
Optimization-Based Control

Richard M. Murray
Control and Dynamical Systems
California Institute of Technology

DRAFT v2.1a, February 15, 2010
© California Institute of Technology
All rights reserved.

This manuscript is for review purposes only and may not be reproduced, in whole or in part, without written consent from the author.

Chapter 4

Stochastic Systems

In this chapter we present a focused overview of stochastic systems, oriented toward the material that is required in Chapters 5 and 6. After a brief review of random variables, we define discrete-time and continuous-time random processes, including the expectation, (co-)variance and correlation functions for a random process. These definitions are used to describe linear stochastic systems (in continuous time) and the stochastic response of a linear system to a random process (e.g., noise). We initially derive the relevant quantities in the state space, followed by a presentation of the equivalent frequency domain concepts.

Prerequisites. Readers should be familiar with basic concepts in probability, including random variables and standard distributions. We do not assume any prior familiarity with random processes.

Caveats. This chapter is written to provide a brief introduction to stochastic processes that can be used to derive the results in the following chapters. In order to keep the presentation compact, we gloss over several mathematical details that are required for rigorous presentation of the results. A more detailed (and mathematically precise) derivation of this material is available in the book by Åström [Åst06].

4.1 Brief Review of Random Variables

To help fix the notation that we will use, we briefly review the key concepts of random variables. A more complete exposition is available in standard books on probability, such as Hoel, Port and Stone [HPS71].

Random variables and processes are defined in terms of an underlying *probability space* that captures the nature of the stochastic system we wish to study. A probability space has three elements:

- a *sample space* Ω that represents the set of all possible outcomes;
- a set of *events* \mathcal{F} that captures combinations of elementary outcomes that are of interest; and
- a *probability measure* \mathcal{P} that describes the likelihood of a given event occurring.

Ω can be any set, either with a finite, countable or infinite number of elements. The event space \mathcal{F} consists of subsets of Ω . There are some mathematical limits on the properties of the sets in \mathcal{F} , but these are not critical for our purposes here. The probability measure \mathcal{P} is a mapping from $\mathcal{P} : \mathcal{F} \rightarrow [0, 1]$ that assigns a probability to each event. It must satisfy the property that given any two disjoint set $A, B \subset \mathcal{F}$,

$P(A \cup B) = P(A) + P(B)$. The term *probability distribution* is also used to describe a probability measure.

With these definitions, we can model many different stochastic phenomena. Given a probability space, we can choose samples $\omega \in \Omega$ and identify each sample with a collection of events chosen from \mathcal{F} . These events should correspond to phenomena of interest and the probability measure \mathcal{P} should capture the likelihood of that event occurring in the system that we are modeling. This definition of a probability space is very general and allows us to consider a number of situations as special cases.

A *random variable* X is a function $X : \Omega \rightarrow S$ that gives a value in S , called the state space, for any sample $\omega \in \Omega$. Given a subset $A \subset S$, we can write the probability that $X \in A$ as

$$P(X \in A) = P(\omega \in \Omega : X(\omega) \in A).$$

We will often find it convenient to omit ω when working with random variables and hence we write $X \in S$ rather than the more correct $X(\omega) \in S$.

A *continuous (real-valued) random variable* X is a variable that can take on any value in the set of real numbers \mathbb{R} . We can model the random variable X according to its *probability distribution* P :

$$P(x_l \leq X \leq x_u) = \text{probability that } x \text{ takes on a value in the range } x_l, x_u.$$

More generally, we write $P(A)$ as the probability that an event A will occur (e.g., $A = \{x_l \leq X \leq x_u\}$). It follows from the definition that if X is a random variable in the range $[L, U]$ then $P(L \leq X \leq U) = 1$. Similarly, if $Y \in [L, U]$ then $P(L \leq X \leq Y) = 1 - P(Y \leq X \leq U)$.

We characterize a random variable in terms of the *probability density function* (pdf) $p(x)$. The density function is defined so that its integral over an interval gives the probability that the random variable takes its value in that interval:

$$P(x_l \leq X \leq x_u) = \int_{x_l}^{x_u} p(x) dx. \quad (4.1)$$

It is also possible to compute $p(x)$ given the distribution P as long as the distribution is suitably smooth:

$$p(x) = \left. \frac{\partial P(x_l \leq x \leq x_u)}{\partial x_u} \right|_{\substack{x_l \text{ fixed,} \\ x_u = x}}, \quad x > x_l.$$

We will sometimes write $p_X(x)$ when we wish to make explicit that the pdf is associated with the random variable X . Note that we use capital letters to refer to a random variable and lower case letters to refer to a specific value.

Probability distributions provide a general way to describe stochastic phenomena. Some standard probability distributions include a *uniform distribution*,

$$p(x) = \frac{1}{U - L}, \quad (4.2)$$

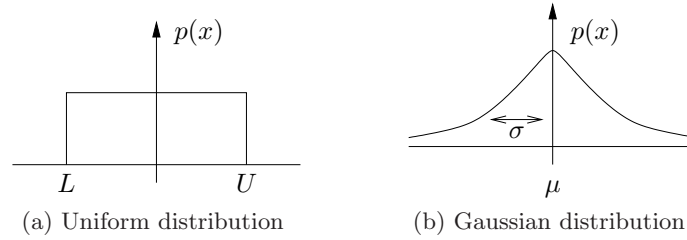


Figure 4.1: Probability density function (pdf) for uniform and Gaussian distributions.

and a *Gaussian distribution* (also called a *normal distribution*),

$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}. \quad (4.3)$$

In the Gaussian distribution, the parameter μ is called the *mean* of the distribution and σ is called the *standard deviation* of the distribution. Figure 4.1 gives a graphical representation of uniform and Gaussian pdfs. There many other distributions that arise in applications, but for the purpose of these notes we focus on uniform distributions and Gaussian distributions.

If two random variables are related, we can talk about their *joint probability distribution*: $P_{X,Y}(A, B)$ is the probability that both event A occurs for X and B occurs for Y . This is sometimes written as $P(A \cap B)$, where we abuse notation by implicitly assuming that A is associated with X and B with Y . For continuous random variables, the joint probability distribution can be characterized in terms of a *joint probability density function*

$$P(x_l \leq X \leq x_u, y_l \leq Y \leq y_u) = \int_{y_l}^{y_u} \int_{x_l}^{x_u} p(x, y) dx dy. \quad (4.4)$$

The joint pdf thus describes the relationship between X and Y , and for sufficiently smooth distributions we have

$$p(x, y) = \frac{\partial^2 P(x_l \leq X \leq x_u, y_l \leq Y \leq y_u)}{\partial x_u \partial y_u} \Bigg|_{\substack{x_l, y_l \text{ fixed,} \\ x_u = x, y_u = y}}, \quad \begin{array}{l} x > x_l, \\ y > y_l. \end{array}$$

We say that X and Y are *independent* if $p(x, y) = p(x)p(y)$, which implies that $P_{X,Y}(A, B) = P_X(A)P_Y(B)$ for events A associated with X and B associated with Y . Equivalently, $P(A \cap B) = P(A)P(B)$ if A and B are independent.

The *conditional probability* for an event A given that an event B has occurred, written as $P(A|B)$, is given by

$$P(A|B) = \frac{P(A \cap B)}{P(B)}. \quad (4.5)$$

If the events A and B are independent, then $P(A|B) = P(A)$. Note that the individual, joint and conditional probability distributions are all different, so we should really write $P_{X,Y}(A \cap B)$, $P_{X|Y}(A|B)$ and $P_Y(B)$.

If X is dependent on Y then Y is also dependent on X . *Bayes' theorem* relates the conditional and individual probabilities:

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}, \quad P(B) \neq 0. \quad (4.6)$$

Bayes' theorem gives the conditional probability of event A on event B given the inverse relationship (B given A). It can be used in situations in which we wish to evaluate a hypothesis H given data D when we have some model for how likely the data is given the hypothesis, along with the unconditioned probabilities for both the hypothesis and the data.

The analog of the probability density function for conditional probability is the *conditional probability density function* $p(x|y)$

$$p(x|y) = \begin{cases} \frac{p(x,y)}{p(y)} & 0 < p(y) < \infty \\ 0 & \text{otherwise.} \end{cases} \quad (4.7)$$

It follows that

$$p(x,y) = p(x|y)p(y) \quad (4.8)$$

and

$$\begin{aligned} P(x_l \leq X \leq x_u|y) &:= P(x_l \leq X \leq x_u|Y = y) \\ &= \int_{x_l}^{x_u} p(x|y)dx = \frac{\int_{x_l}^{x_u} p(x,y)dx}{p(y)}. \end{aligned} \quad (4.9)$$

If X and Y are independent then $p(x|y) = p(x)$ and $p(y|x) = p(y)$. Note that $p(x,y)$ and $p(x|y)$ are different density functions, though they are related through equation (4.8). If X and Y are related with joint probability density function $p(x,y)$ and conditional probability density function $p(x|y)$ then

$$p(x) = \int_{-\infty}^{\infty} p(x,y)dy = \int_{-\infty}^{\infty} p(x|y)p(y)dy.$$

Example 4.1 Conditional probability for sum

Consider three random variables X , Y and Z related by the expression

$$Z = X + Y.$$

In other words, the value of the random variable Z is given by choosing values from two random variables X and Y and adding them. We assume that X and Y are independent Gaussian random variables with mean μ_1 and μ_2 and standard deviation $\sigma = 1$ (the same for both variables).

Clearly the random variable Z is not independent of X (or Y) since if we know the values of X then it provides information about the likely value of Z . To see this, we compute the joint probability between Z and X . Let

$$A = \{x_l \leq x \leq x_u\}, \quad B = \{z_l \leq z \leq z_u\}.$$

The joint probability of both events A and B occurring is given by

$$\begin{aligned} P_{X,Z}(A \cap B) &= P(x_l \leq x \leq x_u, z_l \leq x + y \leq z_u) \\ &= P(x_l \leq x \leq x_u, z_l - x \leq y \leq z_u - x). \end{aligned}$$

We can compute this probability by using the probability density functions for X and Y :

$$\begin{aligned} P(A \cap B) &= \int_{x_l}^{x_u} \left(\int_{z_l-x}^{z_u-x} p_Y(y) dy \right) p_X(x) dx \\ &= \int_{x_l}^{x_u} \int_{z_l}^{z_u} p_Y(z-x) p_X(x) dz dx =: \int_{z_l}^{z_u} \int_{x_l}^{x_u} p_{Z,X}(z,x) dx dz. \end{aligned}$$

Using Gaussians for X and Y we have

$$\begin{aligned} p_{Z,X}(z,x) &= \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(z-x-\mu_Y)^2} \cdot \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(x-\mu_X)^2} \\ &= \frac{1}{2\pi} e^{-\frac{1}{2}((z-x-\mu_Y)^2 + (x-\mu_X)^2)}. \end{aligned}$$

A similar expression holds for $p_{Z,Y}$. ∇

Given a random variable X , we can define various standard measures of the distribution. The *expectation* or *mean* of a random variable is defined as

$$\mathbb{E}[X] = \langle X \rangle = \int_{-\infty}^{\infty} x p(x) dx,$$

and the *mean square* of a random variable is

$$\mathbb{E}[X^2] = \langle X^2 \rangle = \int_{-\infty}^{\infty} x^2 p(x) dx.$$

If we let μ represent the expectation (or mean) of X then we define the *variance* of X as

$$\mathbb{E}[(X-\mu)^2] = \langle (X-\langle X \rangle)^2 \rangle = \int_{-\infty}^{\infty} (x-\mu)^2 p(x) dx.$$

We will often write the variance as σ^2 . As the notation indicates, if we have a Gaussian random variable with mean μ and (stationary) standard deviation σ , then the expectation and variance as computed above return μ and σ^2 .

Several useful properties follow from the definitions.

Proposition 4.1 (Properties of random variables).

1. If X is a random variable with mean μ and variance σ^2 , then αX is random variable with mean $\alpha\mu$ and variance $\alpha^2\sigma^2$.
2. If X and Y are two random variables, then $\mathbb{E}[\alpha X + \beta Y] = \alpha\mathbb{E}[X] + \beta\mathbb{E}[Y]$.
3. If X and Y are Gaussian random variables with means μ_X, μ_Y and variances σ_X^2, σ_Y^2 ,

$$p(x) = \frac{1}{\sqrt{2\pi\sigma_X^2}} e^{-\frac{1}{2}\left(\frac{x-\mu_X}{\sigma_X}\right)^2}, \quad p(y) = \frac{1}{\sqrt{2\pi\sigma_Y^2}} e^{-\frac{1}{2}\left(\frac{y-\mu_Y}{\sigma_Y}\right)^2},$$

then $X+Y$ is a Gaussian random variable with mean $\mu_Z = \mu_X + \mu_Y$ and variance $\sigma_Z^2 = \sigma_X^2 + \sigma_Y^2$,

$$p(x+y) = \frac{1}{\sqrt{2\pi\sigma_Z^2}} e^{-\frac{1}{2}\left(\frac{x+y-\mu_Z}{\sigma_Z}\right)^2}.$$

Proof. The first property follows from the definition of mean and variance:

$$\begin{aligned}\mathbb{E}[\alpha X] &= \int_{-\infty}^{\infty} \alpha x p(x) dx = \alpha \int_{-\infty}^{\infty} x p(x) dx = \alpha \mathbb{E}[X] \\ \mathbb{E}[(\alpha X)^2] &= \int_{-\infty}^{\infty} (\alpha x)^2 p(x) dx = \alpha^2 \int_{-\infty}^{\infty} x^2 p(x) dx = \alpha^2 \mathbb{E}[X^2].\end{aligned}$$

The second property follows similarly, remembering that we must take the expectation using the joint distribution (since we are evaluating a function of two random variables):

$$\begin{aligned}\mathbb{E}[\alpha X + \beta Y] &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (\alpha x + \beta y) p_{X,Y}(x, y) dx dy \\ &= \alpha \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} x p_{X,Y}(x, y) dx dy + \beta \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} y p_{X,Y}(x, y) dx dy \\ &= \alpha \int_{-\infty}^{\infty} x p_X(x) dx + \beta \int_{-\infty}^{\infty} y p_Y(y) dy = \alpha \mathbb{E}[X] + \beta \mathbb{E}[Y].\end{aligned}$$

The third item is left as an exercise. □

4.2 Introduction to Random Processes

A *random process* is a collection of time-indexed random variables. Formally, we consider a random process X to be a joint mapping of sample and a time to a state: $X : \Omega \times \mathcal{T} \rightarrow S$, where \mathcal{T} is an appropriate time set. We view this mapping as a generalized random variable: a sample corresponds to choosing an entire function of time. Of course, we can always fix the time and interpret $X(\omega, t)$ as a regular random variable, with $X(\omega, t')$ representing a different random variable if $t \neq t'$. Our description of random processes will consist of describing how the random variable at a time t relates to the value of the random variable at an earlier time s . To build up some intuition about random processes, we will begin with the discrete time case, where the calculations are a bit more straightforward, and then proceed to the continuous time case.

A *discrete-time random process* is a stochastic system characterized by the *evolution* of a sequence of random variables $X[k]$, where k is an integer. As an example, consider a discrete-time linear system with dynamics

$$x[k+1] = Ax[k] + Bu[k] + Fw[k], \quad y[k] = Cx[k] + v[k]. \quad (4.10)$$

As in ÅM08, $x \in \mathbb{R}^n$ represents the state of the system, $u \in \mathbb{R}^m$ is the vector of inputs and $y \in \mathbb{R}^p$ is the vector of outputs. The (possibly vector-valued) signal w represents disturbances to the process dynamics and v represents noise in the measurements. To try to fix the basic ideas, we will take $u = 0$, $n = 1$ (single state) and $F = 1$ for now.

We wish to describe the evolution of the dynamics when the disturbances and noise are not given as deterministic signals, but rather are chosen from some probability distribution. Thus we will let $W[k]$ be a collection of random variables where

the values at each instant k are chosen from the probability distribution $P[W, k]$. As the notation indicates, the distributions might depend on the time instant k , although the most common case is to have a *stationary* distribution in which the distributions are independent of k (defined more formally below).

In addition to stationarity, we will often also assume that distribution of values of W at time k is independent of the values of W at time l if $k \neq l$. In other words, $W[k]$ and $W[l]$ are two separate random variables that are independent of each other. We say that the corresponding random process is *uncorrelated* (also defined more formally below). As a consequence of our independence assumption, we have that

$$\mathbb{E}[W[k]W[l]] = \mathbb{E}[W^2[k]]\delta(k-l) = \begin{cases} \mathbb{E}[W^2[k]] & k = l \\ 0 & k \neq l. \end{cases}$$

In the case that $W[k]$ is a Gaussian with mean zero and (stationary) standard deviation σ , then $\mathbb{E}[W[k]W[l]] = \sigma^2 \delta(k-l)$.

We next wish to describe the evolution of the state x in equation (4.10) in the case when W is a random variable. In order to do this, we describe the state x as a sequence of random variables $X[k]$, $k = 1, \dots, N$. Looking back at equation (4.10), we see that even if $W[k]$ is an uncorrelated sequence of random variables, then the states $X[k]$ are not uncorrelated since

$$X[k+1] = AX[k] + FW[k],$$

and hence the probability distribution for X at time $k+1$ depends on the value of X at time k (as well as the value of W at time k), similar to the situation in Example 4.1.

Since each $X[k]$ is a random variable, we can define the mean and variance as $\mu[k]$ and $\sigma^2[k]$ using the previous definitions at each time k :

$$\begin{aligned} \mu[k] &:= \mathbb{E}[X[k]] = \int_{-\infty}^{\infty} x p(x, k) dx, \\ \sigma^2[k] &:= \mathbb{E}[(X[k] - \mu[k])^2] = \int_{-\infty}^{\infty} (x - \mu[k])^2 p(x, k) dx. \end{aligned}$$

To capture the relationship between the current state and the future state, we define the *correlation function* for a random process as

$$\rho(k_1, k_2) := \mathbb{E}[X[k_1]X[k_2]] = \int_{-\infty}^{\infty} x_1 x_2 p(x_1, x_2; k_1, k_2) dx_1 dx_2$$

The function $p(x_i, x_j; k_1, k_2)$ is the *joint probability density function*, which depends on the times k_1 and k_2 . A process is *stationary* if $p(x, k+d) = p(x, k)$ for all k , $p(x_i, x_j; k_1+d, k_2+d) = p(x_i, x_j; k_1, k_2)$, etc. In this case we can write $p(x_i, x_j; d)$ for the joint probability distribution. We will almost always restrict to this case. Similarly, we will write $p(k_1, k_2)$ as $p(d) = p(k, k+d)$.

We can compute the correlation function by explicitly computing the joint pdf (see Example 4.1) or by directly computing the expectation. Suppose that we take a random process of the form (4.10) with $x[0] = 0$ and W having zero mean and

standard deviation σ . The correlation function is given by

$$\begin{aligned}\mathbb{E}[X[k_1]X[k_2]] &= E\left\{\left(\sum_{i=0}^{k_1-1} A^{k_1-i} BW[i]\right)\left(\sum_{j=0}^{k_2-1} A^{k_2-j} BW[j]\right)\right\} \\ &= E\left\{\sum_{i=0}^{k_1-1} \sum_{j=0}^{k_2-1} A^{k_1-i} BW[i]W[j]BA^{k_2-j}\right\}.\end{aligned}$$

We can now use the linearity of the expectation operator to pull this inside the summations:

$$\begin{aligned}\mathbb{E}[X[k_1]X[k_2]] &= \sum_{i=0}^{k_1-1} \sum_{j=0}^{k_2-1} A^{k_1-i} B \mathbb{E}[W[i]W[j]] BA^{k_2-j} \\ &= \sum_{i=0}^{k_1-1} \sum_{j=0}^{k_2-1} A^{k_1-i} B \sigma^2 \delta(i-j) BA^{k_2-j} \\ &= \sum_{i=0}^{k_1-1} A^{k_1-i} B \sigma^2 BA^{k_2-i}.\end{aligned}$$

Note that the correlation function depends on k_1 and k_2 .

We can see the dependence of the correlation function on the time more clearly by letting $d = k_2 - k_1$ and writing

$$\begin{aligned}\rho(k, k+d) &= \mathbb{E}[X[k]X[k+d]] = \sum_{i=0}^{k_1-1} A^{k-i} B \sigma^2 BA^{d+k-i} \\ &= \sum_{j=1}^k A^j B \sigma^2 BA^{j+d} = \left(\sum_{j=1}^k A^j B \sigma^2 BA^j\right) A^d.\end{aligned}$$

In particular, if the discrete time system is stable then $|A| < 1$ and the correlation function decays as we take points that are further departed in time (d large). Furthermore, if we let $k \rightarrow \infty$ (i.e., look at the steady state solution) then the correlation function only depends on d (assuming the sum converges) and hence the steady state random process is stationary.

In our derivation so far, we have assumed that $X[k+1]$ only depends on the value of the state at time k (this was implicit in our use of equation (4.10) and the assumption that $W[k]$ is independent of X). This particular assumption is known as the *Markov property* for a random process: a Markovian process is one in which the distribution of possible values of the state at time k depends only on the values of the state at the prior time and not earlier. Written more formally, we say that a discrete random process is Markovian if

$$p_{X,k}(x|X[k-1], X[k-2], \dots, X[0]) = p_{X,k}(x|X[k-1]). \quad (4.11)$$

Markov processes are roughly equivalent to state space dynamical systems, where the future evolution of the system can be completely characterized in terms of the current value of the state (and not its history of values prior to that).

4.3 Continuous-Time, Vector-Valued Random Processes

We now consider the case where our time index is no longer discrete, but instead varies continuously. A fully rigorous derivation requires careful use of measure theory and is beyond the scope of this text, so we focus here on the concepts that will be useful for modeling and analysis of important physical properties.

A *continuous-time random process* is a stochastic system characterized by the evolution of a random variable $X(t)$, $t \in [0, T]$. We are interested in understanding how the (random) state of the system is related at separate times. The process is defined in terms of the “correlation” of $X(t_1)$ with $X(t_2)$.

We call $X(t) \in \mathbb{R}^n$ the *state* of the random process at time t . For the case $n > 1$, we have a vector of random processes:

$$X(t) = \begin{bmatrix} X_1(t) \\ \vdots \\ X_n(t) \end{bmatrix}$$

We can characterize the state in terms of a (vector-valued) time-varying pdf,

$$P(x_l \leq X_i(t) \leq x_u) = \int_{x_l}^{x_u} p_{X_i}(x; t) dx.$$

Note that the state of a random process is not enough to determine the next state (otherwise it would be a deterministic process). We typically omit indexing of the individual states unless the meaning is not clear from context.

We can characterize the dynamics of a random process by its statistical characteristics, written in terms of *joint probability density functions*:

$$\begin{aligned} P(x_{1l} \leq X_i(t_1) \leq x_{1u}, x_{2l} \leq X_j(t_2) \leq x_{2u}) \\ = \int_{x_{2l}}^{x_{2u}} \int_{x_{1l}}^{x_{1u}} p_{X_i, X_j}(x_1, x_2; t_1, t_2) dx_1 dx_2 \end{aligned}$$

The function $p(x_i, x_j; t_1, t_2)$ is called a *joint probability density function* and depends both on the individual states that are being compared and the time instants over which they are compared. Note that if $i = j$, then p_{X_i, X_i} describes how X_i at time t_1 is related to X_i at time t_2 .

In general, the distributions used to describe a random process depend on the specific time or times that we evaluate the random variables. However, in some cases the relationship only depends on the difference in time and not the absolute times (similar to the notion of time invariance in deterministic systems, as described in ÅM08). A process is *stationary* if $p(x, t + \tau) = p(x, t)$ for all τ , $p(x_i, x_j; t_1 + \tau, t_2 + \tau) = p(x_i, x_j; t_1, t_2)$, etc. In this case we can write $p(x_i, x_j; \tau)$ for the joint probability distribution. Stationary distributions roughly correspond to the steady state properties of a random process and we will often restrict our attention to this case.

In looking at biomolecular systems, we are going to be interested in random processes in which the changes in the state occur when a random event occurs (such as a molecular reaction or binding event). In this case, it is natural to describe the

state of the system in terms of a set of times $t_0 < t_1 < t_2 < \dots < t_n$ and $X(t_i)$ is the random variable that corresponds to the possible states of the system at time t_i . Note that time instants do not have to be uniformly spaced and most often (for biomolecular systems) they will not be. All of the definitions above carry through, and the process can now be described by a probability distribution of the form

$$P\left(X(t_i) \in [x_i, x_i + dx_i], i = 1, \dots, n\right) = \int \dots \int p(x_n, x_{n-1}, \dots, x_0; t_n, t_{n-1}, \dots, t_0) dx_n dx_{n-1} dx_1,$$

where dx_i are taken as infinitesimal quantities.

An important class of stochastic systems is those for which the next state of the system depends only on the current state of the system and not the history of the process. Suppose that

$$\begin{aligned} P\left(X(t_n) \in [x_n, x_n + dx_n] | X(t_i) \in [x_i, x_i + dx_i], i = 1, \dots, n-1\right) \\ = P\left(X(t_n) \in [x_n, x_n + dx_n] | X(t_{n-1}) \in [x_{n-1}, x_{n-1} + dx_{n-1}]\right). \end{aligned} \quad (4.12)$$

That is, the probability of being in a given state at time t_n depends *only* on the state that we were in at the previous time instant t_{n-1} and not the entire history of states prior to t_{n-1} . A stochastic process that satisfies this property is called a *Markov process*.

In practice we do not usually specify random processes via the joint probability distribution $p(x_i, x_j; t_1, t_2)$ but instead describe them in terms of a *propogater function*. Let $X(t)$ be a Markov process and define the Markov propogater as

$$\Xi(dt; x, t) = X(t + dt) - X(t), \text{ given } X(t) = x.$$

The propogater function describes how the random variable at time t is related to the random variable at time $t + dt$. Since both $X(t + dt)$ and $X(t)$ are random variables, $\Xi(dt; x, t)$ is also a random variable and hence it can be described by its density function, which we denote as $\Pi(\xi, x; dt, t)$:

$$P(x \leq X(t + dt) \leq x + \xi) = \int_x^{x+\xi} \Pi(dx, x; dt, t) dx.$$

The previous definitions for mean, variance and correlation can be extended to the continuous time, vector-valued case by indexing the individual states:

$$\begin{aligned} E\{X(t)\} &= \begin{bmatrix} E\{X_1(t)\} \\ \vdots \\ E\{X_n(t)\} \end{bmatrix} =: \mu(t) \\ E\{(X(t) - \mu(t))(X(t) - \mu(t))^T\} &= \begin{bmatrix} E\{X_1(t)X_1(t)\} & \dots & E\{X_1(t)X_n(t)\} \\ & \ddots & \vdots \\ & & E\{X_n(t)X_n(t)\} \end{bmatrix} =: \Sigma(t) \\ E\{X(t)X^T(s)\} &= \begin{bmatrix} E\{X_1(t)X_1(s)\} & \dots & E\{X_1(t)X_n(s)\} \\ & \ddots & \vdots \\ & & E\{X_n(t)X_n(s)\} \end{bmatrix} =: R(t, s) \end{aligned}$$

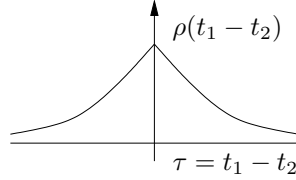


Figure 4.2: Correlation function for a first-order Markov process.

Note that the random variables and their statistical properties are all indexed by the time t (and s). The matrix $R(t, s)$ is called the *correlation matrix* for $X(t) \in \mathbb{R}^n$. If $t = s$ then $R(t, t)$ describes how the elements of x are correlated at time t (with each other) and in the case that the processes have zero mean, $R(t, t) = \Sigma(t)$. The elements on the diagonal of $\Sigma(t)$ are the variances of the corresponding scalar variables. A random process is uncorrelated if $R(t, s) = 0$ for all $t \neq s$. This implies that $X(t)$ and $X(s)$ are independent random events and is equivalent to $p_{X,Y}(x, y) = p_X(x)p_Y(y)$.

If a random process is stationary, then it can be shown that $R(t + \tau, s + \tau) = R(t, s)$ and it follows that the correlation matrix depends only on $t - s$. In this case we will often write $R(t, s) = R(s - t)$ or simple $R(\tau)$ where τ is the correlation time. The correlation matrix in this case is simply $R(0)$.

In the case where X is also scalar random process, the correlation matrix is also a scalar and we will write $\rho(\tau)$, which we refer to as the (scalar) correlation function. Furthermore, for stationary scalar random processes, the correlation function depends only on the absolute value of the correlation function, so $\rho(\tau) = \rho(-\tau) = \rho(|\tau|)$. This property also holds for the diagonal entries of the correlation matrix since $R_{ii}(s, t) = R_{ii}(t, s)$ from the definition.

Example 4.2 Ornstein-Uhlenbeck process

Consider a scalar random process defined by a Gaussian pdf with $\mu = 0$,

$$p(x, t) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2} \frac{x^2}{\sigma^2}},$$

and a correlation function given by

$$\rho(t_1, t_2) = \frac{Q}{2\omega_0} e^{-\omega_0 |t_2 - t_1|}.$$

The correlation function is illustrated in Figure 4.2. This process is also known as an *Ornstein-Uhlenbeck process*, a term that is commonly used in the scientific literature. This is a stationary process. ∇

The terminology and notation for covariance and correlation varies between disciplines. The term covariance is often used to refer to both the relationship between different variables X and Y and the relationship between a single variable at different times, $X(t)$ and $X(s)$. The term “cross-covariance” is used to refer to the covariance between two random vectors X and Y , to distinguish this from the covariance of the elements of X with each other. The term “cross-correlation” is

sometimes also used. Finally, the term “correlation coefficient” refers to the normalized correlation $\bar{\rho}(t, s) = \mathbb{E}[X(t)X(s)]/\mathbb{E}[X(t)X(t)]$.

MATLAB has a number of functions to implement covariance and correlation, which mostly match the terminology here:

- `cov(X)` - this returns the variance of the vector \mathbf{X} that represents samples of a given random variable or the covariance of the columns of a matrix X where the rows represent observations.
- `cov(X, Y)` - equivalent to `cov([X(:), Y(:)])`. Computes the covariance between the columns of X and Y , where the rows are observations.
- `xcorr(X, Y)` - the “cross-correlation” between two random sequences. If these sequences came from a random process, this is correlation function $\rho(t)$.
- `xcov(X, Y)` - this returns the “cross-covariance”, which MATLAB defines as the “mean-removed cross-correlation”.

The MATLAB help pages give the exact formulas used for each, so the main point here is to be careful to make sure you know what you really want.

We will also make use of a special type of random process referred to as “white noise”. A *white noise process* $X(t)$ satisfies $E\{X(t)\} = 0$ and $R(t, s) = W\delta(s - t)$, where $\delta(\tau)$ is the impulse function and W is called the *noise intensity*. White noise is an idealized process, similar to the impulse function or Heaviside (step) function in deterministic systems. In particular, we note that $\rho(0) = E\{X^2(t)\} = \infty$, so the covariance is infinite and we never see this signal in practice. However, like the step function, it is very useful for characterizing the response of a linear system, as described in the following proposition. It can be shown that the integral of a white noise process is a Wiener process, and so often white noise is described as the derivative of a Wiener process.

4.4 Linear Stochastic Systems with Gaussian Noise

We now consider the problem of how to compute the response of a linear system to a random process. We assume we have a linear system described in state space as

$$\dot{X} = AX + FW, \quad Y = CX \quad (4.13)$$

Given an “input” W , which is itself a random process with mean $\mu(t)$, variance $\sigma^2(t)$ and correlation $\rho(t, t + \tau)$, what is the description of the random process Y ?

Let W be a white noise process, with zero mean and noise intensity Q :

$$\rho(\tau) = Q\delta(\tau).$$

We can write the output of the system in terms of the convolution integral

$$Y(t) = \int_0^t h(t - \tau)W(\tau) d\tau,$$

where $h(t - \tau)$ is the impulse response for the system

$$h(t - \tau) = Ce^{A(t-\tau)}B + D\delta(t - \tau).$$

We now compute the statistics of the output, starting with the mean:

$$\begin{aligned}\mathbb{E}[Y(t)] &= E\left\{\int_0^t h(t-\eta)W(\eta) d\eta\right\} \\ &= \int_0^t h(t-\eta)E\{W(\eta)\} d\eta = 0.\end{aligned}$$

Note here that we have relied on the linearity of the convolution integral to pull the expectation inside the integral.

We can compute the covariance of the output by computing the correlation $\rho(\tau)$ and setting $\sigma^2 = \rho(0)$. The correlation function for y is

$$\begin{aligned}\rho_Y(t, s) &= E\{Y(t)Y(s)\} = E\left\{\int_0^t h(t-\eta)W(\eta) d\eta \cdot \int_0^s h(s-\xi)W(\xi) d\xi\right\} \\ &= E\left\{\int_0^t \int_0^s h(t-\eta)W(\eta)W(\xi)h(s-\xi) d\eta d\xi\right\}\end{aligned}$$

Once again linearity allows us to exchange expectation and integration

$$\begin{aligned}\rho_Y(t, s) &= \int_0^t \int_0^s h(t-\eta)E\{W(\eta)W(\xi)\}h(s-\xi) d\eta d\xi \\ &= \int_0^t \int_0^s h(t-\eta)Q\delta(\eta-\xi)h(s-\xi) d\eta d\xi \\ &= \int_0^t h(t-\eta)Qh(s-\eta) d\eta\end{aligned}$$

Now let $\tau = s - t$ and write

$$\begin{aligned}\rho_Y(\tau) &= \rho_Y(t, t + \tau) = \int_0^t h(t-\eta)Qh(t + \tau - \eta) d\eta \\ &= \int_0^t h(\xi)Qh(\xi + \tau) d\xi \quad (\text{setting } \xi = t - \eta)\end{aligned}$$

Finally, we let $t \rightarrow \infty$ (steady state)

$$\lim_{t \rightarrow \infty} \rho_Y(t, t + \tau) = \bar{\rho}_Y(\tau) = \int_0^\infty h(\xi)Qh(\xi + \tau) d\xi \quad (4.14)$$

If this integral exists, then we can compute the second order statistics for the output Y .

We can provide a more explicit formula for the correlation function ρ in terms of the matrices A , F and C by expanding equation (4.14). We will consider the general case where $W \in \mathbb{R}^m$ and $Y \in \mathbb{R}^p$ and use the correlation matrix $R(t, s)$ instead of the correlation function $\rho(t, s)$. Define the *state transition matrix* $\Phi(t, t_0) = e^{A(t-t_0)}$ so that the solution of system (4.13) is given by

$$x(t) = \Phi(t, t_0)x(t_0) + \int_{t_0}^t \Phi(t, \lambda)Fw(\lambda)d\lambda$$

Proposition 4.2 (Stochastic response to white noise). *Let $E\{X(t_0)X^T(t_0)\} = P(t_0)$ and W be white noise with $E\{W(\lambda)W^T(\xi)\} = R_W\delta(\lambda - \xi)$. Then the correlation matrix for X is given by*

$$R_X(t, s) = P(t)\Phi^T(s, t)$$

where $P(t)$ satisfies the linear matrix differential equation

$$\dot{P}(t) = AP + PA^T + FR_W F, \quad P(0) = P_0.$$

Proof. Using the definition of the correlation matrix, we have

$$\begin{aligned} E\{X(t)X^T(s)\} &= E\left\{\Phi(t, 0)X(0)X^T(0)\Phi^T(t, 0) + \text{cross terms}\right. \\ &\quad \left.+ \int_0^t \Phi(t, \xi)FW(\xi) d\xi \int_0^s W^T(\lambda)F^T\Phi(s, \lambda) d\lambda\right\} \\ &= \Phi(t, 0)E\{X(0)X^T(0)\}\Phi(s, 0) \\ &\quad + \int_0^t \int_0^s \Phi(t, \xi)FE\{W(\xi)W^T(\lambda)\}F^T\Phi(s, \lambda) d\xi d\lambda \\ &= \Phi(t, 0)P(0)\Phi^T(s, 0) + \int_0^t \Phi(t, \lambda)FR_W(\lambda)F^T\Phi(s, \lambda) d\lambda. \end{aligned}$$

Now use the fact that $\Phi(s, 0) = \Phi(s, t)\Phi(t, 0)$ (and similar relations) to obtain

$$R_X(t, s) = P(t)\Phi^T(s, t)$$

where

$$P(t) = \Phi(t, 0)P(0)\Phi^T(t, 0) + \int_0^t \Phi(t, \lambda)FR_W F^T(\lambda)\Phi^T(t, \lambda)d\lambda$$

Finally, differentiate to obtain

$$\dot{P}(t) = AP + PA^T + FR_W F, \quad P(0) = P_0$$

(see Friedland for details). □

The correlation matrix for the output Y can be computed using the fact that $Y = CX$ and hence $R_Y = C^T R_X C$. We will often be interested in the steady state properties of the output, which are given by the following proposition.

Proposition 4.3 (Steady state response to white noise). *For a time-invariant linear system driven by white noise, the correlation matrices for the state and output converge in steady state to*

$$R_X(\tau) = R_X(t, t + \tau) = Pe^{A^T \tau}, \quad R_Y(\tau) = CR_X(\tau)C^T$$

where P satisfies the algebraic equation

$$AP + PA^T + FR_W F^T = 0 \quad P > 0. \quad (4.15)$$

Equation (4.15) is called the *Lyapunov equation* and can be solved in MATLAB using the function `lyap`.

Example 4.3 First-order system

Consider a scalar linear process

$$\dot{X} = -aX + W, \quad Y = cX,$$

where W is a white, Gaussian random process with noise intensity σ^2 . Using the results of Proposition 4.2, the correlation function for X is given by

$$R_X(t, t + \tau) = p(t)e^{-a\tau}$$

where $p(t) > 0$ satisfies

$$p(t) = -2ap + \sigma^2.$$

We can solve explicitly for $p(t)$ since it is a (non-homogeneous) linear differential equation:

$$p(t) = e^{-2at}p(0) + (1 - e^{-2at})\frac{\sigma^2}{2a}.$$

Finally, making use of the fact that $Y = cX$ we have

$$\rho(t, t + \tau) = c^2(e^{-2at}p(0) + (1 - e^{-2at})\frac{\sigma^2}{2a})e^{-a\tau}.$$

In steady state, the correlation function for the output becomes

$$\rho(\tau) = \frac{c^2\sigma^2}{2a}e^{-a\tau}.$$

Note correlation function has the same form as the Ornstein-Uhlenbeck process in Example 4.2 (with $Q = c^2\sigma^2$). ∇

4.5 Random Processes in the Frequency Domain

As in the case of deterministic linear systems, we can analyze a stochastic linear system either in the state space or the frequency domain. The frequency domain approach provides a very rich set of tools for modeling and analysis of interconnected systems, relying on the frequency response and transfer functions to represent the flow of signals around the system.

Given a random process $X(t)$, we can look at the frequency content of the properties of the response. In particular, if we let $\rho(\tau)$ be the correlation function for a (scalar) random process, then we define the *power spectral density function* as the Fourier transform of ρ :

$$S(\omega) = \int_{-\infty}^{\infty} \rho(\tau)e^{-j\omega\tau} d\tau, \quad \rho(\tau) = \frac{1}{2\pi} \int_{-\infty}^{\infty} S(\omega)e^{j\omega\tau} d\omega.$$

The power spectral density provides an indication of how quickly the values of a random process can change through the frequency content: if there is high frequency content in the power spectral density, the values of the random variable can change quickly in time.

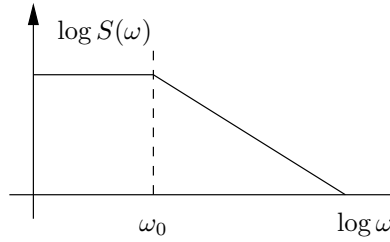


Figure 4.3: Spectral power density for a first-order Markov process.

Example 4.4 First-order Markov process

To illustrate the use of these measures, consider a first-order Markov process as defined in Example 4.2. The correlation function is

$$\rho(\tau) = \frac{Q}{2\omega_0} e^{-\omega_0|\tau|}.$$

The power spectral density becomes

$$\begin{aligned} S(\omega) &= \int_{-\infty}^{\infty} \frac{Q}{2\omega_0} e^{-\omega|\tau|} e^{-j\omega\tau} d\tau \\ &= \int_{-\infty}^0 \frac{Q}{2\omega_0} e^{(\omega-j\omega)\tau} d\tau + \int_0^{\infty} \frac{Q}{2\omega_0} e^{(-\omega-j\omega)\tau} d\tau = \frac{Q}{\omega^2 + \omega_0^2}. \end{aligned}$$

We see that the power spectral density is similar to a transfer function and we can plot $S(\omega)$ as a function of ω in a manner similar to a Bode plot, as shown in Figure 4.3. Note that although $S(\omega)$ has a form similar to a transfer function, it is a real-valued function and is not defined for complex s . ∇

Using the power spectral density, we can more formally define “white noise”: a *white noise process* is a zero-mean, random process with power spectral density $S(\omega) = W = \text{constant}$ for all ω . If $X(t) \in \mathbb{R}^n$ (a random vector), then $W \in \mathbb{R}^{n \times n}$. We see that a random process is white if all frequencies are equally represented in its power spectral density; this spectral property is the reason for the terminology “white”. The following proposition verifies that this formal definition agrees with our previous (time domain) definition.

Proposition 4.4. *For a white noise process,*

$$\rho(\tau) = \frac{1}{2\pi} \int_{-\infty}^{\infty} S(\omega) e^{j\omega\tau} d\omega = W\delta(\tau),$$

where $\delta(\tau)$ is the unit impulse function.

Proof. If $\tau \neq 0$ then

$$\rho(\tau) = \frac{1}{2\pi} \int_{-\infty}^{\infty} W(\cos(\omega\tau) + j \sin(\omega\tau)) d\omega = 0$$

If $\tau = 0$ then $\rho(\tau) = \infty$. Can show that

$$\rho(0) = \lim_{\epsilon \rightarrow 0} \int_{-\epsilon}^{\epsilon} \int_{-\infty}^{\infty} (\dots) d\omega d\tau = W\delta(0)$$

□

Given a linear system

$$\dot{X} = AX + FW, \quad Y = CX,$$

with W given by white noise, we can compute the spectral density function corresponding to the output Y . We start by computing the Fourier transform of the steady state correlation function (4.14):

$$\begin{aligned} S_Y(\omega) &= \int_{-\infty}^{\infty} \left[\int_0^{\infty} h(\xi) Q h(\xi + \tau) d\xi \right] e^{-j\omega\tau} d\tau \\ &= \int_0^{\infty} h(\xi) Q \left[\int_{-\infty}^{\infty} h(\xi + \tau) e^{-j\omega\tau} d\tau \right] d\xi \\ &= \int_0^{\infty} h(\xi) Q \left[\int_0^{\infty} h(\lambda) e^{-j\omega(\lambda - \xi)} d\lambda \right] d\xi \\ &= \int_0^{\infty} h(\xi) e^{j\omega\xi} d\xi \cdot QH(j\omega) = H(-j\omega)Q_uH(j\omega) \end{aligned}$$

This is then the (steady state) response of a linear system to white noise.

As with transfer functions, one of the advantages of computations in the frequency domain is that the composition of two linear systems can be represented by multiplication. In the case of the power spectral density, if we pass white noise through a system with transfer function $H_1(s)$ followed by transfer function $H_2(s)$, the resulting power spectral density of the output is given by

$$S_Y(\omega) = H_1(-j\omega)H_2(-j\omega)Q_uH_2(j\omega)H_1(j\omega).$$

As stated earlier, white noise is an idealized signal that is not seen in practice. One of the ways to produce more realistic models of noise and disturbances is to apply a filter to white noise that matches a measured power spectral density function. Thus, we wish to find a covariance W and filter $H(s)$ such that we match the statistics $S(\omega)$ of a measured noise or disturbance signal. In other words, given $S(\omega)$, find $W > 0$ and $H(s)$ such that $S(\omega) = H(-j\omega)WH(j\omega)$. This problem is known as the *spectral factorization problem*.

Figure 4.4 summarizes the relationship between the time and frequency domains.

4.6 Further Reading

There are several excellent books on stochastic systems that cover the results in this chapter in much more detail. For discrete-time systems, the textbook by Kumar and Varaiya [KV86] provides an derivation of the key results. Results for continuous-time systems can be found in the textbook by Friedland [Fri04]. Åström [Åst06] gives a very elegant derivation in a unified framework that integrates discrete-time and continuous-time systems.

$$\begin{array}{ccc}
 p(v) = \frac{1}{\sqrt{2\pi R_V}} e^{-\frac{v^2}{2R_V}} & V \longrightarrow \boxed{H} \longrightarrow Y & p(y) = \frac{1}{\sqrt{2\pi R_Y}} e^{-\frac{y^2}{2R_Y}} \\
 S_V(\omega) = R_V & & S_Y(\omega) = H(-j\omega)R_V H(j\omega) \\
 \rho_V(\tau) = R_V \delta(\tau) & \dot{X} = AX + FV & \rho_Y(\tau) = R_Y(\tau) = CPe^{-A|\tau|}C^T \\
 & Y = CX & AP + PA^T + FR_VF^T = 0
 \end{array}$$

Figure 4.4: Summary of steady state stochastic response.

Exercises

4.1 Let Z be a random variable that is the sum of two independent normally (Gaussian) distributed random variables X_1 and X_2 having means m_1 , m_2 and variances σ_1^2 , σ_2^2 respectively. Show that the probability density function for Z is

$$p(z) = \frac{1}{2\pi\sigma_1\sigma_2} \int_{-\infty}^{\infty} \exp\left\{-\frac{(z-x-m_1)^2}{2\sigma_1^2} - \frac{(x-m_2)^2}{2\sigma_2^2}\right\} dx$$

and confirm that this is normal (Gaussian) with mean m_1+m_2 and variance $\sigma_1^2+\sigma_2^2$. (Hint: Use the fact that $p(z|x_2) = p_{X_1}(x_1) = p_{X_1}(z-x_2)$.)

4.2 (ÅM08, Exercise 7.13) Consider the motion of a particle that is undergoing a random walk in one dimension (i.e., along a line). We model the position of the particle as

$$x[k+1] = x[k] + u[k],$$

where x is the position of the particle and u is a white noise processes with $E\{u[i]\} = 0$ and $E\{u[i]u[j]\} = R_u\delta(i-j)$. We assume that we can measure x subject to additive, zero-mean, Gaussian white noise with covariance 1. Show that the expected value of the particle as a function of k is given by

$$E\{x[k]\} = E\{x[0]\} + \sum_{i=0}^{k-1} E\{u[i]\} = E\{x[0]\} =: \mu_x$$

and the covariance $E\{(x[k] - \mu_x)^2\}$ is given by

$$E\{(x[k] - \mu_x)^2\} = \sum_{i=0}^{k-1} E\{u^2[i]\} = kR_u$$

4.3 Consider a second order system with dynamics

$$\begin{bmatrix} \dot{X}_1 \\ \dot{X}_2 \end{bmatrix} = \begin{bmatrix} -a & 0 \\ 0 & -b \end{bmatrix} \begin{bmatrix} X_1 \\ X_2 \end{bmatrix} + \begin{bmatrix} 1 \\ 1 \end{bmatrix} v, \quad Y = [1 \quad 1] \begin{bmatrix} X_1 \\ X_2 \end{bmatrix}$$

that is forced by Gaussian white noise with zero mean and variance σ^2 . Assume $a, b > 0$.

(a) Compute the correlation function $\rho(\tau)$ for the output of the system. Your answer should be an explicit formula in terms of a , b and σ .

(b) Assuming that the input transients have died out, compute the mean and variance of the output.

4.4 Find a constant matrix A and vectors F and C such that for

$$\dot{X} = AX + FW, \quad Y = CX$$

the power spectrum of Y is given by

$$S(\omega) = \frac{1 + \omega^2}{(1 - 7\omega^2)^2 + 1}$$

Describe the sense in which your answer is unique.

