} \exp\{j\boldsymbol{\mu}^T \gamma - \frac{1}{2} \gamma^T \mathbf{P}^{-1} \gamma\} d\gamma_1 \cdots d\gamma_n$$

Having translated the origin of the coordinates to \mathbf{m} , rotate the axes into the principal directions. Since \mathbf{P} is symmetric, \mathbf{P}^{-1} is also symmetric. Therefore, there exists an orthogonal transformation matrix \mathbf{A} which diagonalizes \mathbf{P}^{-1} : there exists a matrix such that $\mathbf{A}^T = \mathbf{A}^{-1}$ and

$$\mathbf{A}^T \mathbf{P}^{-1} \mathbf{A} = \begin{bmatrix} \sigma_1^{-2} & & 0 \\ & \sigma_2^{-2} & \\ 0 & & \sigma_n^{-2} \end{bmatrix} \quad \text{or} \quad \mathbf{P}^{-1} = \mathbf{A} \begin{bmatrix} \sigma_1^{-2} & & 0 \\ & \sigma_2^{-2} & \\ 0 & & \sigma_n^{-2} \end{bmatrix} \mathbf{A}^T$$

Thus,

$$\phi_{\mathbf{x}}(\boldsymbol{\mu}) = a \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} \exp\left\{j\boldsymbol{\mu}^T \gamma - \frac{1}{2} \gamma^T \mathbf{A} \begin{bmatrix} \sigma_1^{-2} & & 0 \\ & \sigma_2^{-2} & \\ 0 & & \sigma_n^{-2} \end{bmatrix} \mathbf{A}^T \gamma\right\} d\gamma_1 \cdots d\gamma_n$$

Define a new set of variables through a coordinate rotation as

$$\boldsymbol{\rho} = \mathbf{A}^T \boldsymbol{\mu} \leftrightarrow \boldsymbol{\mu} = \mathbf{A} \boldsymbol{\rho}, \quad \boldsymbol{\zeta} = \mathbf{A}^T \gamma \leftrightarrow \gamma = \mathbf{A} \boldsymbol{\zeta}$$

When we change from integrating over $d\gamma_1 \cdots d\gamma_n$ to $d\zeta_1 \cdots d\zeta_n$, the Jacobian determinant is equal to one because of the orthogonality of \mathbf{A} :

$$d\zeta_1 \cdots d\zeta_n = \begin{vmatrix} \partial\zeta_1/\partial\gamma_1 & \cdots & \partial\zeta_1/\partial\gamma_n \\ \vdots & & \vdots \\ \partial\zeta_n/\partial\gamma_1 & \cdots & \partial\zeta_n/\partial\gamma_n \end{vmatrix} d\gamma_1 \cdots d\gamma_n = (1) d\gamma_1 \cdots d\gamma_n$$

With this change of variables, we can write

$$\begin{aligned} \phi_{\mathbf{x}}(\boldsymbol{\mu}) &= a \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} \exp\left\{j\boldsymbol{\rho}^T \mathbf{A}^T \mathbf{A} \boldsymbol{\zeta} - \frac{1}{2} \boldsymbol{\zeta}^T \begin{bmatrix} \sigma_1^{-2} & & 0 \\ & \sigma_2^{-2} & \\ 0 & & \sigma_n^{-2} \end{bmatrix} \boldsymbol{\zeta}\right\} d\zeta_1 \cdots d\zeta_n \\ &= a \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} \exp\left\{\sum_{i=1}^n (j\rho_i \zeta_i - \frac{1}{2} \zeta_i^2 / \sigma_i^2)\right\} d\zeta_1 \cdots d\zeta_n \\ &= a \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} \prod_{i=1}^n \exp\{j\rho_i \zeta_i - \frac{1}{2} \zeta_i^2 / \sigma_i^2\} d\zeta_1 \cdots d\zeta_n \\ &= a \prod_{i=1}^n \int_{-\infty}^{\infty} \exp\{j\rho_i \zeta_i - \frac{1}{2} \zeta_i^2 / \sigma_i^2\} d\zeta_i \end{aligned}$$

This is now in the form of a product of n separate one-dimensional integrals. To evaluate each integral, complete the square in the exponential to form

$$\begin{aligned}\int_{-\infty}^{\infty} \exp\{ \} d\zeta_i &= \int_{-\infty}^{\infty} \exp\{ -\frac{1}{2}(\zeta_i - j\rho_i\sigma_i^2)^2/\sigma_i^2 - \frac{1}{2}\rho_i^2\sigma_i^2 \} d\zeta_i \\ &= \int_{-\infty}^{\infty} \exp\{ -\frac{1}{2}\beta_i^2/\sigma_i^2 - \frac{1}{2}\rho_i^2\sigma_i^2 \} d\beta_i \\ &= \exp\{ -\frac{1}{2}\rho_i^2\sigma_i^2 \} \int_{-\infty}^{\infty} \exp\{ -\frac{1}{2}\beta_i^2/\sigma_i^2 \} d\beta_i\end{aligned}$$

But now the integral term can be recognized as the integral of a scaled Gaussian density of mean zero and variance σ_i^2 , so by putting in the proper coefficient, we know the integral equals one:

$$\begin{aligned}\int_{-\infty}^{\infty} \exp\{ \} d\zeta_i &= \sqrt{2\pi}\sigma_i \exp\{ -\frac{1}{2}\rho_i^2\sigma_i^2 \} \left[\frac{1}{\sqrt{2\pi}\sigma_i} \int_{-\infty}^{\infty} \exp\{ -\frac{1}{2}\beta_i^2/\sigma_i^2 \} d\beta_i \right] \\ &= \sqrt{2\pi}\sigma_i \exp\{ -\frac{1}{2}\rho_i^2\sigma_i^2 \} [1]\end{aligned}$$

The characteristic function is now the product of n such integrals as just shown. Writing a explicitly in the expression yields

$$\phi_{\mathbf{x}}(\boldsymbol{\mu}) = \left[\frac{\exp\{j\boldsymbol{\mu}^T \mathbf{m}\}}{(2\pi)^{n/2} |\mathbf{P}|^{1/2}} \right] \prod_{i=1}^n (\sqrt{2\pi}\sigma_i \exp\{ -\frac{1}{2}\rho_i^2\sigma_i^2 \})$$

The $(2\pi)^{n/2}$ in the denominator cancels $\prod_{i=1}^n \sqrt{2\pi}$. Moreover, since $\sigma_1^2, \sigma_2^2, \dots, \sigma_n^2$ are the eigenvalues of \mathbf{P} , $|\mathbf{P}| = \sigma_1^2 \sigma_2^2 \cdots \sigma_n^2$ and so $|\mathbf{P}|^{1/2} = \sigma_1 \sigma_2 \cdots \sigma_n$, and this then cancels $\prod_{i=1}^n \sigma_i$. Thus,

$$\begin{aligned}\phi_{\mathbf{x}}(\boldsymbol{\mu}) &= \exp\{j\boldsymbol{\mu}^T \mathbf{m}\} \exp\left\{ -\frac{1}{2} \sum_{i=1}^n \rho_i^2 \sigma_i^2 \right\} \\ &= \exp\{j\boldsymbol{\mu}^T \mathbf{m}\} \exp\left\{ -\frac{1}{2} \boldsymbol{\rho}^T \begin{bmatrix} \sigma_1^2 & & 0 \\ & \ddots & \\ 0 & & \sigma_n^2 \end{bmatrix} \boldsymbol{\rho} \right\} \\ &= \exp\{j\boldsymbol{\mu}^T \mathbf{m}\} \exp\left\{ -\frac{1}{2} \boldsymbol{\mu}^T \mathbf{A} \begin{bmatrix} \sigma_1^2 & & 0 \\ & \ddots & \\ 0 & & \sigma_n^2 \end{bmatrix} \mathbf{A}^T \boldsymbol{\mu} \right\}\end{aligned}$$

Since $\mathbf{P}^{-1} = \mathbf{A}[\sigma_i^{-2}] \mathbf{A}^T$,

$$\mathbf{P} = (\mathbf{A}[\sigma_i^{-2}] \mathbf{A}^T)^{-1} = (\mathbf{A}^T)^{-1} [\sigma_i^{-2}]^{-1} \mathbf{A}^{-1} = \mathbf{A}[\sigma_i^2] \mathbf{A}^T$$

Substituting this into the quadratic form in $\phi_{\mathbf{x}}(\boldsymbol{\mu})$ yields

$$\begin{aligned}\phi_{\mathbf{x}}(\boldsymbol{\mu}) &= \exp\{j\boldsymbol{\mu}^T \mathbf{m}\} \exp\{ -\frac{1}{2} \boldsymbol{\mu}^T \mathbf{P} \boldsymbol{\mu} \} \\ &= \exp\{j\boldsymbol{\mu}^T \mathbf{m} - \frac{1}{2} \boldsymbol{\mu}^T \mathbf{P} \boldsymbol{\mu}\}\end{aligned}$$

as claimed in (3-105). ■

Note that the characteristic function does not involve \mathbf{P}^{-1} , and therefore it does not inherently require \mathbf{P} to be positive definite. In fact, (3-105) is valid for positive semidefinite \mathbf{P} .

The characteristic function can now be used to generate the moments of the Gaussian random vector \mathbf{x} . It will be shown now that the parameters \mathbf{m} and

\mathbf{P} in the density and characteristic functions are the *mean* and *covariance* of \mathbf{x} , respectively:

$$E[\mathbf{x}] = \mathbf{m} \quad (3-106)$$

$$E[\mathbf{x}\mathbf{x}^T] = \mathbf{P} + \mathbf{m}\mathbf{m}^T \quad (3-107)$$

$$E[(\mathbf{x} - \mathbf{m})(\mathbf{x} - \mathbf{m})^T] = \mathbf{P} \quad (3-108)$$

DERIVATION OF MEAN AND COVARIANCE The characteristic function for Gaussian \mathbf{x} is

$$\begin{aligned} \phi_{\mathbf{x}}(\boldsymbol{\mu}) &= \exp\{j\boldsymbol{\mu}^T\mathbf{m} - \tfrac{1}{2}\boldsymbol{\mu}^T\mathbf{P}\boldsymbol{\mu}\} \\ &= \exp\left\{j \sum_{i=1}^n \mu_i m_i - \frac{1}{2} \sum_{i=1}^n \sum_{q=1}^n \mu_i \mu_q P_{iq}\right\} \end{aligned}$$

To generate the mean of the k th component of \mathbf{x} , we take the partial of $\phi_{\mathbf{x}}(\boldsymbol{\mu})$ with respect to μ_k .

$$\frac{\partial \phi_{\mathbf{x}}(\boldsymbol{\mu})}{\partial \mu_k} = \left(jm_k - \sum_{q=1}^n P_{kq} \mu_q \right) \phi_{\mathbf{x}}(\boldsymbol{\mu})$$

so we can write

$$E[x_k] = \frac{1}{j} \frac{\partial \phi_{\mathbf{x}}(\boldsymbol{\mu})}{\partial \mu_k} \bigg|_{\boldsymbol{\mu}=0} = m_k$$

Since this is true for all k , $k = 1, 2, \dots, n$, $E[\mathbf{x}] = \mathbf{m}$.

Using the first partial expression just given, the second partials are

$$\begin{aligned} \frac{\partial^2 \phi_{\mathbf{x}}(\boldsymbol{\mu})}{\partial \mu_k \partial \mu_l} &= \left[\frac{\partial}{\partial \mu_l} \left(jm_k - \sum_{q=1}^n P_{kq} \mu_q \right) \right] \phi_{\mathbf{x}}(\boldsymbol{\mu}) + \left(jm_k - \sum_{q=1}^n P_{kq} \mu_q \right) \frac{\partial \phi_{\mathbf{x}}(\boldsymbol{\mu})}{\partial \mu_l} \\ &= (-P_{kl}) \phi_{\mathbf{x}}(\boldsymbol{\mu}) + \left(jm_k - \sum_{q=1}^n P_{kq} \mu_q \right) \left(jm_l - \sum_{q=1}^n P_{lq} \mu_q \right) \phi_{\mathbf{x}}(\boldsymbol{\mu}) \end{aligned}$$

so the second moment $E[x_k x_l]$ is

$$E[x_k x_l] = \frac{1}{j^2} \frac{\partial^2 \phi_{\mathbf{x}}(\boldsymbol{\mu})}{\partial \mu_k \partial \mu_l} \bigg|_{\boldsymbol{\mu}=0} = P_{kl} + m_k m_l$$

This is true for all k and l , so we obtain $E[\mathbf{x}\mathbf{x}^T] = \mathbf{P} + \mathbf{m}\mathbf{m}^T$. ■

EXAMPLE 3.17 Consider a zero-mean Gaussian random vector, with

$$f_{\mathbf{x}}(\boldsymbol{\xi}) = [(2\pi)^n |\mathbf{P}|^{1/2}]^{-1} \exp\{-\tfrac{1}{2}\boldsymbol{\xi}^T \mathbf{P}^{-1} \boldsymbol{\xi}\}, \quad \phi_{\mathbf{x}}(\boldsymbol{\mu}) = \exp\{-\tfrac{1}{2}\boldsymbol{\mu}^T \mathbf{P} \boldsymbol{\mu}\}$$

For such a random variable, the characteristic function can be used to generate the first four moments as

$$\begin{aligned} E[x_k] &= 0 & E[x_k x_l x_m] &= 0 \\ E[x_k x_l] &= P_{kl} & E[x_k x_l x_m x_n] &= P_{kl} P_{mn} + P_{km} P_{ln} + P_{kn} P_{lm} \end{aligned} \quad \blacksquare$$

Generalizing the results of the previous example, all odd central moments of a Gaussian random vector are zero (due to symmetry). Moreover, all even central moments can be expressed in terms of the covariance. This is just

another way of saying that the mean and covariance completely define the Gaussian density function.

Previously it was shown that independence implies uncorrelatedness but not necessarily vice versa. It will now be shown that *two jointly Gaussian (normal) random vectors which are uncorrelated are also independent*. We must assume *jointly Gaussian* (defined in the following) random vectors for this to be true— \mathbf{x} and \mathbf{y} can be Gaussian vectors that are *not* jointly Gaussian, and the implication is then not true.

DEMONSTRATION THAT UNCORRELATED \rightarrow INDEPENDENT IF JOINTLY GAUSSIAN Suppose \mathbf{x} and \mathbf{y} are random vectors, of dimensions n and m , respectively, that are jointly Gaussian and uncorrelated. Define \mathbf{z} to be a random vector of dimension $(n + m)$ composed of components \mathbf{x} and \mathbf{y} :

$$\mathbf{z} = \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix}$$

To say \mathbf{x} and \mathbf{y} are jointly Gaussian is equivalent to saying that \mathbf{z} is Gaussian. The first two moments of \mathbf{z} would be

$$\mathbf{m}_z = E[\mathbf{z}] = \begin{bmatrix} E[\mathbf{x}] \\ E[\mathbf{y}] \end{bmatrix} = \begin{bmatrix} \mathbf{m}_x \\ \mathbf{m}_y \end{bmatrix}$$

$$E[\mathbf{z}\mathbf{z}^T] = \begin{bmatrix} E[\mathbf{x}\mathbf{x}^T] & E[\mathbf{x}\mathbf{y}^T] \\ E[\mathbf{y}\mathbf{x}^T] & E[\mathbf{y}\mathbf{y}^T] \end{bmatrix}; \quad \mathbf{P}_{zz} = E[\mathbf{z}\mathbf{z}^T] - \mathbf{m}_z\mathbf{m}_z^T$$

Now, since \mathbf{x} and \mathbf{y} are uncorrelated,

$$E[\mathbf{x}\mathbf{y}^T] = \mathbf{m}_x\mathbf{m}_y^T; \quad E[\mathbf{y}\mathbf{x}^T] = \mathbf{m}_y\mathbf{m}_x^T$$

Thus, the covariance \mathbf{P}_{zz} becomes block diagonal:

$$\mathbf{P}_{zz} = \begin{bmatrix} E[\mathbf{x}\mathbf{x}^T] & \mathbf{m}_x\mathbf{m}_y^T \\ \mathbf{m}_y\mathbf{m}_x^T & E[\mathbf{y}\mathbf{y}^T] \end{bmatrix} - \begin{bmatrix} \mathbf{m}_x\mathbf{m}_x^T & \mathbf{m}_x\mathbf{m}_y^T \\ \mathbf{m}_y\mathbf{m}_x^T & \mathbf{m}_y\mathbf{m}_y^T \end{bmatrix} = \begin{bmatrix} \mathbf{P}_{xx} & \mathbf{0} \\ \mathbf{0} & \mathbf{P}_{yy} \end{bmatrix}$$

Letting ζ be composed of partitions ξ (\mathbf{x} values) and ρ (\mathbf{y} values):

$$\zeta = \begin{bmatrix} \xi \\ \rho \end{bmatrix}$$

the density $f_z(\zeta) = f_{\mathbf{x},\mathbf{y}}(\xi, \rho)$ can now be written as:

$$\begin{aligned} f_z(\zeta) &= [(2\pi)^{(n+m)/2} |\mathbf{P}_{zz}|^{1/2}]^{-1} \exp\{-\frac{1}{2}[\zeta - \mathbf{m}_z]^T \mathbf{P}_{zz}^{-1} [\zeta - \mathbf{m}_z]\} \\ &= \left[(2\pi)^{(n+m)/2} \left| \begin{bmatrix} \mathbf{P}_{xx} & \mathbf{0} \\ \mathbf{0} & \mathbf{P}_{yy} \end{bmatrix} \right|^{1/2} \right]^{-1} \exp\left\{-\frac{1}{2} \begin{bmatrix} \xi - \mathbf{m}_x \\ \rho - \mathbf{m}_y \end{bmatrix}^T \begin{bmatrix} \mathbf{P}_{xx}^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{P}_{yy}^{-1} \end{bmatrix} \begin{bmatrix} \xi - \mathbf{m}_x \\ \rho - \mathbf{m}_y \end{bmatrix}\right\} \\ &= [(2\pi)^{n/2} (2\pi)^{m/2} |\mathbf{P}_{xx}|^{1/2} |\mathbf{P}_{yy}|^{1/2}]^{-1} \\ &\quad \cdot \exp\{-\frac{1}{2}[\xi - \mathbf{m}_x]^T \mathbf{P}_{xx}^{-1} [\xi - \mathbf{m}_x] - \frac{1}{2}[\rho - \mathbf{m}_y]^T \mathbf{P}_{yy}^{-1} [\rho - \mathbf{m}_y]\} \\ &= [(2\pi)^{n/2} |\mathbf{P}_{xx}|^{1/2}]^{-1} \exp\{-\frac{1}{2}[\xi - \mathbf{m}_x]^T \mathbf{P}_{xx}^{-1} [\xi - \mathbf{m}_x]\} \\ &\quad \cdot [(2\pi)^{m/2} |\mathbf{P}_{yy}|^{1/2}]^{-1} \exp\{-\frac{1}{2}[\rho - \mathbf{m}_y]^T \mathbf{P}_{yy}^{-1} [\rho - \mathbf{m}_y]\} \\ &= [f_x(\xi)] \cdot [f_y(\rho)] \end{aligned}$$

Thus, \mathbf{x} and \mathbf{y} are independent. \blacksquare

It was mentioned previously that Gaussian random variables are of engineering importance because they provide adequate models of many random phenomena observed empirically. The basic justification for this statement is embodied in the *central limit theorem*; one of its numerous precise statements (differing in specific assumptions and details, but all essentially the same) is now stated.

CENTRAL LIMIT THEOREM Let \mathbf{x}_i , $i = 1, 2, \dots, N$, be a set of independent random n -vectors which are identically distributed with means and covariance matrices \mathbf{m}_i and \mathbf{P}_i , respectively. Define the random vector \mathbf{y}_N as their sum:

$$\mathbf{y}_N = \sum_{i=1}^N \mathbf{x}_i$$

and also define \mathbf{z}_N as the (zero-mean) normalized sum random variable:

$$\mathbf{z}_N = [\mathbf{P}_{y_N y_N}]^{-1/2} [\mathbf{y}_N - E[\mathbf{y}_N]]$$

where

$$E[\mathbf{y}_N] = \sum_{i=1}^N \mathbf{m}_i, \quad \mathbf{P}_{y_N y_N} = \sum_{i=1}^N \mathbf{P}_i, \quad \text{and} \quad \mathbf{P}^{-1/2} = (\mathbf{P}^{1/2})^{-1}$$

where $\mathbf{P}^{1/2}$ is defined as the n -by- n matrix such that $\mathbf{P}^{1/2} \mathbf{P}^{1/2 T} = \mathbf{P}$. Then, in the limit as $N \rightarrow \infty$, \mathbf{z}_N becomes a zero-mean Gaussian random n -vector with a covariance matrix equal to the identity matrix:

$$\lim_{N \rightarrow \infty} f_{\mathbf{z}_N}(\boldsymbol{\zeta}) = [(2\pi)^{n/2}]^{-1} \exp\left\{-\frac{1}{2} \boldsymbol{\zeta}^T \boldsymbol{\zeta}\right\} \quad \blacksquare$$

Actually, more general statements can be made, such as not requiring identical distributions for the random variables being summed and then adding some additional, though not very restrictive, assumptions [1-5, 7-11, 14].

Essentially, the theorem states that if the random phenomenon we observe is generated as the sum of effects of many independent random phenomena, then the distribution of the observed phenomenon approaches a Gaussian distribution as more random effects are summed, *regardless* of the distribution of each individual phenomenon. In practice, the assumptions in the theorem are seldom verifiable. Rather, if there are a large number of additive contributing effects to a random phenomenon (as is usually the case when one probes beyond a macroscopic view of a phenomenon), then one suspects that a Gaussian distribution is a reasonable approximation to the actual distribution.

The theorem claims only that a Gaussian distribution is approached as N grows without bound. One would then logically ask, how large does N have to be before the Gaussian approximation is reasonable? The following example due to Papoulis [9] demonstrates a surprisingly good approximation for $N = 3$ and scalar x_i 's uniformly distributed (each thus having a distribution very different from Gaussian).

EXAMPLE 3.18 Let x_1, x_2 , and x_3 each be uniformly distributed on the interval $[0, T]$, as in Fig. 3.20a. If $y_2 = x_1 + x_2$, then y_2 has a triangular density function (verifiable by convolution) as in Fig. 3.20b, with mean T and variance $T^2/6$; also plotted is the Gaussian density function with the same first two moments. If $y_3 = x_1 + x_2 + x_3$, its density function consists of three parabolic pieces as in Fig. 3.20c, with mean $3T/2$ and variance $T^2/4$. The normal density with these same statistics is a very good approximation to the true density. ■

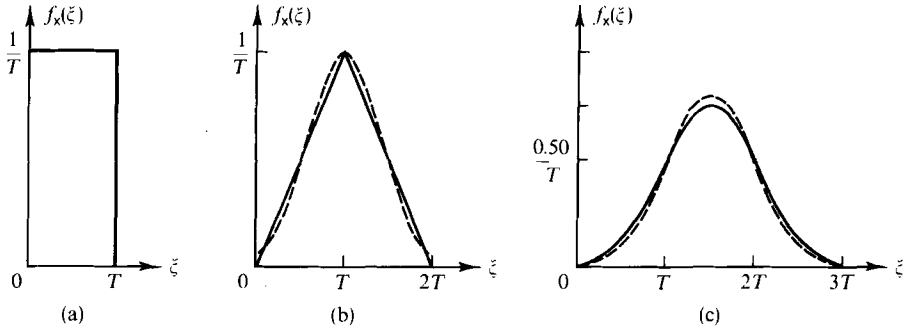


FIG. 3.20 Central limit theorem exemplified. (a) $f_{x_1}(\xi)$. (b) $x = x_1 + x_2$. Solid line indicates $f_x(\xi)$; dashed, $(1/T)\sqrt{3/\pi} \exp[-3(x - T)^2/T^2]$. (c) $x = x_1 + x_2 + x_3$. Solid line indicates $f_x(\xi)$; dashed, $(1/T)\sqrt{2/\pi} \exp[-2(x - 1.5T)^2/T^2]$. From *Probability, Random Variables, and Stochastic Processes* by A. Papoulis. © 1965. Used with permission of McGraw-Hill Book Co.

Later when estimation is discussed, the *conditional Gaussian density* will be of primary interest. Therefore, it will be characterized more fully at this point. Let \mathbf{x} and \mathbf{y} be jointly Gaussian vectors mapping Ω into R^n and R^m , respectively, so that $f_{\mathbf{x},\mathbf{y}}(\xi, \rho)$ can be written as

$$f_{\mathbf{x},\mathbf{y}}(\xi, \rho) = \left[(2\pi)^{(n+m)/2} \left| \begin{bmatrix} \mathbf{P}_{xx} & \mathbf{P}_{xy} \\ \mathbf{P}_{yx} & \mathbf{P}_{yy} \end{bmatrix} \right| \right]^{-1} \times \exp \left\{ -\frac{1}{2} \begin{bmatrix} \xi - \mathbf{m}_x \\ \rho - \mathbf{m}_y \end{bmatrix}^T \begin{bmatrix} \mathbf{P}_{xx} & \mathbf{P}_{xy} \\ \mathbf{P}_{yx} & \mathbf{P}_{yy} \end{bmatrix}^{-1} \begin{bmatrix} \xi - \mathbf{m}_x \\ \rho - \mathbf{m}_y \end{bmatrix} \right\} \quad (3-109)$$

where we assume that the covariance matrix in (3-109) is positive definite. We claim here, and will prove in the next section, that \mathbf{x} is thus a Gaussian n -vector of mean \mathbf{m}_x and covariance \mathbf{P}_{xx} , and \mathbf{y} is a Gaussian m -vector of mean \mathbf{m}_y and covariance \mathbf{P}_{yy} . To obtain the conditional density $f_{\mathbf{x}|\mathbf{y}}(\xi|\rho)$, Bayes' rule can be used to write

$$f_{\mathbf{x}|\mathbf{y}}(\xi|\rho) = f_{\mathbf{x},\mathbf{y}}(\xi, \rho) / f_{\mathbf{y}}(\rho) \quad (3-110)$$

where $f_{\mathbf{x},\mathbf{y}}(\xi, \rho)$ is given by (3-109) and $f_{\mathbf{y}}(\rho)$ is Gaussian, with moments \mathbf{m}_y and \mathbf{P}_{yy} . Performing algebraic reduction yields the result as

$$f_{\mathbf{x}|\mathbf{y}}(\xi|\rho) = \frac{1}{(2\pi)^{n/2} |\mathbf{P}_{x|y}|^{1/2}} \exp \left\{ -\frac{1}{2} [\xi - \mathbf{m}_{x|y}]^T \mathbf{P}_{x|y}^{-1} [\xi - \mathbf{m}_{x|y}] \right\} \quad (3-111)$$

where

$$\mathbf{m}_{x|y} = \mathbf{m}_x + \mathbf{P}_{xy} \mathbf{P}_{yy}^{-1} (\rho - \mathbf{m}_y) \quad (3-112a)$$

$$\mathbf{P}_{x|y} = \mathbf{P}_{xx} - \mathbf{P}_{xy} \mathbf{P}_{yy}^{-1} \mathbf{P}_{yx} \quad (3-112b)$$

Thus, if \mathbf{x} and \mathbf{y} are jointly Gaussian, with joint density given by (3-109), then $f_{x|y}(\xi|\rho)$ is Gaussian with moments $\mathbf{m}_{x|y}$ and $\mathbf{P}_{x|y}$ as just given. The conditional mean of \mathbf{x} , given that $\mathbf{y}(\omega) = \mathbf{y}$, is then

$$E_{\mathbf{x}}[\mathbf{x}|\mathbf{y} = \mathbf{y}] \triangleq \mathbf{m}_{x|y} = \mathbf{m}_x + \mathbf{P}_{xy} \mathbf{P}_{yy}^{-1} (\mathbf{y} - \mathbf{m}_y) \quad (3-113)$$

From this expression, $E_{\mathbf{x}}[\mathbf{x}|\mathbf{y} = \cdot]$ can be seen to be an *explicit function* of the realizations \mathbf{y} of \mathbf{y} , as stated previously in Section 3.7 for conditional expectations in general. Furthermore, the conditional covariance is:

$$E_{\mathbf{x}}\{[\mathbf{x} - \mathbf{m}_{x|y}][\mathbf{x} - \mathbf{m}_{x|y}]^T | \mathbf{y} = \mathbf{y}\} \triangleq \mathbf{P}_{x|y} = \mathbf{P}_{xx} - \mathbf{P}_{xy} \mathbf{P}_{yy}^{-1} \mathbf{P}_{yx} \quad (3-114)$$

Finally, since $f_{x|y}(\xi|\rho)$ is Gaussian with mean and covariance as described, the conditional characteristic function is:

$$\phi_{x|y}(\boldsymbol{\mu}|\rho) = \exp\{j\boldsymbol{\mu}^T \mathbf{m}_{x|y} - \frac{1}{2} \boldsymbol{\mu}^T \mathbf{P}_{x|y} \boldsymbol{\mu}\} \quad (3-115)$$

As mentioned before, this function is properly defined for $\mathbf{P}_{x|y}$ positive semi-definite, whereas the density function (3-111) requires $\mathbf{P}_{x|y}$ to be positive definite to ensure the existence of $\mathbf{P}_{x|y}^{-1}$.

If \mathbf{x} represents variables of interest and \mathbf{y} models the measurements available to us, then $f_{x|y}(\xi|\rho)$ represents the conditional density for the variables of interest, conditioned on knowledge that \mathbf{y} has assumed a particular realization, i.e., conditioned on knowledge of the numerical output of the measuring devices. Should this density be Gaussian, the conditional mean is obviously a valid choice as an estimator for \mathbf{x} . (In fact, it will be shown to be an excellent choice under more general conditions as well.) Under such circumstances, the estimator $E_{\mathbf{x}}[\mathbf{x}|\mathbf{y} = \mathbf{y}(\cdot)]$ is itself a Gaussian random variable which is a linear combination of the components of $\mathbf{y}(\cdot)$:

$$E_{\mathbf{x}}[\mathbf{x}|\mathbf{y} = \mathbf{y}(\cdot)] = \mathbf{m}_x + \mathbf{P}_{xy} \mathbf{P}_{yy}^{-1} [\mathbf{y}(\cdot) - \mathbf{m}_y] \quad (3-116)$$

Moreover, the error in this estimate, $\{\mathbf{x} - E_{\mathbf{x}}[\mathbf{x}|\mathbf{y} = \mathbf{y}(\cdot)]\}$, can be shown to be a Gaussian random variable that is *independent* of any random vector obtained by a linear transformation on \mathbf{y} ; there is no information left in the measurements \mathbf{y} that would yield better insights into the value assumed by \mathbf{x} .

3.10 LINEAR OPERATIONS ON GAUSSIAN RANDOM VARIABLES

In developing models for random phenomena and processes, we will want to perform various operations on random variables. If we allow general non-linear operations on random variables with arbitrary distributions, little can be said in general about the distribution of the transformed variables. However, if we consider linear operations on Gaussian random variables, we can claim that the Gaussian nature is preserved.

First of all, *linear transformations of Gaussian random variables are also Gaussian random variables*. If \mathbf{x} is a Gaussian random n -vector with mean \mathbf{m}_x and covariance \mathbf{P}_{xx} , and \mathbf{A} is a known $(m \times n)$ matrix (not random), then the random m -vector \mathbf{y} defined by

$$\mathbf{y} = \mathbf{A}\mathbf{x} \quad (3-117)$$

is Gaussian with mean and covariance given by

$$\mathbf{m}_y = \mathbf{A}\mathbf{m}_x \quad (3-118a)$$

$$\mathbf{P}_{yy} = \mathbf{A}\mathbf{P}_{xx}\mathbf{A}^T \quad (3-118b)$$

Proof By definition, the characteristic function for \mathbf{y} is

$$\phi_y(\boldsymbol{\mu}) = E[e^{j\boldsymbol{\mu}^T\mathbf{y}}] = E[e^{j\boldsymbol{\mu}^T\mathbf{A}\mathbf{x}}] = E[e^{j(\mathbf{A}^T\boldsymbol{\mu})^T\mathbf{x}}] = \phi_x(\mathbf{A}^T\boldsymbol{\mu})$$

where the last step is by the definition of $\phi_x(\cdot)$. Since \mathbf{x} is Gaussian, we can write $\phi_x(\mathbf{A}^T\boldsymbol{\mu})$ explicitly as

$$\phi_x(\mathbf{A}^T\boldsymbol{\mu}) = \exp\{j(\mathbf{A}^T\boldsymbol{\mu})^T\mathbf{m}_x - \frac{1}{2}(\mathbf{A}^T\boldsymbol{\mu})^T\mathbf{P}_{xx}(\mathbf{A}^T\boldsymbol{\mu})\} = \exp\{j(\boldsymbol{\mu}^T\mathbf{A}\mathbf{m}_x) - \frac{1}{2}\boldsymbol{\mu}^T\mathbf{A}\mathbf{P}_{xx}\mathbf{A}^T\boldsymbol{\mu}\}$$

Thus,

$$\phi_y(\boldsymbol{\mu}) = \exp\{j\boldsymbol{\mu}^T(\mathbf{A}\mathbf{m}_x) - \frac{1}{2}\boldsymbol{\mu}^T(\mathbf{A}\mathbf{P}_{xx}\mathbf{A}^T)\boldsymbol{\mu}\}$$

which is recognized as the characteristic function of a Gaussian random variable with mean $\mathbf{A}\mathbf{m}_x$ and covariance $\mathbf{A}\mathbf{P}_{xx}\mathbf{A}^T$.

Linear combinations of jointly Gaussian random variables are also Gaussian random variables. Note specifically that we are assuming jointly Gaussian variables here. If \mathbf{x} and \mathbf{y} are jointly Gaussian n - and m -vectors, respectively, and \mathbf{A} and \mathbf{B} are known $(p \times n)$ and $(p \times m)$ matrices, respectively, then the random p -vector \mathbf{z} defined by

$$\mathbf{z} = \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{y} \quad (3-119)$$

is Gaussian, characterized by mean and covariance

$$\mathbf{m}_z = \mathbf{A}\mathbf{m}_x + \mathbf{B}\mathbf{m}_y \quad (3-120a)$$

$$\mathbf{P}_{zz} = \mathbf{A}\mathbf{P}_{xx}\mathbf{A}^T + \mathbf{A}\mathbf{P}_{xy}\mathbf{B}^T + \mathbf{B}\mathbf{P}_{yx}\mathbf{A}^T + \mathbf{B}\mathbf{P}_{yy}\mathbf{B}^T \quad (3-120b)$$

Proof Form the $(n + m)$ -dimensional Gaussian random variable \mathbf{w}

$$\mathbf{w} = \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix}$$

which is characterized by mean and covariance

$$\mathbf{m}_w = \begin{bmatrix} \mathbf{m}_x \\ \mathbf{m}_y \end{bmatrix}, \quad \mathbf{P}_{ww} = \begin{bmatrix} \mathbf{P}_{xx} & \mathbf{P}_{xy} \\ \mathbf{P}_{yx} & \mathbf{P}_{yy} \end{bmatrix}$$

Also form the matrix $\mathbf{C} = [\mathbf{A} \mid \mathbf{B}]$. Then $\mathbf{z} = \mathbf{C}\mathbf{w}$, and the result of (3-117) and (3-118) can be invoked. ■

A useful extension of this result is that *linear combinations of jointly Gaussian random variables and nonrandom vectors are also Gaussian random variables*.

If we modify (3-119) to write

$$\mathbf{z} = \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{y} + \mathbf{c} \quad (3-121)$$

where \mathbf{c} is a known nonrandom p -vector, then \mathbf{z} is a Gaussian random p -vector, with mean and covariance

$$\mathbf{m}_z = \mathbf{A}\mathbf{m}_x + \mathbf{B}\mathbf{m}_y + \mathbf{c} \quad (3-122a)$$

$$\mathbf{P}_{zz} = \mathbf{A}\mathbf{P}_{xx}\mathbf{A}^T + \mathbf{A}\mathbf{P}_{xy}\mathbf{B}^T + \mathbf{B}\mathbf{P}_{yx}\mathbf{A}^T + \mathbf{B}\mathbf{P}_{yy}\mathbf{B}^T \quad (3-122b)$$

Proof This is easily proven by following the same steps as in the proof following (3-118), but writing

$$\phi_z(\boldsymbol{\mu}) = E[e^{j\boldsymbol{\mu}^T \mathbf{z}}] = E[e^{j\boldsymbol{\mu}^T (\mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{y} + \mathbf{c})}] = e^{j\boldsymbol{\mu}^T \mathbf{c}} E[e^{j\boldsymbol{\mu}^T (\mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{y})}]$$

The proof is then as before, but with the additional $e^{j\boldsymbol{\mu}^T \mathbf{c}}$ contributing to the mean in $\phi_z(\boldsymbol{\mu})$. ■

Note that only the mean is affected by the addition of \mathbf{c} : this makes sense since no uncertainty is contributed by the addition of a known vector. This result will be useful for adding deterministic control inputs to the dynamics model to be developed in the next chapter.

As a final point of interest, the results of (3-117)–(3-118) can be used to show that *any portion of a Gaussian random vector is itself Gaussian*, or equivalently, *if \mathbf{x} and \mathbf{y} are jointly normal, then their individual marginal densities are also Gaussian*. Let \mathbf{z} be defined as the Gaussian random variable

$$\mathbf{z} = \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} \quad (3-123)$$

where \mathbf{x} and \mathbf{y} are n - and m -dimensional partitions, respectively. Assume its mean and covariance are \mathbf{m}_z and \mathbf{P}_{zz} , partitioned as

$$\mathbf{m}_z = \begin{bmatrix} \mathbf{m}_x \\ \mathbf{m}_y \end{bmatrix}, \quad \mathbf{P}_{zz} = \begin{bmatrix} \mathbf{P}_{xx} & \mathbf{P}_{xy} \\ \mathbf{P}_{yx} & \mathbf{P}_{yy} \end{bmatrix} \quad (3-124)$$

Then \mathbf{x} is Gaussian, with mean \mathbf{m}_x and covariance \mathbf{P}_{xx} , and \mathbf{y} is Gaussian, with first moments \mathbf{m}_y and \mathbf{P}_{yy} . This was stated previously [after Eq. (3-109)], but not proven.

Proof The result is obviously true if \mathbf{x} and \mathbf{y} are independent, so that $\mathbf{P}_{xy} = \mathbf{0}$ and $\mathbf{P}_{yx} = \mathbf{0}$. However, it is also true in general, which may not be so obvious. To prove validity in the general case, let $\mathbf{A} = [\mathbf{I} \mid \mathbf{0}]$ so that

$$\mathbf{x} = \mathbf{A}\mathbf{z} = \mathbf{I}\mathbf{x} + \mathbf{0}\mathbf{y}$$

Then invoke (3-117) through (3-118) to claim \mathbf{x} is a Gaussian random variable with mean and covariance

$$\begin{aligned} E[\mathbf{x}] &= \mathbf{A}\mathbf{m}_z = \mathbf{I}\mathbf{m}_x + \mathbf{0}\mathbf{m}_y = \mathbf{m}_x \\ E[(\mathbf{x} - \mathbf{m}_x)(\mathbf{x} - \mathbf{m}_x)^T] &= \mathbf{A}\mathbf{P}_{zz}\mathbf{A}^T = \mathbf{I}\mathbf{P}_{xx}\mathbf{I}^T + \mathbf{0} = \mathbf{P}_{xx} \end{aligned}$$

and similarly for \mathbf{y} . ■

3.11 ESTIMATION WITH STATIC LINEAR GAUSSIAN SYSTEM MODELS

A general *estimation problem* can be posed in the following manner. Suppose there are some quantities of interest whose value you do not know exactly. Measuring devices can provide you with data that is functionally related to these variables, but which is also generally noise corrupted. What you would like to do is use this data, and whatever knowledge you have about its relationship to the variables of interest and about its noise corruption, to generate an estimate of the variables under consideration. Furthermore, you would like this estimate to be “optimal” in some sense, where you define the criteria for optimality.

Thus, there are five fundamental components of an estimation problem:

- (1) the *variables to be estimated*,
- (2) the *measurements* or observations available,
- (3) the *mathematical model* describing how the measurements are *related* to the variables of interest,
- (4) the *mathematical model* of the *uncertainties* present, and
- (5) the *performance evaluation criterion* to judge which estimation algorithms are “best.”

We are now able to consider the problem of estimation with static linear Gaussian system models. In other words, we are addressing ourselves to a problem which does not involve dynamics and for which linear system models and Gaussian noise models provide an adequate description. For this case, let us explicitly describe the five problem components just listed.

- (1) The variables to be estimated will be put into the form of the components of the n -dimensional vector \mathbf{x} . The true values of these quantities will remain constant, but we do not know exactly what the values are.
- (2) There will be m measurements available to us, and these will be made the components of an m -dimensional vector, \mathbf{z} .
- (3) The set of measurement data \mathbf{z} will be assumed to be a linear combination of the variables of interest, corrupted by an uncertain measurement disturbance \mathbf{v} of dimension m :

$$\mathbf{z} = \mathbf{H}\mathbf{x} + \mathbf{v} \quad (3-125)$$

where \mathbf{H} is a known $(m \times n)$ matrix.

- (4) Probabilistic models will be proposed in the form of random variables to describe the uncertainties (there are other approaches, such as unknown but bounded set descriptions of uncertainties, or “completely unknown” descriptions of disturbances). Thus, our a priori knowledge of the variables of interest can be used in describing the possible values \mathbf{x} as the realizations of a random variable \mathbf{x} , assumed to be a Gaussian random variable with mean $\hat{\mathbf{x}}^-$ and

covariance \mathbf{P}^- . (The superscript $-$ denotes a value at a time before incorporation of a measurement; $+$ will denote the corresponding value after such incorporation.)

Similarly, a random variable model is used to describe the noise corruption. We let \mathbf{v} be a Gaussian random variable, characterized by mean $\mathbf{0}$ and covariance \mathbf{R} , and assume that \mathbf{v} and \mathbf{x} are independent. Equation (3-125) can then be viewed as an equation relating the realizations of random variables: for a particular outcome ω , the realization \mathbf{v} of the random variable \mathbf{v} is added to the linear combination $\mathbf{H}\mathbf{x}$ of the realization \mathbf{x} of \mathbf{x} (the particular realization being the “true” value of the variables of interest) to generate the measurement data \mathbf{z} . This data \mathbf{z} can itself then be interpreted as the realization of a random variable, denoted as \mathbf{z} . Consequently, a random variable model would be

$$\mathbf{z} = \mathbf{H}\mathbf{x} + \mathbf{v} \quad (3-126)$$

where \mathbf{x} and \mathbf{v} are as previously described.

(5) With respect to performance criteria, we will adopt the Bayesian viewpoint that the true objective of our efforts is to generate a complete description of the probability distribution for values of the variables of interest. Since we are interested in estimating the value assumed by a continuous random variable \mathbf{x} , knowing the value of the measurement $\mathbf{z}(\omega) = \mathbf{z}$, we are thus really interested in explicitly generating the conditional density function $f_{\mathbf{x}|\mathbf{z}}(\xi|\mathbf{z})$. Once such a density function were established, it would provide all the information necessary to define an “optimal” estimate, regardless of the optimality criterion. Consider the general asymmetrical, multi peaked density $f_{\mathbf{x}|\mathbf{z}}(\xi|\mathbf{z})$ in Fig. 3.21. Reasonable definitions of an optimal estimate might include the median (having equal probability “weight” on either side), the mode (the peak or maximum likelihood value; hard to distinguish computationally from local peaks), or the mean

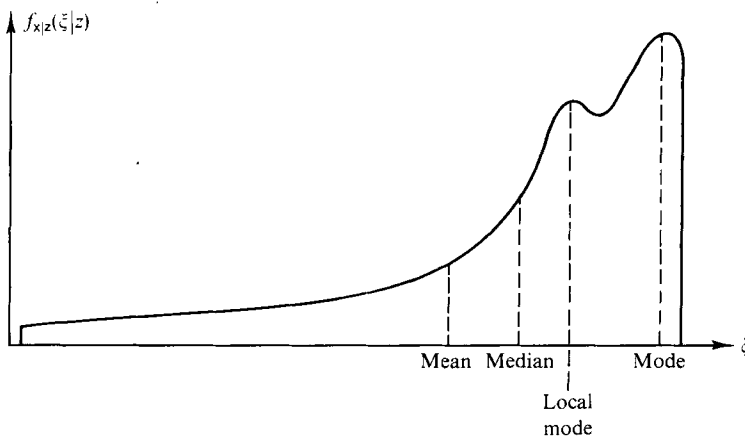


FIG. 3.21 Choice of estimator.

(the "center of probability mass" estimate). By generating the density function, some judgment can be made as to which criterion defines the most reasonable estimate for our purposes, an insight lost by first defining the criterion.

To obtain an explicit evaluation of $f_{\mathbf{x}|\mathbf{z}}(\xi|\mathbf{z})$, we want to show that it is a Gaussian density. In view of the development of (3-109)–(3-113), this entails demonstrating that \mathbf{x} and \mathbf{z} are *jointly* Gaussian random variables. First, since \mathbf{x} and \mathbf{v} are independent Gaussian random variables, they are jointly Gaussian and uncorrelated. This can be shown by reversing the steps taken in Section 3.10 to prove that uncorrelatedness implies independence for jointly Gaussian random variables, writing $f_{\mathbf{x},\mathbf{v}}(\xi, \eta)$ as

$$f_{\mathbf{x},\mathbf{v}}(\xi, \eta) = f_{\mathbf{x}}(\xi)f_{\mathbf{v}}(\eta) \quad (3-127)$$

If we define \mathbf{u} and γ as

$$\mathbf{u} = \begin{bmatrix} \mathbf{x} \\ \mathbf{v} \end{bmatrix}, \quad \gamma = \begin{bmatrix} \xi \\ \eta \end{bmatrix} \quad (3-128)$$

then $f_{\mathbf{x},\mathbf{v}}(\xi, \eta)$ can be written equivalently as $f_{\mathbf{u}}(\gamma)$, a Gaussian density described by mean \mathbf{m}_u and covariance \mathbf{P}_{uu} given by

$$\mathbf{m}_u = \begin{bmatrix} \hat{\mathbf{x}}^- \\ \mathbf{0} \end{bmatrix}, \quad \mathbf{P}_{uu} = \begin{bmatrix} \mathbf{P}^- & \mathbf{0} \\ \mathbf{0} & \mathbf{R} \end{bmatrix} \quad (3-129)$$

So far we have shown \mathbf{u} to be Gaussian. But linear transformations of Gaussian random variables are themselves Gaussian, so \mathbf{w} defined by

$$\mathbf{w} = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{H} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \mathbf{v} \end{bmatrix} = \begin{bmatrix} \mathbf{x} \\ \mathbf{H}\mathbf{x} + \mathbf{v} \end{bmatrix} = \begin{bmatrix} \mathbf{x} \\ \mathbf{z} \end{bmatrix} \quad (3-130)$$

is a Gaussian random variable with mean and covariance given by (3-118) as

$$\mathbf{m}_w = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{H} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \hat{\mathbf{x}}^- \\ \mathbf{0} \end{bmatrix} = \begin{bmatrix} \hat{\mathbf{x}}^- \\ \mathbf{H}\hat{\mathbf{x}}^- \end{bmatrix} \quad (3-131a)$$

$$\begin{aligned} \mathbf{P}_{ww} &= \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{H} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{P}^- & \mathbf{0} \\ \mathbf{0} & \mathbf{R} \end{bmatrix} \begin{bmatrix} \mathbf{I} & \mathbf{H}^T \\ \mathbf{0} & \mathbf{I} \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{P}^- & \mathbf{P}^- \mathbf{H}^T \\ \mathbf{H} \mathbf{P}^- & \mathbf{H} \mathbf{P}^- \mathbf{H}^T + \mathbf{R} \end{bmatrix} \end{aligned} \quad (3-131b)$$

Thus, \mathbf{x} and \mathbf{z} are jointly Gaussian random variables, of dimensions n and m , respectively, with their joint density $f_{\mathbf{x},\mathbf{z}}(\xi, \zeta)$ a Gaussian density characterized by \mathbf{m}_w and \mathbf{P}_{ww} .

At this point, we can say that $f_{\mathbf{x}|\mathbf{z}}(\xi|\zeta)$ is a Gaussian conditional density function and define it completely through its mean and covariance. Recall the

result of Eqs. (3-109)–(3-113), and make the following replacements:

Random variable: $\mathbf{y} \rightarrow \mathbf{z}$

Realization: $\mathbf{y} \rightarrow \mathbf{z}$

Dummy variable: $\rho \rightarrow \zeta$

Mean of joint density: $\begin{bmatrix} \mathbf{m}_x \\ \mathbf{m}_y \end{bmatrix} \rightarrow \begin{bmatrix} \hat{\mathbf{x}}^- \\ \mathbf{H}\hat{\mathbf{x}}^- \end{bmatrix}$

Covariance of joint density: $\begin{bmatrix} \mathbf{P}_{xx} & \mathbf{P}_{xy} \\ \mathbf{P}_{yx} & \mathbf{P}_{yy} \end{bmatrix} \rightarrow \begin{bmatrix} \mathbf{P}^- & \mathbf{P}^- \mathbf{H}^T \\ \mathbf{H} \mathbf{P}^- & \mathbf{H} \mathbf{P}^- \mathbf{H}^T + \mathbf{R} \end{bmatrix}$

Then the conditional mean, denoted as $\hat{\mathbf{x}}^+$, is given by (3-113) as

$$\hat{\mathbf{x}}^+ = E_{\mathbf{x}}[\mathbf{x} | \mathbf{z} = \mathbf{z}] = \hat{\mathbf{x}}^- + [\mathbf{P}^- \mathbf{H}^T][\mathbf{H} \mathbf{P}^- \mathbf{H}^T + \mathbf{R}]^{-1}[\mathbf{z} - \mathbf{H}\hat{\mathbf{x}}^-] \quad (3-132)$$

Similarly (3-114) yields the conditional covariance, denoted by \mathbf{P}^+ , as

$$\mathbf{P}^+ = \mathbf{P}^- - [\mathbf{P}^- \mathbf{H}^T][\mathbf{H} \mathbf{P}^- \mathbf{H}^T + \mathbf{R}]^{-1}[\mathbf{H} \mathbf{P}^-] \quad (3-133)$$

Note that if we define the gain matrix \mathbf{K} as

$$\mathbf{K} = \mathbf{P}^- \mathbf{H}^T [\mathbf{H} \mathbf{P}^- \mathbf{H}^T + \mathbf{R}]^{-1} \quad (3-134)$$

then (3-132) and (3-133) can be written as

$$\hat{\mathbf{x}}^+ = \hat{\mathbf{x}}^- + \mathbf{K}[\mathbf{z} - \mathbf{H}\hat{\mathbf{x}}^-] \quad (3-135)$$

$$\mathbf{P}^+ = \mathbf{P}^- - \mathbf{K} \mathbf{H} \mathbf{P}^- \quad (3-136)$$

Since $\hat{\mathbf{x}}^+$ is the mean of the symmetric Gaussian conditional density $f_{\mathbf{x}|\mathbf{z}}(\boldsymbol{\xi}|\mathbf{z})$, it is also the mode. Consequently, we choose it as an optimal estimate of the variables of interest. As discussed at the end of Section 3.5, (3-135) is an equation for a vector $\hat{\mathbf{x}}^+(\mathbf{z})$ in R^n ; the mapping $\hat{\mathbf{x}}^+(\cdot)$ from R^m into R^n

$$\hat{\mathbf{x}}^+(\cdot) = \hat{\mathbf{x}}^- + \mathbf{K}[\cdot - \mathbf{H}\hat{\mathbf{x}}^-] \quad (3-137)$$

is an estimator, and the composite mapping $\hat{\mathbf{x}}^+[\mathbf{z}(\cdot)]$ is a random variable

$$\hat{\mathbf{x}}^+ = \hat{\mathbf{x}}^+[\mathbf{z}(\cdot)] = \hat{\mathbf{x}}^- + \mathbf{K}[\mathbf{z}(\cdot) - \mathbf{H}\hat{\mathbf{x}}^-] \quad (3-138)$$

By choosing $\hat{\mathbf{x}}^+$ as an estimate of \mathbf{x} , the vector $[\mathbf{x} - \hat{\mathbf{x}}^+]$ is a Gaussian random variable that describes the error in the estimate, denoted as

$$\mathbf{e} = \mathbf{x} - \hat{\mathbf{x}}^+ \quad (3-139)$$

The conditional mean of \mathbf{e} is zero, and the conditional covariance (see Problem 3.20) is

$$E_{\mathbf{x}}[\mathbf{e}\mathbf{e}^T | \mathbf{z} = \mathbf{z}] = E_{\mathbf{x}}[(\mathbf{x} - \hat{\mathbf{x}}^+)(\mathbf{x} - \hat{\mathbf{x}}^+)^T | \mathbf{z} = \mathbf{z}] = \mathbf{P}^+ \quad (3-140)$$

Thus, if we choose $\hat{\mathbf{x}}^+$ as an estimate of \mathbf{x} , the \mathbf{P}^+ calculated through (3-136) assumes additional significance: it is the covariance to describe the Gaussian error committed by the estimate. Note that this covariance matrix can be computed *without* knowledge of the actual measurement realization, $\mathbf{z}(\omega) = \mathbf{z}$. Consequently, both \mathbf{P}^+ and the gain matrix \mathbf{K} can be *precomputed*.

Equations (3-134)–(3-136) can be written in the algebraically equivalent form of

$$\hat{\mathbf{x}}^+ = [\mathbf{P}^+(\mathbf{P}^-)^{-1}] \hat{\mathbf{x}}^- + [\mathbf{P}^+ \mathbf{H}^T \mathbf{R}^{-1}] \mathbf{z} \quad (3-141)$$

$$\mathbf{P}^+ = [(\mathbf{P}^-)^{-1} + \mathbf{H}^T \mathbf{R}^{-1} \mathbf{H}]^{-1} \quad (3-142)$$

These expressions involve $(n \times n)$ matrix inversions rather than $(m \times m)$ inversions as in the previous equations, so are attractive computationally only if $m > n$; this will be developed further in Sections 5.7 and 7.8. However, this equivalent set of expressions will be more readily manipulated for the case of little or no a priori state information, as will be seen in examples to follow.

EXAMPLE 3.19 Recall the scalar example of the two simultaneous star sightings discussed in Section 1.5. There, x was the one-dimensional position, and z was the location measured by means of the star sightings, modeled as

$$z_1 = x + v_1, \quad z_2 = x + v_2$$

We assumed that we had no a priori information about x , that v_1 and v_2 could be modeled as zero-mean Gaussian random variables with variances $\sigma_{z_1}^2$ and $\sigma_{z_2}^2$, respectively, and that x , v_1 , and v_2 were independent random variables.

One means of solving for the best estimate of position would be to consider $z_1(\omega) = z_1$ and the variance $\sigma_{z_1}^2$ to provide the “a priori” information about x before the second measurement is taken. This was the approach taken in Chapter 1: a sequential, or recursive, estimation procedure. Thus, we use this a priori knowledge to describe the random variable x as a Gaussian random variable with mean z_1 and variance $\sigma_{z_1}^2$ (“ $\hat{\mathbf{x}}^- = z_1$,” “ $\mathbf{P}^- = \sigma_{z_1}^2$ ”). Consequently, we consider $z_2(\omega) = z_2$ as the “available measurement” to be incorporated into the estimate of position. Since we model z_2 as $(x + v_2)$ with v_2 zero-mean, Gaussian, and with variance $\sigma_{z_2}^2$, we have “ $\mathbf{R} = \sigma_{z_2}^2$.”

The optimal estimate is then the mean (and mode) of the conditional density $f_{x|z_1, z_2}(\xi|z_1, z_2)$:

$$\begin{aligned} \hat{x}^+ &= \hat{x}^- + \mathbf{P}^- \mathbf{H}^T [\mathbf{H} \mathbf{P}^- \mathbf{H}^T + \mathbf{R}]^{-1} (z_2 - \mathbf{H} \hat{x}^-) \\ &= z_1 + \sigma_{z_1}^2 [\sigma_{z_1}^2 + \sigma_{z_2}^2]^{-1} (z_2 - z_1) \end{aligned}$$

which is, in fact, the result obtained in Eq. (1-6).

The error variance associated with using \hat{x}^+ as an estimate, as generated by the estimation algorithm itself, is then \mathbf{P}^+ :

$$\begin{aligned} \mathbf{P}^+ &= \mathbf{P}^- - \mathbf{P}^- \mathbf{H}^T [\mathbf{H} \mathbf{P}^- \mathbf{H}^T + \mathbf{R}]^{-1} \mathbf{H} \mathbf{P}^- \\ &= \sigma_{z_1}^2 - \sigma_{z_1}^2 [\sigma_{z_1}^2 + \sigma_{z_2}^2]^{-1} \sigma_{z_1}^2 \end{aligned}$$

which was also the result obtained previously, Eq. (1-9). ■

EXAMPLE 3.20 Another means of solving for the best estimate of position in the previous example would be to assume no a priori information about x , and to incorporate the two measurements simultaneously, i.e., in a batch. If there is no a priori information, we could model this through a Gaussian random variable with infinite variance, $\mathbf{P}^- = \infty$, or equivalently, $(\mathbf{P}^-)^{-1} = 0$.

The measurement is

$$\mathbf{z} = \begin{bmatrix} z_1 \\ z_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \mathbf{x} + \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = \mathbf{H}\mathbf{x} + \mathbf{v}$$

where \mathbf{v} is modeled as a zero-mean Gaussian noise of covariance \mathbf{R} :

$$\mathbf{R} = E\{\mathbf{v}\mathbf{v}^T\} = \begin{bmatrix} E\{v_1^2\} & E\{v_1 v_2\} \\ E\{v_1 v_2\} & E\{v_2^2\} \end{bmatrix} = \begin{bmatrix} \sigma_{z_1}^2 & 0 \\ 0 & \sigma_{z_2}^2 \end{bmatrix}$$

where the off-diagonal zeros are due to v_1 and v_2 being independent and thus uncorrelated.

Now, using (3-142), we can write P^+ as

$$\begin{aligned} P^+ &= [(P^-)^{-1} + \mathbf{H}^T \mathbf{R}^{-1} \mathbf{H}]^{-1} = [\mathbf{H}^T \mathbf{R}^{-1} \mathbf{H}]^{-1} \\ &= \left\{ \begin{bmatrix} 1 & 1 \end{bmatrix} \begin{bmatrix} 1/\sigma_{z_1}^2 & 0 \\ 0 & 1/\sigma_{z_2}^2 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} \right\}^{-1} = \frac{1}{[(1/\sigma_{z_1}^2) + (1/\sigma_{z_2}^2)]} \end{aligned}$$

which is identical to the result of the previous example, and equivalent to the result of Eq. (1-4).

The state estimate can be written using (3-141):

$$\begin{aligned} \hat{\mathbf{x}}^+ &= [P^+(P^-)^{-1}] \hat{\mathbf{x}}^- + [P^+ \mathbf{H}^T \mathbf{R}^{-1}] \mathbf{z} \\ &= P^+ \mathbf{H}^T \mathbf{R}^{-1} \mathbf{z} = [\mathbf{H}^T \mathbf{R}^{-1} \mathbf{H}]^{-1} \mathbf{H}^T \mathbf{R}^{-1} \mathbf{z} \\ &= \frac{\sigma_{z_1}^2 \sigma_{z_2}^2}{\sigma_{z_1}^2 + \sigma_{z_2}^2} \begin{bmatrix} 1 & 1 \end{bmatrix} \begin{bmatrix} 1/\sigma_{z_1}^2 & 0 \\ 0 & 1/\sigma_{z_2}^2 \end{bmatrix} \begin{bmatrix} z_1 \\ z_2 \end{bmatrix} \\ &= \frac{\sigma_{z_2}^2}{\sigma_{z_1}^2 + \sigma_{z_2}^2} z_1 + \frac{\sigma_{z_1}^2}{\sigma_{z_1}^2 + \sigma_{z_2}^2} z_2 \end{aligned}$$

Again this is identical to both the previous result and that of Chapter 1. ■

The previous examples demonstrated two methods of processing measurements. In *batch* processing, \mathbf{z} is the vector of all measurements that are available, and thus all measurements are simultaneously incorporated into the estimate. For *recursive* processing, \mathbf{z} is partitioned into components:

$$\mathbf{z} = \begin{bmatrix} \mathbf{z}_1 \\ \mathbf{z}_2 \\ \vdots \\ \mathbf{z}_K \end{bmatrix} \quad (3-143)$$

First, the estimate $\hat{\mathbf{x}}^+$ and covariance \mathbf{P}^+ based upon $\mathbf{z}_1(\omega) = \mathbf{z}_1$ alone are computed. Then the Gaussian random variable so obtained, with mean $\hat{\mathbf{x}}^+$ and covariance \mathbf{P}^+ , is considered to be the information available about \mathbf{x} *prior to the next measurement*, $\mathbf{z}_2(\omega) = \mathbf{z}_2$. The update process is then repeated until all partitions of $\mathbf{z}(\omega) = \mathbf{z}$ are incorporated.

As illustrated by the previous simple examples, if \mathbf{R} is an $(m \times m)$ diagonal matrix, the batch processing of the m -dimensional measurement realization \mathbf{z} and the recursive processing of the m scalar measurements z_1, z_2, \dots, z_m yield *equivalent* results. To generalize this statement, let \mathbf{R} be block diagonal and

let \mathbf{z} be partitioned corresponding to the diagonal blocks of \mathbf{R} :

$$\mathbf{R} = \begin{bmatrix} \mathbf{R}_1 & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \mathbf{R}_2 & \cdots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{R}_K \end{bmatrix}, \quad \mathbf{z} = \begin{bmatrix} \mathbf{z}_1 \\ \mathbf{z}_2 \\ \vdots \\ \mathbf{z}_K \end{bmatrix} \quad (3-144)$$

Then batch processing of \mathbf{z} and recursive processing of $\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_K$ will yield identical results (see Maybeck [6] and Problem 3.17 for proof). Chapter 7 will extend these results by explicitly generating a transformation of variables to convert any system model into an equivalent form, but with a diagonal \mathbf{R} so that m scalar updates can *always* be used.

This equivalence can be exploited in the design of online estimators. First of all, the recursive form entails the inversion of smaller dimensioned matrices, yielding simpler algorithms. In addition, online estimators are often implemented in general purpose computers, so that only a certain time is allotted to the algorithm, determined in part by the number of high priority interrupts received by the computer to perform other functions. Thus, there *may* not be sufficient time to perform a single batch processing of \mathbf{z} in a given period if the computer is heavily loaded. However, there would be time to process at least *some* of the partitions \mathbf{z}_1 to \mathbf{z}_K . Since a partially updated estimate would be preferable to one not updated at all, the recursive form might be a substantially better implementation.

The estimation result of Eqs. (3-134)–(3-136) or (3-141)–(3-142) can be directly related to *weighted least squares estimation*. Least squares estimation is a classical technique used extensively, especially in curve fitting applications in which it is desired to obtain the polynomial of given order (or some other chosen functional form) that “best” fits a set of data points. “Best” is defined in terms of minimizing the sum of squares of the differences between the actual measurement data and the proposed, or estimated, function or curve. If one wants to match certain data points more closely than others, a weighting coefficient can be assigned to each term in the sum to be minimized, more heavily weighting the “cost” of differing from the more critical points, yielding what is termed weighted least squares estimation.

Suppose that the measurement model is

$$\mathbf{z} = \mathbf{H}\mathbf{x} + \mathbf{v} \quad (3-145)$$

where \mathbf{v} is an m -vector of measurement noise, whose statistical characteristics are *not* defined. We then want to use our knowledge of the measured value \mathbf{z} to generate an estimate $\hat{\mathbf{x}}$ of the unknown \mathbf{x} . Thus, we want to find the value of $\hat{\mathbf{x}}$ that minimizes the weighted sum of squares of the m components of the vector $[\mathbf{z} - \mathbf{H}\hat{\mathbf{x}}]$. If we let \mathbf{W} be a general $(m \times m)$ weighting matrix, then we

want to find the vector $\hat{\mathbf{x}}$ that minimizes the scalar cost J :

$$J = \frac{1}{2} [\mathbf{z} - \mathbf{H}\hat{\mathbf{x}}]^T \mathbf{W} [\mathbf{z} - \mathbf{H}\hat{\mathbf{x}}] \quad (3-146)$$

Note that if $\mathbf{W} = \mathbf{I}$, this is standard least squares, with

$$J = \frac{1}{2} \sum_{i=1}^m [z - H\hat{x}]_i^2 \quad (3-147)$$

If \mathbf{W} is a diagonal matrix with diagonal terms w_1, w_2, \dots, w_m , then

$$J = \frac{1}{2} \sum_{i=1}^m w_i [z - H\hat{x}]_i^2 \quad (3-148)$$

This minimization is accomplished by the value $\hat{\mathbf{x}}_{\text{WLS}}$, where the subscript denotes weighted least squares, if

$$\left. \frac{\partial J}{\partial \hat{\mathbf{x}}} \right|_{\hat{\mathbf{x}} = \hat{\mathbf{x}}_{\text{WLS}}} \triangleq \left[\frac{\partial J}{\partial \hat{x}_1} \frac{\partial J}{\partial \hat{x}_2} \dots \frac{\partial J}{\partial \hat{x}_n} \right] \bigg|_{\hat{\mathbf{x}} = \hat{\mathbf{x}}_{\text{WLS}}} = \mathbf{0}^T \quad (3-149a)$$

and

$$\left. \frac{\partial^2 J}{\partial \hat{\mathbf{x}}^2} \right|_{\hat{\mathbf{x}} = \hat{\mathbf{x}}_{\text{WLS}}} \geq \mathbf{0} \quad (3-149b)$$

i.e., the second derivative matrix is positive semidefinite. Performing the indicated differentiation on (3-146) yields

$$-[\mathbf{z} - \mathbf{H}\hat{\mathbf{x}}]^T \mathbf{W} \mathbf{H} \big|_{\hat{\mathbf{x}} = \hat{\mathbf{x}}_{\text{WLS}}} = \mathbf{0}^T$$

or, $\hat{\mathbf{x}}_{\text{WLS}}$ is the vector that satisfies

$$\mathbf{H}^T \mathbf{W} \mathbf{H} \hat{\mathbf{x}}_{\text{WLS}} = \mathbf{H}^T \mathbf{W} \mathbf{z} \quad (3-150)$$

if $[\mathbf{H}^T \mathbf{W} \mathbf{H}]$ is positive semidefinite. If $[\mathbf{H}^T \mathbf{W} \mathbf{H}]$ is in fact positive definite, and thus has a unique inverse, then

$$\hat{\mathbf{x}}_{\text{WLS}} = [\mathbf{H}^T \mathbf{W} \mathbf{H}]^{-1} \mathbf{H}^T \mathbf{W} \mathbf{z} \quad (3-151)$$

This can be compared to the result obtained in Example 3.20 for the case of no a priori information about \mathbf{x} , i.e., letting $(\mathbf{P}^-)^{-1} = \mathbf{0}$:

$$\hat{\mathbf{x}}^+ = [\mathbf{H}^T \mathbf{R}^{-1} \mathbf{H}]^{-1} \mathbf{H}^T \mathbf{R}^{-1} \mathbf{z} \quad (3-152)$$

The two results are identical if we choose \mathbf{W} to be \mathbf{R}^{-1} (positive definite). However, least squares theory gives *no* insights into such a choice of weighting matrix, since no statistical characterization (as \mathbf{v} being Gaussian, zero-mean, of covariance \mathbf{R}) was assumed to be known. Analogous to the previous discussion, if \mathbf{W} is block diagonal, then an equivalent result can be achieved by recursive least squares estimation, the form of which is developed in Problem 3.18.

3.12 SUMMARY

This chapter has presented basic concepts of probability theory in a progression that is logical for addressing the problem of estimating some quantities of interest based upon noise-corrupted measurements of related variables. From a Bayesian point of view, this problem is solved by establishing a complete description of the conditional density function of the random vector \mathbf{x} modeling the quantities of interest, conditioned on knowledge of the measured values $\mathbf{z}(\omega_i) = \mathbf{z}$ available: $f_{\mathbf{x}|\mathbf{z}}(\xi|\mathbf{z})$. Thus, it was necessary to develop the concepts of random variable models as real-valued functions or mappings, and descriptions of random variables through associated probability distribution and density functions (assuming the existence of the latter). Both unconditioned and conditioned probability functions were discussed, and conditioning allowed the observed realizations of one random variable to provide information about the possible realizations of another, related variable.

The expected value of some function of a random variable is simply the ensemble average value of that function, as the random variable assumes all of its possible realizations. Expectations of particular functions, called moments of a random variable, provided in general a partial description of that random variable. Conditional expectations and moments are of special significance to estimation since it is considerably more feasible to generate and implement algorithms to compute conditional moments than those intended to construct the entire description explicitly as $F_{\mathbf{x}|\mathbf{z}}(\xi|\mathbf{z})$ or $f_{\mathbf{x}|\mathbf{z}}(\xi|\mathbf{z})$.

In the special case in which $f_{\mathbf{x}|\mathbf{z}}(\xi|\mathbf{z})$ is Gaussian, numerical computation of the first two moments, the mean and covariance, provides a complete depiction of the density function rather than just a partial description. Thus, a computationally feasible estimation algorithm can be developed that satisfies the Bayesian objective of portraying this conditional density. Because linear operations on Gaussian random vectors again yield Gaussian random vectors, the class of problems to which such an algorithm is directly applicable is rather large. This chapter concluded with the detailed development of such an algorithm for estimation with static linear Gaussian system models.

The following chapter will extend these concepts to the case in which quantities of interest can undergo dynamic changes in time. The fundamental ideas of probability theory will be instrumental in developing not only such stochastic process models, but also the estimation and control algorithms that will later exploit these system models.

REFERENCES

1. Bury, K. V., *Statistical Models in Applied Science*. Wiley, New York, 1975.
2. Cramér, H., *Mathematical Methods of Statistics*. Princeton Univ. Press, Princeton, New Jersey, 1966.
3. Davenport, W. B. Jr., *Probability and Random Processes*. McGraw-Hill, New York, 1970.
4. Feller, W., *An Introduction to Probability and Its Applications*, Vols I and II. Wiley, New York, 1950 (Vol. 1), 1966 (Vol. 2).

5. Loeve, M., *Probability Theory*. Van Nostrand-Reinhold, Princeton, New Jersey, 1963.
6. Maybeck, P. S., "Combined Estimation of States and Parameters for On-Line Applications," Ph.D. dissertation, M.I.T., Cambridge, Massachusetts, February 1972.
7. McGarty, T. P., *Stochastic Systems and State Estimation*. Wiley, New York, 1974.
8. Mendenhall, W., and Schaeffer, R. L., *Mathematical Statistics with Applications*. Duxbury Press, North Scituate, Massachusetts, 1973.
9. Papoulis, A., *Probability, Random Variables, and Stochastic Processes*. McGraw-Hill, New York, 1965.
10. Parzen, E., *Modern Probability Theory and Its Applications*. Wiley, New York, 1960.
11. Rao, C. R., *Linear Statistical Inference and Its Applications*. Wiley, New York, 1965.
12. Royden, H. L., *Real Analysis*. Macmillan, New York, 1968.
13. Rudin, W., *Principles of Mathematical Analysis*. McGraw-Hill, New York, 1964.
14. Wilks, S. S., *Mathematical Statistics*. Wiley, New York, 1962.
15. Wong, E., *Stochastic Processes in Information and Dynamical Systems*. McGraw-Hill, New York, 1971.

PROBLEMS

3.1 Prove that for any set $A \subset \Omega$, $A \in \mathcal{F}$, the probability $P(A)$ is bounded as $0 \leq P(A) \leq 1$.

3.2 Consider two tosses of a fair coin. Completely define the appropriate probability space $\{\Omega, \mathcal{F}, P\}$. Let us say that we are interested only in the number of heads in the two tosses: can a different probability space $\{\Omega, \mathcal{F}^1, P^1\}$ be defined with \mathcal{F}^1 a smaller collection of sets? Define an appropriate random variable $x(\cdot)$ to consider the number of heads appearing in two tosses. Obtain the probability distribution function for this $x(\cdot)$.

3.3 If the joint probability density of x_1 and x_2 is

$$f_{x_1, x_2}(\xi_1, \xi_2) = \begin{cases} \frac{1}{2\pi} \exp[-\frac{1}{2}(\xi_1 - 3)^2 - \xi_2] & \xi_2 > 0 \\ 0 & \xi_2 \leq 0 \end{cases}$$

what is the characteristic function for a random variable y , where

$$y = x_1 + x_2 ?$$

Determine the mean of y .

3.4 By definition, the i th component of the n -dimensional mean vector is

$$m_i = \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} \xi_i f_n(\xi) d\xi_1 \cdots d\xi_n$$

Is this the same as

$$m_i = \int_{-\infty}^{\infty} \xi_i f_{x_i}(\xi_i) d\xi_i ?$$

Show why.

3.5 Let $x(\cdot)$ and $y(\cdot)$ be independent random variables that are each uniformly distributed on the interval $[0, 1]$. Define the random variable $z(\cdot)$ as

$$z(\cdot) = x(\cdot)y(\cdot)$$

- (a) What are the mean, mean squared value, and variance of $x(\cdot)$ or $y(\cdot)$? For $z(\cdot)$?
- (b) What is the probability that $z(\cdot)$ assumes a value less than 0.5? Less than or equal to 0.5?

3.6 Prove that the random vector $\{\mathbf{x} - E_{\mathbf{x}}[\mathbf{x}|\mathbf{y} = \mathbf{y}(\cdot)]\}$ is orthogonal to the random vector \mathbf{y} : that

$$E_{\mathbf{xy}}\{\mathbf{y}(\mathbf{x} - E_{\mathbf{x}}[\mathbf{x}|\mathbf{y} = \mathbf{y}(\cdot)])^T\} = \mathbf{0}$$

Show that this can be generalized—that $\{\mathbf{x} - E_{\mathbf{x}}[\mathbf{x}|\mathbf{y} = \mathbf{y}(\cdot)]\}$ is orthogonal to any function of \mathbf{y} . This concept is instrumental in deriving the Kalman filter by means of “orthogonal projections,” which was the original means of derivation.

3.7 At the end of Section 3.6, it was stated that if either \mathbf{x} or \mathbf{y} (or both) is zero-mean, then orthogonality and uncorrelatedness of \mathbf{x} and \mathbf{y} imply each other. Prove this. Also prove that, if neither \mathbf{x} nor \mathbf{y} is zero-mean, then they cannot be both uncorrelated and orthogonal.

3.8 Let \mathbf{x} and \mathbf{y} be random n -vectors with $\mathbf{y} = \boldsymbol{\theta}(\mathbf{x})$. Suppose $\boldsymbol{\theta}^{-1}$ exists and that both $\boldsymbol{\theta}$ and $\boldsymbol{\theta}^{-1}$ are continuously differentiable. Then

$$f_{\mathbf{y}}(\boldsymbol{\rho}) = f_{\mathbf{x}}[\boldsymbol{\theta}^{-1}(\boldsymbol{\rho})] \|\partial\boldsymbol{\theta}^{-1}(\boldsymbol{\rho})/\partial\boldsymbol{\rho}\|$$

where $\|\partial\boldsymbol{\theta}^{-1}(\boldsymbol{\rho})/\partial\boldsymbol{\rho}\| > 0$ is the absolute value of the Jacobian determinant. Prove this theorem using the conditional density relationship

$$f_{\mathbf{y}|\mathbf{x}}(\boldsymbol{\rho}|\boldsymbol{\xi}) = f_{\mathbf{x},\mathbf{y}}(\boldsymbol{\xi}, \boldsymbol{\rho})/f_{\mathbf{x}}(\boldsymbol{\xi})$$

as a beginning. Write $f_{\mathbf{y}}(\boldsymbol{\rho})$ in terms of $f_{\mathbf{x},\mathbf{y}}(\boldsymbol{\xi}, \boldsymbol{\rho})$ and continue the proof.

3.9 The scalar Gaussian random variable $x(\cdot)$ is defined on $\Omega = R^1$ (i.e., the sample space is the real line). The statistics of $x(\cdot)$ are

$$E\{x\} = m, \quad E\{[x - m]^2\} = P$$

The scalar random variables $y(\cdot)$ and $z(\cdot)$ are defined by

$$y(\cdot) = x^5(\cdot), \quad z(\cdot) = x^2(\cdot)$$

(a) Find the probability density for $y(\cdot)$ by using fundamental set definitions to establish $F_y(\rho)$ and then find its derivative.

(b) Find $f_y(\rho)$ by a method analogous to the proof of the preceding problem.

(c) Find $f_y(\rho)$ by direct application of the result of the last problem.

(d) Find the probability density for $z(\cdot)$ by the method used in part (a). Why would the methods of (b) and (c) not be directly applicable?

3.10 Let x , y , and z be pairwise independent. Show that they need not be triplewise independent. What if $[x, y, z]^T$ is a Gaussian random vector?

3.11 Consider a three-dimensional Gaussian random vector, $\mathbf{x}(\cdot)$, one whose probability density is described by

$$f_{\mathbf{x}}(\boldsymbol{\xi}) = [(2\pi)^{3/2}|\mathbf{P}|^{1/2}]^{-1} \exp\{-\frac{1}{2}[\boldsymbol{\xi} - \mathbf{m}]^T \mathbf{P}^{-1}[\boldsymbol{\xi} - \mathbf{m}]\}$$

where the mean \mathbf{m} and covariance \mathbf{P} are

$$\mathbf{m} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}, \quad \mathbf{P} = \begin{bmatrix} 9 & 0 & 0 \\ 0 & 2.5 & 0.5 \\ 0 & 0.5 & 2.5 \end{bmatrix}$$

Surfaces of constant probability density are called surfaces of constant likelihood. They are ellipsoids with principal axes not generally aligned with the coordinate axes.

(a) Determine a transformation of variables $\mathbf{x}' = \mathbf{T}\mathbf{x}$ so that it is possible to use the principal axes of the ellipsoid as the coordinate axes. When this is done, \mathbf{P}' becomes diagonal, i.e.,

$$\mathbf{P}' = \begin{bmatrix} \sigma_{11}'^2 & 0 & 0 \\ 0 & \sigma_{22}'^2 & 0 \\ 0 & 0 & \sigma_{33}'^2 \end{bmatrix}$$

Obtain this form for the given matrix \mathbf{P} .

(b) Show that now the surface of constant likelihood is an ellipsoid of the form

$$(\xi_1'^2/\sigma_{11}'^2) + (\xi_2'^2/\sigma_{22}'^2) + (\xi_3'^2/\sigma_{33}'^2) = c^2$$

Write an expression for the probability that x_1 , x_2 , and x_3 take values within the ellipsoid.

(c) Show that our ellipsoid becomes a sphere by defining new variables

$$x_1'' = x_1'/\sigma_{11}', \quad x_2'' = x_2'/\sigma_{22}', \quad x_3'' = x_3'/\sigma_{33}'$$

and that the probability can be written as a volume integral over the ellipsoid:

$$\text{Prob}\{(x_1, x_2, x_3) \text{ lies within ellipsoid}\} = \iiint \frac{e^{-r^2/2}}{(2\pi)^{3/2}} d\xi_1'' d\xi_2'' d\xi_3''$$

where $r^2 = \xi_1''^2 + \xi_2''^2 + \xi_3''^2$, or in another form as

$$\text{Prob}\{(x_1, x_2, x_3) \text{ lies within ellipsoid}\} = \int_0^c \frac{s(r)e^{-r^2/2}}{(2\pi)^{3/2}} dr$$

where $s(r)$ is the surface area of a sphere of radius r .

(d) Calculate the probability for $c = 1$ and $c = 2$.

3.12 Prove that for a zero-mean Gaussian random vector \mathbf{x} , with covariance \mathbf{P} ,

$$E[x_k x_l x_m x_n] = P_{kl}P_{mn} + P_{km}P_{ln} + P_{kn}P_{lm}$$

3.13 At the end of Section 3.9, it was claimed that the error $\{\mathbf{x} - E_{\mathbf{x}}[\mathbf{x}|\mathbf{y} = \mathbf{y}(\cdot)]\}$ is a Gaussian random vector that is independent of any random vector obtained as a linear transformation on \mathbf{y} , under the assumptions made in that section. Prove this.

3.14 A parameter x is to be estimated on the basis of a priori information and a single noisy measurement. The quality of the a priori information is expressed by the probability density function in Fig. 3.P1. The measurement is assumed to be of the form

$$z = x + n$$

where n is a noise, independent of x , which has a probability density of the form given in Fig. 3.P2. The actual measurement taken had the value of $\frac{1}{2}$. Find the conditional probability density $f_{x|z}(\xi|\frac{1}{2})$ for $\xi = \frac{1}{2}$, i.e., $f_{x|z}(\xi|\frac{1}{2})$. Plot this density as a function of ξ .

One reasonable estimate of x would be the value of ξ that maximizes the density $f_{x|z}(\xi|\frac{1}{2})$. This is a "maximum likelihood estimate," and we will denote it here as \hat{x}_{ML} . (It is also called the maximum a posteriori, or MAP, estimate; see Section 5.5.) Find its value.

Another reasonable estimate of x would be the conditional mean, $E_{\mathbf{x}}[x|z = \frac{1}{2}]$, which we will denote as \hat{x} . Find its value.

Now determine some statistical information about the error committed by these two estimates. Define the error in the maximum likelihood estimate as

$$e_{\text{ML}} = \hat{x}_{\text{ML}} - x$$

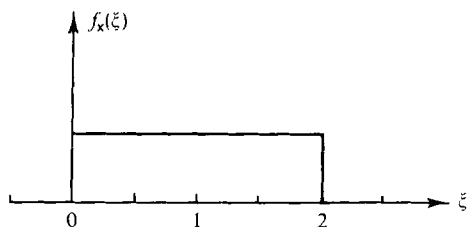


FIG. 3.P1 A priori information for Problem 3.14.

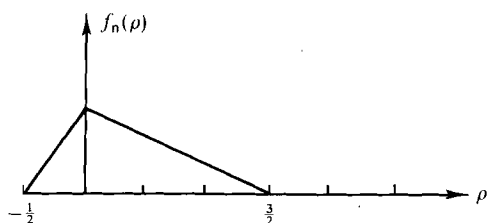


FIG. 3.P2 Measurement noise description for Problem 3.14.

Obtain the conditional mean and conditional variance of this error, conditioned on the fact that $z = \xi = \frac{1}{2}$. Similarly define the error in the conditional mean estimate of x , and obtain the conditional mean and variance of this error. The conditional mean can be shown to be the estimator, out of the class of linear estimators with zero-mean error, that has minimum error variance: this does *not* mean that other estimators cannot *duplicate* this error variance (as in this problem), or that estimators *outside* of the class under consideration cannot *outperform* the conditional mean.

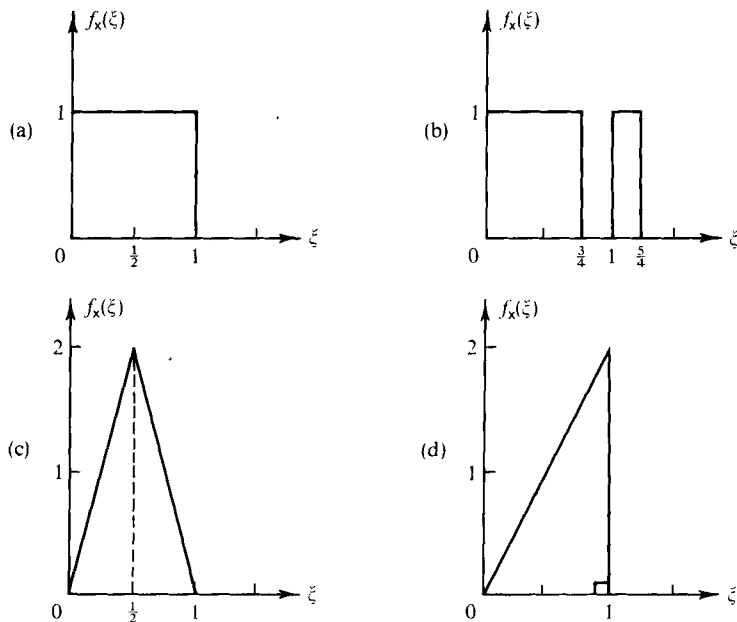


FIG. 3.P3 Density functions for Problem 3.15.

3.15 This problem demonstrates the differences among mean, mode, and median estimates. Find these three estimates for the four density functions depicted in Fig. 3.P3.

3.16 (a) Generate the matrix inverse indicated in Eq. (3-109) by letting

$$\mathbf{P}^{-1} \triangleq \left[\begin{array}{c|c} \mathbf{P}_{xx} & \mathbf{P}_{xy} \\ \hline \mathbf{P}_{yx} & \mathbf{P}_{yy} \end{array} \right]^{-1} \triangleq \left[\begin{array}{c|c} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \hline \mathbf{A}_{12}^T & \mathbf{A}_{22} \end{array} \right]$$

and solving $\mathbf{P}^{-1}\mathbf{P} = \mathbf{I}$ for \mathbf{A}_{11} , \mathbf{A}_{12} , and \mathbf{A}_{22} as

$$\begin{aligned} \mathbf{A}_{11} &= (\mathbf{P}_{xx} - \mathbf{P}_{xy}\mathbf{P}_{yy}^{-1}\mathbf{P}_{yx})^{-1}, & \mathbf{A}_{22} &= (\mathbf{P}_{yy} - \mathbf{P}_{yx}\mathbf{P}_{xx}^{-1}\mathbf{P}_{xy})^{-1} \\ \mathbf{A}_{12} &= -\mathbf{A}_{11}\mathbf{P}_{xy}\mathbf{P}_{yy}^{-1} = -\mathbf{P}_{xx}^{-1}\mathbf{P}_{xy}\mathbf{A}_{22} \end{aligned}$$

(b) Show that equivalent expressions for \mathbf{A}_{11} and \mathbf{A}_{22} are

$$\mathbf{A}_{11} = \mathbf{P}_{xx}^{-1} + \mathbf{P}_{xx}^{-1}\mathbf{P}_{xy}\mathbf{A}_{22}\mathbf{P}_{yx}\mathbf{P}_{xx}^{-1}, \quad \mathbf{A}_{22} = \mathbf{P}_{yy}^{-1} + \mathbf{P}_{yy}^{-1}\mathbf{P}_{yx}\mathbf{A}_{11}\mathbf{P}_{xy}\mathbf{P}_{yy}^{-1}$$

Why might this be of use?

(c) Use these results and (3-110) to develop (3-111)–(3-113).

3.17 Prove the claim associated with (3-144) for the case of two measurement vector partitions (an inductive proof for the general case is a simple extension). Let

$$\mathbf{z} = \begin{bmatrix} \mathbf{z}_1 \\ \mathbf{z}_2 \end{bmatrix} = \begin{bmatrix} \mathbf{H}_1 \\ \mathbf{H}_2 \end{bmatrix} \mathbf{x} + \begin{bmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \end{bmatrix}, \quad \mathbf{R} = \begin{bmatrix} \mathbf{R}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{R}_2 \end{bmatrix}$$

Show that two recursions of (3-141) and (3-142) for \mathbf{z}_1 and \mathbf{z}_2 , respectively, yields equivalent results to one application of these equations to incorporate \mathbf{z} .

3.18 The end of Section 3.11 generated a weighted least squares estimate of \mathbf{x} as given by (3-151):

$$\hat{\mathbf{x}}_{\text{WLS}} = [\mathbf{H}^T \mathbf{W} \mathbf{H}]^{-1} \mathbf{H}^T \mathbf{W} \mathbf{z}$$

Now convert this into a recursive technique. Let \mathbf{z} be m -dimensional and write the result of the above equation as $\hat{\mathbf{x}}_m$. Now assume an additional scalar measurement value z_{m+1} becomes available; $\hat{\mathbf{x}}_{m+1}$ could be generated in the same manner for an $(m+1)$ -dimensional measurement. However, if the “new” values of the $(m+1)$ -by- n \mathbf{H}_{m+1} and the $(m+1)$ -by- $(m+1)$ \mathbf{W}_{m+1} are

$$\mathbf{H}_{m+1} = \begin{bmatrix} \mathbf{H}_m \\ \mathbf{h}_{m+1}^T \end{bmatrix}, \quad \mathbf{W}_{m+1} = \begin{bmatrix} \mathbf{W}_m & \mathbf{0} \\ \mathbf{0} & w_{m+1} \end{bmatrix}$$

then show the result can be written equivalently as

$$\begin{aligned} \tilde{\mathbf{P}}_m &= [\mathbf{H}^T \mathbf{W} \mathbf{H}]^{-1} \\ \hat{\mathbf{x}}_{m+1} &= \hat{\mathbf{x}}_m + \tilde{\mathbf{P}}_m \mathbf{h}_{m+1} [\mathbf{h}_{m+1}^T \tilde{\mathbf{P}}_m \mathbf{h}_{m+1} + (1/w_{m+1})]^{-1} [z_{m+1} - \mathbf{h}_{m+1}^T \hat{\mathbf{x}}_m] \\ \tilde{\mathbf{P}}_{m+1} &= \tilde{\mathbf{P}}_m - \tilde{\mathbf{P}}_m \mathbf{h}_{m+1} [\mathbf{h}_{m+1}^T \tilde{\mathbf{P}}_m \mathbf{h}_{m+1} + (1/w_{m+1})]^{-1} \mathbf{h}_{m+1}^T \tilde{\mathbf{P}}_m \end{aligned}$$

3.19 Consider the estimation of a vector, \mathbf{x} , composed of n constant parameters, based upon m measurements. Let the relationship between the measurements and parameters be given as

$$\mathbf{z} = \mathbf{H}\mathbf{x} + \mathbf{v}$$

Let \mathbf{x} and \mathbf{v} be modeled as jointly Gaussian, zero-mean vectors with

$$E[\mathbf{x}\mathbf{x}^T] = \mathbf{X}, \quad E[\mathbf{v}\mathbf{v}^T] = \mathbf{R}, \quad E[\mathbf{x}\mathbf{v}^T] = \mathbf{S}$$

What is the conditional mean $E[\mathbf{x}|\mathbf{z} = \mathbf{z}]$? If this were used as an estimate of \mathbf{x} , what is the corresponding conditional error covariance?

Problem 3.19 is modified from *Uncertain Dynamic Systems* by F. C. Schweppe. © 1973. Used with permission of Prentice-Hall, Inc.

3.20 In Section 3.9 an expression was developed for the conditional covariance of a Gaussian random variable \mathbf{x} , conditioned on the fact that a random variable \mathbf{z} , jointly Gaussian with \mathbf{x} , assumed some value, \mathbf{z} . We called this conditional covariance $\mathbf{P}_{\mathbf{x}|\mathbf{z}}$.

Later, in Section 3.11, this $\mathbf{P}_{\mathbf{x}|\mathbf{z}}$ was called the covariance of the *error* associated with using $E[\mathbf{x}|\mathbf{z} = \mathbf{z}]$ as an *estimate* of the value of \mathbf{x} . Explain this new interpretation of $\mathbf{P}_{\mathbf{x}|\mathbf{z}}$ explicitly. Show that it is a valid interpretation by showing that if \mathbf{y} is a continuous (Baire) function of \mathbf{z} , $\mathbf{y} = \boldsymbol{\theta}(\mathbf{z})$, then

$$E\{\mathbf{xy}^T | \mathbf{z}(\omega_j) = \mathbf{z}\} = E\{\mathbf{x} | \mathbf{z}(\omega_j) = \mathbf{z}\} [\boldsymbol{\theta}^T(\mathbf{z})]$$

and use this result to prove (3-140).

Is $\mathbf{P}_{\mathbf{x}|\mathbf{z}}$ a function of the value \mathbf{z} that \mathbf{z} assumes? So what?

3.21 Assume you are responsible for increasing the position-tracking precision of a flight test range. Presently, you have a radar tracking system capable of measuring position ± 10 ft (peak expected errors due to very wideband noise). You desire to double the precision to ± 5 ft approximately. For equal cost you could either

(a) optimally combine the present radar data with a new radar system's data, where the new system alone provides position ± 6 ft (peak errors due to very wide band noise) or

(b) triplicate the original system and combine data optimally.

Which would you propose to do and why? State all modeling assumptions explicitly.

3.22 It is desired to estimate the value assumed by some zero-mean scalar random variable x using the conditional mean of x , given the values of three other zero-mean scalar random variables, z_1 , z_2 , and z_3 . Prove by counterexample that the following two "reasonable" statements are actually false.

(a) If $E[xz_1] \neq 0$, then the value of z_1 , i.e., z_1 , is always a part of the best estimate of $x(\omega_k) = x$.

(b) If $E[xz_1] = 0$, then the value of $z_1(z_1)$ is never of use in estimating x .

Note what this implies about the common practice in economics and other fields of judging whether a variable should be included in an analysis based on its correlation to the variable of interest.

Problems 3.22 and 3.23 are modified from *Uncertain Dynamic Systems* by F. C. Schweppe. © 1973. Used with permission of Prentice-Hall, Inc.

3.23 The conditional error covariance matrix \mathbf{P}^+ derived in Section 3.11 is independent of the values \mathbf{z} assumed by the measurements \mathbf{z} . Thus, an error analysis can be performed before the measurements are taken, to decide how accurate the estimate will be when (if) it is actually calculated.

Assume that two meters are to provide measurements of a parameter, x , that you want to determine. Let the a priori knowledge of x indicate that it is well modeled as Gaussian, zero-mean, and having a variance of 8. Let the measurements be of the form

$$z_1 = x + v_1, \quad z_2 = x + v_2$$

One meter has been built and is assumed to be such that v_1 is Gaussian, zero-mean with variance of unity. The other meter has not yet been built, and v_2 can be assumed to be Gaussian, zero-mean, and of variance R , where R is a design parameter.

Assume system specifications require that the final estimate must have an error with variance less than or equal to $\frac{1}{2}$. Since accurate meters cost money, it is reasonable to try to find the maximum value of R that is acceptable. Find this R . With that R , determine the equations for the estimator that incorporates both z_1 and z_2 simultaneously. Now find the equations used to incorporate z_1 to obtain an estimate, and then recursively incorporate z_2 into the estimate. Show that these are the same estimates with the same error variances.

3.24 This and subsequent problems are concerned with estimation of the moments [1, 2, 8, 11, 14] of a random variable $x(\cdot)$, based only upon N realized values, x_1, x_2, \dots, x_N , that can be considered to be empirical data. The distribution and/or density function are unknown, and assume that the true (unknown) mean and variance of $x(\cdot)$ are μ_x and σ_x^2 , respectively.

A logical choice of estimate of the mean would be

$$\hat{M}_x = \frac{1}{N} \sum_{i=1}^N x_i$$

To consider the error committed by using this estimate, let $x_1(\cdot), x_2(\cdot), \dots, x_N(\cdot)$ be N random variables, each with distribution identical to that of $x(\cdot)$. Thus, we can generate the estimator

$$\hat{M}_x = \frac{1}{N} \sum_{i=1}^N x_i$$

and conceive of conducting an experiment of generating the N data points repeatedly, the j th such experiment yielding a single realization of $\hat{M}_x, \hat{M}_x(\omega_j)$. We want to characterize the distribution of these estimate values.

- Show that $E\{\hat{M}_x\} = \mu_x$: that \hat{M}_x is an unbiased estimator of the mean of $x(\cdot)$.
- Show that the variance of \hat{M}_x , denoted as $\sigma_{\hat{m}}^2$, is

$$\sigma_{\hat{m}}^2 = \frac{1}{N} E\{x^2\} + \frac{2}{N^2} \sum_{i=1}^{N-1} \sum_{j=i+1}^N E\{x_i x_j\} - \mu_x^2$$

so that, if the observations are independent of each other (i.e., x_1, x_2, \dots, x_N are a set of independent random variables), then

$$\sigma_{\hat{m}}^2 = (1/N)\sigma_x^2$$

Why is $\sigma_{\hat{m}}^2$ also the variance of the error committed by using \hat{M}_x to estimate the mean of x ?

- As N is increased, not only does the estimate become more precise, \hat{M}_x becomes more and more Gaussian, regardless of the distribution of x . (Why?) If \hat{M}_x is assumed to be Gaussian, how many observations should be made (i.e., at least how large should N be) so that the probability that the error in the estimate is less than 10% of σ_x is 0.954? (Answer = 400.)

3.25 Since the variance of $x(\cdot)$ is defined as

$$\sigma_x^2 \triangleq E\{[x - E\{x\}]^2\} = E\{x^2\} - [E\{x\}]^2$$

a reasonable estimator of the variance of $x(\cdot)$ would be

$$\hat{V}_x' = \frac{1}{N} \sum_{i=1}^N (x_i - \hat{M}_x)^2 = \frac{1}{N} \left\{ \sum_{i=1}^N x_i^2 \right\} - \hat{M}_x^2$$

with \hat{M}_x defined in the previous problem.

- Demonstrate the equality of the two forms of \hat{V}_x' .

(b) Although this is a reasonable estimate (and the maximum likelihood estimate), show that it is a biased estimate, that

$$E\{\hat{V}_x'\} = \sigma_x^2 - \sigma_{\hat{m}}^2 \neq \sigma_x^2$$

with $\sigma_{\hat{m}}^2$ defined in the previous problem. For independent observations, show that this becomes

$$E\{\hat{V}_x'\} = [(N-1)/N]\sigma_x^2$$

(c) Thus, for independent observations, a variance estimate of $[N/(N-1)]V_x'$, or

$$\hat{V}_x = \frac{1}{N-1} \sum_{i=1}^N (x_i - \hat{M}_x)^2 = \frac{1}{N-1} \sum_{i=1}^N x_i^2 - \frac{N}{N-1} \hat{M}_x^2$$

will yield an unbiased estimate of the variance of $x(\cdot)$, σ_x^2 . If observations are not independent, show that

$$E\{\hat{V}_x\} = [N/(N-1)]\{\sigma_x^2 - \sigma_{\hat{m}}^2\}$$

and that the mean error is less than that committed by \hat{V}_x' . This estimator is defined only for $N > 1$; would an estimate of variance for $N = 1$ be meaningful? Why would the second form of \hat{V}_x be preferable computationally?

(d) Show that the variance of \hat{V}_x , for independent observations and $N > 1$, is

$$\sigma_{\hat{v}}^2 = (1/N)[E\{x - \mu_x\}^4] - \{(N-3)/(N-1)\}\sigma_x^4]$$

If $x(\cdot)$ were assumed Gaussian, show that the fourth central moment is equal to $3\sigma_x^4$, and thus

$$\sigma_{\hat{v}}^2 = [2/(N-1)]\sigma_x^4 \quad [\text{if } x(\cdot) \text{ is Gaussian}]$$

If $x(\cdot)$ were instead assumed to be uniform, show that the fourth central moment is equal to $\frac{9}{5}\sigma_x^4$, and so

$$\sigma_{\hat{v}}^2 = [(4N+6)/(5N(N-1))]\sigma_x^4 \quad [\text{if } x(\cdot) \text{ is uniform}]$$

Thus, under very different assumptions about $x(\cdot)$, the quality of the estimate provided by \hat{V}_x is very similar:

$$\sigma_{\hat{v}} \quad (x \text{ Gaussian}) \cong 1.4\sigma_{\hat{v}} \quad (x \text{ uniform})$$

3.26 Analogous to the variance estimator of the previous problem, a good estimate of the covariance between $x(\cdot)$ and $y(\cdot)$ is

$$\hat{C}_{xy} = \frac{1}{N-1} \sum_{i=1}^N (x_i - \hat{M}_x)(y_i - \hat{M}_y) = \frac{1}{N-1} \sum_{i=1}^N x_i y_i - \frac{N}{N-1} \hat{M}_x \hat{M}_y$$

where $\hat{M}_x = (1/N) \sum_{i=1}^N x_i$ and $\hat{M}_y = (1/N) \sum_{i=1}^N y_i$, and the second form of \hat{C}_{xy} is more convenient computationally.

(a) Show that, for independent observations, \hat{C}_{xy} is an unbiased estimate:

$$E\{\hat{C}_{xy}\} = \sigma_{xy}^2 = \text{true covariance of } x(\cdot) \text{ and } y(\cdot)$$

(b) For independent measurements, the variance of \hat{C}_{xy} is

$$\sigma_{\hat{c}}^2 = (1/N)[E\{[x - \mu_x]^2[y - \mu_y]^2\} + [1/(N-1)]\{\sigma_x^2\sigma_y^2 - (N-2)\sigma_{xy}^4\}]$$

Show that this measure of the quality of the \hat{C}_{xy} estimate becomes, if $x(\cdot)$ and $y(\cdot)$ are assumed jointly Gaussian,

$$\sigma_{\hat{c}}^2 = [1/(N-1)]\sigma_x^2\sigma_y^2[1 + r_{xy}^2]$$

where r_{xy} is the correlation coefficient between $x(\cdot)$ and $y(\cdot)$.

(c) From the definition of correlation coefficient, (3-75), a good estimator of the (linear) correlation coefficient can be produced as

$$\hat{r}_{xy} = \frac{\hat{C}_{xy}}{(\hat{V}_x \hat{V}_y)^{1/2}} = \frac{N \sum_{i=1}^N x_i y_i - (\sum_{i=1}^N x_i)(\sum_{i=1}^N y_i)}{([N \sum_{i=1}^N x_i^2 - (\sum_{i=1}^N x_i)^2][N \sum_{i=1}^N y_i^2 - (\sum_{i=1}^N y_i)^2])^{1/2}}$$

A “scatter diagram” is a two-dimensional plot of the N realizations $[x(\omega_1), y(\omega_1)], [x(\omega_2), y(\omega_2)], \dots, [x(\omega_N), y(\omega_N)]$, as shown in Fig. 3.P4. Plot (a) shows perfect correlation between the two variables, (b) portrays smaller positive correlation, (c) depicts no correlation, and (d) shows negative correlation. From plots (a) and (b) it can be seen that the magnitude of the correlation coefficient depicts the dispersion of the points from the least squares (regression) line fit to the data, and not the slope of the line itself. Verify these claims by calculating \hat{r}_{xy} for the four plots of Fig. 3.P4.

The least squares regression line of y on x (least squares fit of a line to the data, with “residuals” being the distance between data points and the line measured in the y direction) is given by

$$y - \hat{M}_y = [\hat{C}_{xy}/\hat{V}_x](x - \hat{M}_x)$$

On the other hand, the least squares regression line of x on y is

$$x - \hat{M}_x = [\hat{C}_{xy}/\hat{V}_y](y - \hat{M}_y)$$

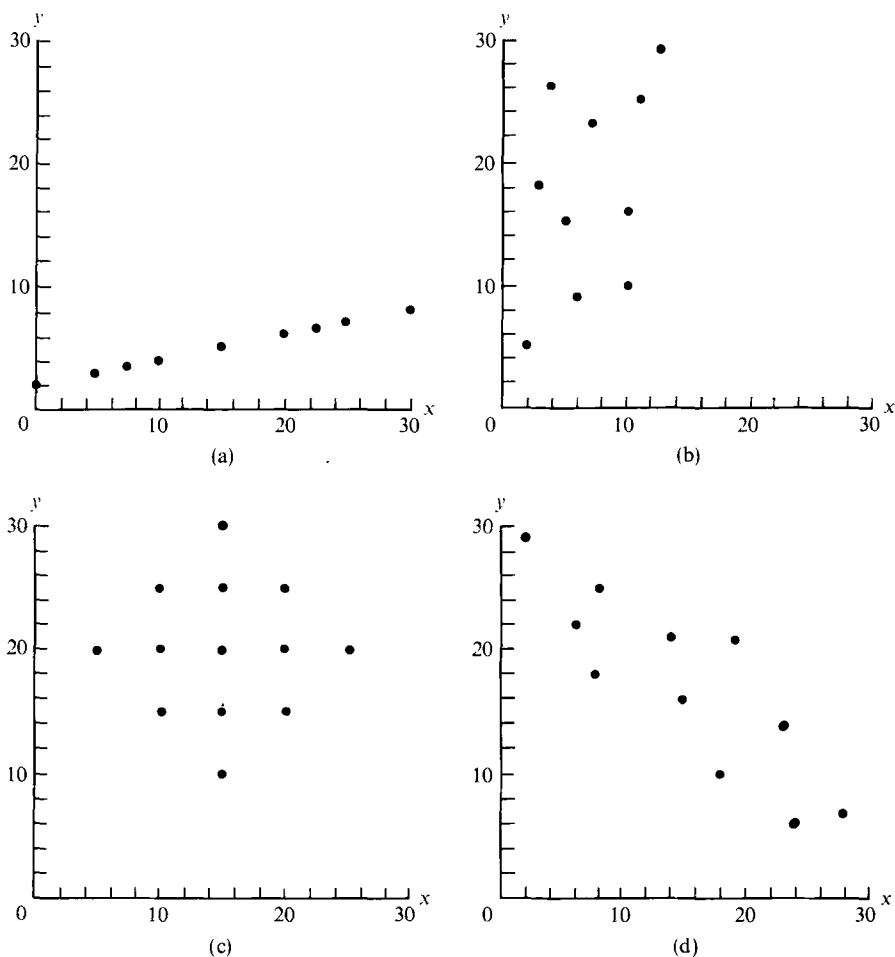


FIG. 3.P4 Scatter diagrams for Problem 3.26.

Thus, each line passes through the "centroid" (\hat{M}_x, \hat{M}_y) , and the product of their slopes is \hat{r}_{xy}^2 . The linear correlation coefficient is a measure of the departure of the two regression lines, with colinearity indicated by $\hat{r}_{xy} = \pm 1$ (perfect linear correlation) and orthogonality by $\hat{r}_{xy} = 0$ (no linear correlation). Verify these interpretations for the four cases depicted in Fig. 3.P4.

3.27 The previous three problems assumed perfect measurements of realized values. Now assume noise corruption of the measuring device, so that what are available are realizations of

$$y(\cdot) = x(\cdot) + w(\cdot)$$

where $w(\cdot)$ is zero mean, of variance σ_w^2 , and independent of $x(\cdot)$. Assume independent observations.

(a) Show that $\hat{M}_x = (1/N) \sum_{i=1}^N y_i$ is still an unbiased mean estimator, but with increased variance: $\sigma_m^2 = (1/N)(\sigma_x^2 + \sigma_w^2)$.

(b) Show that $\hat{V}_x = [1/(N-1)] \sum_{i=1}^N y_i^2 - [N/(N-1)] \hat{M}_x^2$, however, is a biased estimator: $E\{\hat{V}_x\} = \sigma_x^2 + \sigma_w^2$; thus, the best estimator would be $[\hat{V}_x - \sigma_w^2]$. How would σ_w^2 be estimated? If x and w are assumed Gaussian, show that

$$\sigma_{\hat{x}} \text{ (with meas. noise)} = [1 + (\sigma_w^2/\sigma_x^2)]\sigma_{\hat{x}} \text{ (without meas. noise)}.$$