# 2

## Markov Models

The results presented in this book have been written in the desire that practitioners will use them. We have tried therefore to illustrate the use of the theory in a systematic and accessible way, and so this book concentrates not only on the theory of general space Markov chains, but on the application of that theory in considerable detail.

We will apply the results which we develop across a range of specific applications: typically, after developing a theoretical construct, we apply it to models of increasing complexity in the areas of systems and control theory, both linear and nonlinear, both scalar and vector-valued; traditional "applied probability" or operations research models, such as random walks, storage and queueing models, and other regenerative schemes; and models which are in both domains, such as classical and recent time-series models.

These are not given merely as "examples" of the theory: in many cases, the application is difficult and deep of itself, whilst applications across such a diversity of areas have often driven the definition of general properties and the links between them. Our goal has been to develop the analysis of applications on a step by step basis as the theory becomes richer throughout the book.

To motivate the general concepts, then, and to introduce the various areas of application, we leave until Chapter 3 the normal and necessary foundations of the subject, and first introduce a cross-section of the models for which we shall be developing those foundations.

These models are still described in a somewhat heuristic way. The full mathematical description of their dynamics must await the development in the next chapter of the concepts of transition probabilities, and the reader may on occasion benefit by moving to some of those descriptions in parallel with the outlines here.

It is also worth observing immediately that the descriptive definitions here are from time to time supplemented by other assumptions in order to achieve specific results: these assumptions, and those in this chapter and the last, are collected for ease of reference in Appendix C.

As the definitions are developed, it will be apparent immediately that very many of these models have a random additive component, such as the i.i.d. sequence $\{W_n\}$ in both the linear state space model and the random walk model. Such a component goes by various names, such as error, noise, innovation, disturbance or increment sequence, across the various model areas we consider. We shall use the nomenclature relevant to the context of each model.

We will save considerable repetitive definition if we adopt a global convention immediately to cover these sequences.

---

**Error, Noise, Innovation, Disturbance and Increments**

Suppose $\mathbf{W} = \{W_n\}$ is labeled as an error, noise, innovation, disturbance or increment sequence. Then this has the interpretation that the random variables $\{W_n\}$ are independent and identically distributed, with distribution identical to that of a generic variable denoted $W$.

We will systematically denote the probability law of such a variable $W$ by $\Gamma$.

---

It will also be apparent that many models are defined inductively from their own past in combination with such innovation sequences. In order to commence the induction, initial values are needed. We adopt a second convention immediately to avoid repetition in defining our models.

---

**Initialization**

Unless specifically defined otherwise, the initial state $\{\Phi_0\}$ of a Markov model will be taken as independent of the error, noise, innovation, disturbance or increments process, and will have an arbitrary distribution.

---

## 2.1 Markov Models In Time Series

The theory of time series has been developed to model a set of observations developing in time: in this sense, the fundamental starting point for time series and for more general Markov models is virtually identical. However, whilst the Markov theory immediately assumes a short-term dependence structure on the variables at each time point, time series theory concentrates rather on the parametric form of dependence between the variables.

The time series literature has historically concentrated on linear models (that is, those for which past disturbances and observations are combined to form the present

observation through some linear transformation) although recently there has been greater emphasis on nonlinear models. We first survey a number of general classes of linear models and turn to some recent nonlinear time series models in Section 2.2.

It is traditional to denote time series models as a sequence $\mathbf{X} = \{X_n : n \in \mathbb{Z}_+\}$, and we shall follow this tradition.

### 2.1.1 Simple linear models

The first class of models we discuss has direct links with deterministic linear models, state space models and the random walk models we have already introduced in Chapter 1.

We begin with the simplest possible "time series" model, the scalar autoregression of order one, or AR(1) model on $\mathbb{R}^1$.

---

**Simple Linear Model**

The process $\mathbf{X} = \{X_n, n \in \mathbb{Z}_+\}$ is called the *simple linear model*, or *AR(1) model* if

(SLM1)    for each $n \in \mathbb{Z}_+$, $X_n$ and $W_n$ are random variables on $\mathbb{R}$, satisfying
$$X_{n+1} = \alpha X_n + W_{n+1},$$
for some $\alpha \in \mathbb{R}$;

(SLM2)    the random variables $\{W_n\}$ are an error sequence with distribution $\Gamma$ on $\mathbb{R}$.

---

The simple linear model is trivially Markovian: the independence of $X_{n+1}$ from $X_{n-1}, X_{n-2}, \ldots$ given $X_n = x$ follows from the construction rule (SLM1), since the value of $W_n$ does not depend on any of $\{X_{n-1}, X_{n-2} \ldots\}$ from (SLM2).

The simple linear model can be viewed in one sense as an extension of the random walk model, where now we take some proportion or multiple of the previous value, not necessarily equal to the previous value, and again add a new random amount (the "noise" or "error") onto this scaled random value. Equally, it can be viewed as the simplest special case of the linear state space model LSS($F$,$G$), in the scalar case with $F = \alpha$ and $G = 1$.

In Figure 2.1 and Figure 2.2 we give sets of sample paths of linear models with different values of the parameter $\alpha$.

The choice of this parameter critically determines the behavior of the chain. If $|\alpha| < 1$ then the sample paths remain bounded in ways which we describe in detail in
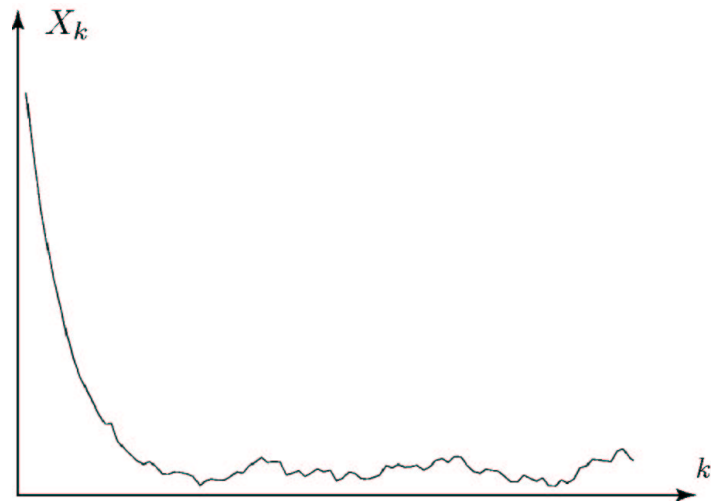
**Fig. 2.1.** Linear model path with $\alpha = 0.85$, increment distribution $N(0,1)$
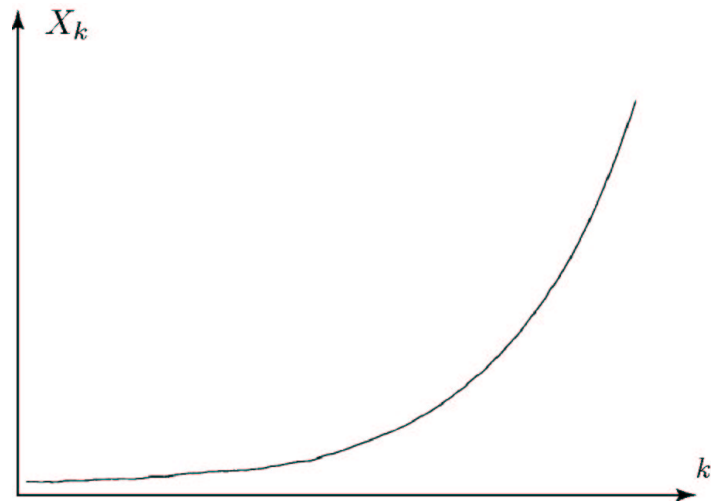


**Fig. 2.2.** Linear model path with $\alpha = 1.05$, increment distribution $N(0,1)$

later chapters, and the process $\mathbf{X}$ is inherently "stable": in fact, ergodic in the sense of Section 1.3.1 (III) and Theorem 1.3.1, for reasonable distributions $\Gamma$. But if $|\alpha| > 1$ then $\mathbf{X}$ is unstable, in a well-defined way: in fact, evanescent with probability one, in the sense of Section 1.3.1 (II), if the noise distribution $\Gamma$ is again reasonable.

### 2.1.2 Linear autoregressions and ARMA models

In the development of time series theory, simple linear models are usually analyzed as a subset of the class of autoregressive models, which depend in a linear manner on their past history for a fixed number $k \geq 1$ of steps in the past.

---

**Autoregressive Model**

A process $\mathbf{Y} = \{Y_n\}$ is called a (scalar) *autoregression of order $k$*, or AR($k$) model, if it satisfies, for each set of initial values $(Y_0, \ldots, Y_{-k+1})$,

(AR1)    for each $n \in \mathbb{Z}_+$, $Y_n$ and $W_n$ are random variables on $\mathbb{R}$ satisfying inductively for $n \geq 1$

$$Y_n = \alpha_1 Y_{n-1} + \alpha_2 Y_{n-2} + \ldots + \alpha_k Y_{n-k} + W_n,$$

for some $\alpha_1, \ldots, \alpha_k \in \mathbb{R}$;

(AR2)    the sequence $\mathbf{W}$ is an error sequence on $\mathbb{R}$.

---

The collection $\mathbf{Y} = \{Y_n\}$ is generally not Markovian if $k > 1$, since information on the past (or at least the past in terms of the variables $Y_{n-1}, Y_{n-2}, \ldots, Y_{n-k}$) provides information on the current value $Y_n$ of the process. But by the device mentioned in Section 1.2.1, of constructing the multivariate sequence

$$X_n = (Y_n, \ldots, Y_{n-k+1})^\top$$

and setting $\mathbf{X} = \{X_n, n \geq 0\}$, we define $\mathbf{X}$ as a Markov chain whose first component has exactly the sample paths of the autoregressive process. Note that the general convention that $X_0$ has an arbitrary distribution implies that the first $k$ variables $(Y_0, \ldots, Y_{-k+1})$ are also considered arbitrary.

The autoregressive model can then be viewed as a specific version of the vector-valued linear state space model LSS($F$,$G$). For by (AR1),

$$X_n = \begin{bmatrix} \alpha_1 & \cdots & \cdots & \alpha_k \\ 1 & & & 0 \\ & \ddots & & \vdots \\ 0 & & 1 & 0 \end{bmatrix} X_{n-1} + \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} W_n. \tag{2.1}$$

The same technique for producing a Markov model can be used for any linear model which admits a finite dimensional description. In particular, we take the following general model:

---

### Autoregressive-Moving Average Models

The process $\mathbf{Y} = \{Y_n\}$ is called an *autoregressive-moving average process of order* $(k, \ell)$, or ARMA$(k, \ell)$ model, if it satisfies, for each set of initial values $(Y_0, \ldots, Y_{-k+1}, W_0, \ldots, W_{-\ell+1})$,

(ARMA1)    for each $n \in \mathbb{Z}_+$, $Y_n$ and $W_n$ are random variables on $\mathbb{R}$, satisfying, inductively for $n \geq 1$,

$$Y_n \;\; = \;\; \alpha_1 Y_{n-1} + \alpha_2 Y_{n-2} + \ldots + \alpha_k Y_{n-k}$$

$$+ W_n + \beta_1 W_{n-1} + \beta_2 W_{n-2} + \ldots + \beta_\ell W_{n-\ell},$$

for some $\alpha_1, \ldots, \alpha_k, \beta_1, \ldots, \beta_\ell \in \mathbb{R}$;

(ARMA2)    the sequence $\mathbf{W}$ is an error sequence on $\mathbb{R}$.

---

In this case more care must be taken to obtain a suitable Markovian description of the process. One approach is to take

$$X_n = (Y_n, \ldots, Y_{n-k+1}, W_n, \ldots, W_{n-\ell+1})^\top.$$

Although the resulting state process $\mathbf{X}$ is Markovian, the dimension of this realization may be overly large for effective analysis. A realization of lower dimension may be obtained by defining the stochastic process $\mathbf{Z}$ inductively by

$$Z_n = \alpha_1 Z_{n-1} + \alpha_2 Z_{n-2} + \ldots + \alpha_k Z_{n-k} + W_n. \tag{2.2}$$

When the initial conditions are defined appropriately, it is a matter of simple algebra and an inductive argument to show that

$$Y_n = Z_n + \beta_1 Z_{n-1} + \beta_2 Z_{n-2} + \ldots + \beta_\ell Z_{n-\ell},$$

Hence the probabilistic structure of the ARMA$(k, \ell)$ process is completely determined by the Markov chain $\{(Z_n, \ldots, Z_{n-k+1})^\top : n \in \mathbb{Z}_+\}$ which takes values in $\mathbb{R}^k$.

The behavior of the general ARMA$(k, \ell)$ model can thus be placed in the Markovian context, and we will develop the stability theory of this, and more complex versions of this model, in the sequel.

## 2.2  Nonlinear State Space Models

In discrete time, a general (semi) dynamical system on $\mathbb{R}$ is defined, as in Section 1.3.2, through a recursion of the form

$$x_{n+1} = F(x_n), \qquad n \in \mathbb{Z}_+ \tag{2.3}$$

for some continuous function $F: \mathbb{R} \to \mathbb{R}$. Hence the simple linear model defined in (SLM1) may be interpreted as a linear dynamical system perturbed by the "noise" sequence $\mathbf{W}$.

The theory of time series is in this sense closely related to the general theory of dynamical systems: it has developed essentially as that subset of stochastic dynamical systems theory for which the relationships between the variables are linear, and even with the nonlinear models from the time series literature which we consider below, there is still a large emphasis on linear substructures.

The theory of dynamical systems, in contrast to time series theory, has grown from a deterministic base, considering initially the type of linear relationship in (1.3) with which we started our examples in Section 1.2, but progressing to models allowing a very general (but still deterministic) relationship between the variables in the present and in the past, as in (2.3). It is in the more recent development that "noise" variables, allowing the system to be random in some part of its evolution, have been introduced.

*Nonlinear state space models* are stochastic versions of dynamical systems where a Markovian realization of the model is both feasible and explicit: thus they satisfy a generalization of (2.3) such as

$$X_{n+1} = F(X_n, W_{n+1}), \qquad k \in \mathbb{Z}_+ \tag{2.4}$$

where $\mathbf{W}$ is a noise sequence and the function $F: \mathbb{R}^n \times \mathbb{R}^p \to \mathbb{R}^n$ is smooth ($C^\infty$): that is, all derivatives of $F$ exist and are continuous.

### 2.2.1  Scalar nonlinear models

We begin with the simpler version of (2.4) in which the random variables are scalar.

---

**Scalar Nonlinear State Space Model**

The chain $\mathbf{X} = \{X_n\}$ is called a scalar *nonlinear state space model on* $\mathbb{R}$ *driven by* $F$, or SNSS($F$) model, if it satisfies

(SNSS1)   for each $n \geq 0$, $X_n$ and $W_n$ are random variables on $\mathbb{R}$, satisfying, inductively for $n \geq 1$,

$$X_n = F(X_{n-1}, W_n),$$

for some smooth $(C^\infty)$ function $F : \mathbb{R} \times \mathbb{R} \to \mathbb{R}$;

(SNSS2)   the sequence $\mathbf{W}$ is a disturbance sequence on $\mathbb{R}$, whose marginal distribution $\Gamma$ possesses a density $\gamma_w$ supported on an open set $O_w$.

---

The independence of $X_{n+1}$ from $X_{n-1}, X_{n-2}, \ldots$ given $X_n = x$ follows from the rules (SNSS1) and (SNSS2), and ensures as previously that $\mathbf{X}$ is a Markov chain.

As with the linear control model (LCM1) associated with the linear state space model (LSS1), we will analyze nonlinear state space models through the associated deterministic "control models". Define the sequence of maps $\{F_k : \mathbb{R} \times \mathbb{R}^k \to \mathbb{R} : k \geq 0\}$ inductively by setting $F_0(x) = x$, $F_1(x_0, u_1) = F(x_0, u_1)$ and for $k > 1$

$$F_k(x_0, u_1, \ldots, u_k) = F(F_{k-1}(x_0, u_1, \ldots, u_{k-1}), u_k). \qquad (2.5)$$

We call the deterministic system with trajectories

$$x_k = F_k(x_0, u_1, \ldots, u_k), \qquad k \in \mathbb{Z}_+ \qquad (2.6)$$

the *associated control model* CM($F$) for the SNSS($F$) model, provided the deterministic control sequence $\{u_1, \ldots, u_k, k \in \mathbb{Z}_+\}$ lies in the set $O_w$, which we call the *control set* for the scalar nonlinear state space model.

To make these definitions more concrete we define two particular classes of scalar nonlinear models with specific structure which we shall use as examples on a number of occasions.

The first of these is the *bilinear model*, so called because it is linear in each of its input variables, namely the immediate past of the process and a noise component, whenever the other is fixed: but their joint action is multiplicative as well as additive.
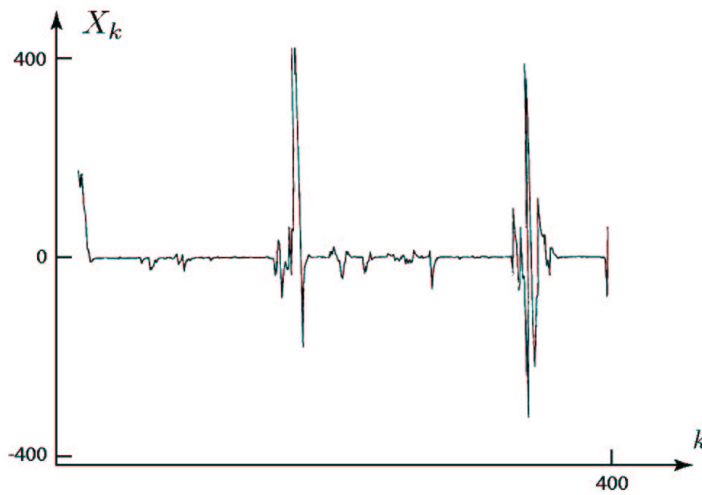
**Fig. 2.3.** Simple bilinear model path with $F(x, w) = (0.707 + w)x + w$

---

Simple Bilinear Model

The chain $\mathbf{X} = \{X_n\}$ is called the *simple bilinear model* if it satisfies

(SBL1)    for each $n \geq 0$, $X_n$ and $W_n$ are random variables on $\mathbb{R}$, satisfying for $n \geq 1$,

$$X_n = \theta X_{n-1} + b X_{n-1} W_n + W_n$$

where $\theta$ and $b$ are scalars, and the sequence $\mathbf{W}$ is an error sequence on $\mathbb{R}$.

---

The bilinear process is thus a SNSS($F$) model with $F$ given by

$$F(x, w) = \theta x + bxw + w, \qquad (2.7)$$

where the control set $O_w \subseteq \mathbb{R}$ depends upon the specific distribution of $W$.

In Figure 2.3 we give a sample path of a scalar nonlinear model with

$$F(x, w) = (0.707 + w)x + w$$

and with $\Gamma = N(0, \frac{1}{2})$. This is the simple bilinear model with $\theta = 0.707$ and $b = 1$. One can see from this simulation that the behavior of this model is quite different from that of any linear model.

The second specific nonlinear model we shall analyze is the scalar first-order *SETAR model*. This is piecewise linear in contiguous regions of $\mathbb{R}$, and thus while it may serve as an approximation to a completely nonlinear process, we shall see that much of its analysis is still tractable because of the linearity of its component parts.

---

SETAR Models

The chain $\mathbf{X} = \{X_n\}$ is called a *scalar self-exciting threshold autoregression (SETAR) model* if it satisfies

(SETAR1)    for each $1 \leq j \leq M$, $X_n$ and $W_n(j)$ are random variables on $\mathbb{R}$, satisfying, inductively for $n \geq 1$,

$$X_n = \phi(j) + \theta(j)X_{n-1} + W_n(j), \qquad r_{j-1} < X_{n-1} \leq r_j,$$

where $-\infty = r_0 < r_1 < \cdots < r_M = \infty$ and $\{W_n(j)\}$ forms an i.i.d. zero-mean error sequence for each $j$, independent of $\{W_n(i)\}$ for $i \neq j$.

---

Because of lack of continuity, the SETAR models do not fall into the class of nonlinear state space models, although they can often be analyzed using essentially the same methods. The SETAR model will prove to be a useful example on which to test the various stability criteria we develop, and the overall outcome of that analysis is gathered together in Section B.2.

### 2.2.2  Multi-dimensional nonlinear models

Many nonlinear processes cannot be modeled by a scalar Markovian model such as the SNSS($F$) model. The more general multi-dimensional model is defined quite analogously.

Nonlinear State Space Models

Suppose $\mathbf{X} = \{X_k\}$, where

(NSS1)   for each $k \geq 0$ $X_k$ and $W_k$ are random variables on $\mathbb{R}^n$, $\mathbb{R}^p$ respectively, satisfying inductively for $k \geq 1$,

$$X_k = F(X_{k-1}, W_k),$$

for some smooth $(C^\infty)$ function $F: \mathsf{X} \times O_w \to \mathsf{X}$, where $\mathsf{X}$ is an open subset of $\mathbb{R}^n$, and $O_w$ is an open subset of $\mathbb{R}^p$;

(NSS2)   the random variables $\{W_k\}$ are a disturbance sequence on $\mathbb{R}^p$, whose marginal distribution $\Gamma$ possesses a density $\gamma_w$ which is supported on an open set $O_w$.

Then $\mathbf{X}$ is called a *nonlinear state space model driven by $F$, or NSS(F) model, with control set $O_w$*.

The general nonlinear state space model can often be analyzed by the same methods that are used for the scalar SNSS($F$) model, under appropriate conditions on the disturbance process $\mathbf{W}$ and the function $F$.

It is a central observation of such analysis that the structure of the NSS($F$) model (and of course its scalar counterpart) is governed under suitable conditions by an associated deterministic control model, defined analogously to the linear control model and the linear state space model.

The Associated Control Model CM($\boldsymbol{F}$)

(CM1)   The deterministic system

$$x_k = F_k(x_0, u_1, \ldots, u_k), \qquad k \in \mathbb{Z}_+, \qquad (2.8)$$

where the sequence of maps $\{F_k : \mathsf{X} \times O_w^k \to \mathsf{X} : k \geq 0\}$ is defined by (2.5), is called the *associated control system* for the NSS($F$) model and is denoted CM($F$) provided the deterministic control sequence $\{u_1, \ldots, u_k, k \in \mathbb{Z}_+\}$ lies in the control set $O_w \subseteq \mathbb{R}^p$.

The general ARMA model may be generalized to obtain a class of nonlinear models, all of which may be "Markovianized", as in the linear case.

---

**Nonlinear Autoregressive-Moving Average Models**

The process $\mathbf{Y} = \{Y_n\}$ is called a *nonlinear autoregressive-moving average process of order* $(k, \ell)$ if the values $Y_0, \ldots, Y_{k-1}$ are arbitrary and

(NARMA1)   for each $n \geq 0$, $Y_n$ and $W_n$ are random variables on $\mathbb{R}$, satisfying, inductively for $n \geq k$,

$$Y_n = G(Y_{n-1}, Y_{n-2}, \ldots, Y_{n-k}, W_n, W_{n-1}, W_{n-2}, \ldots, W_{n-\ell})$$

where the function $G \colon \mathbb{R}^{k+\ell+1} \to \mathbb{R}$ is smooth $(C^\infty)$;

(NARMA2)   the sequence $\mathbf{W}$ is an error sequence on $\mathbb{R}$.

---

As in the linear case, we may define

$$X_n = (Y_n, \ldots, Y_{n-k+1}, W_n, \ldots, W_{n-\ell+1})^\top$$

to obtain a Markovian realization of the process $\mathbf{Y}$. The process $\mathbf{X}$ is Markovian, with state space $\mathsf{X} = \mathbb{R}^{k+\ell}$, and has the general form of an NSS($F$) model, with

$$X_n = F(X_{n-1}, W_n), \qquad n \in \mathbb{Z}_+. \tag{2.9}$$

### 2.2.3  The gumleaf attractor

The gumleaf attractor is an example of a nonlinear model such as those which frequently occur in the analysis of control algorithms for nonlinear systems, some of which are briefly described below in Section 2.3. In an investigation of the pathologies which can reveal themselves in adaptive control, a specific control methodology which is described in Section 2.3.2, Mareels and Bitmead [161] found that the closed loop system dynamics in an adaptive control application can be described by the simple recursion

$$v_n = -\frac{1}{v_{n-1}} + \frac{1}{v_{n-2}}, \qquad n \in \mathbb{Z}_+.$$

Here $v_n$ is a "closed loop system gain" which is a simple function of the output of the system which is to be controlled. By setting $x_n = \begin{pmatrix} x_n^a \\ x_n^b \end{pmatrix} = \begin{pmatrix} v_n \\ v_{n-1} \end{pmatrix}$ we obtain a nonlinear state space model with
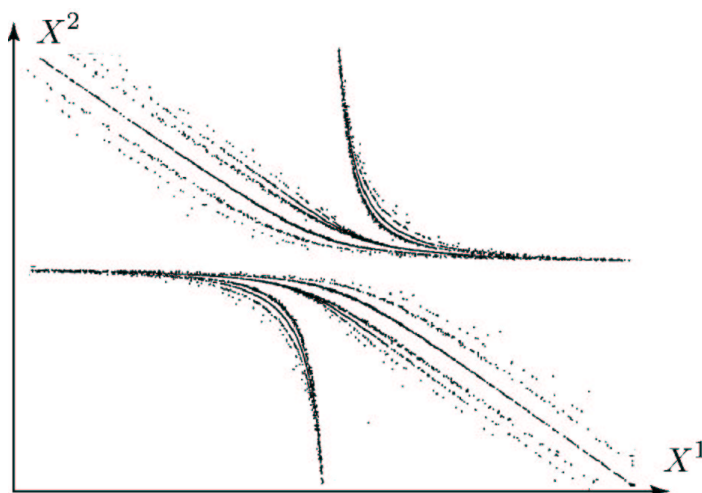
**Fig. 2.4.** The gumleaf attractor

$$F\begin{pmatrix} x^a \\ x^b \end{pmatrix} = \begin{pmatrix} -1/x^a + 1/x^b \\ x^a \end{pmatrix}$$

so that

$$x_n = \begin{pmatrix} x_n^a \\ x_n^b \end{pmatrix} = F\begin{pmatrix} x_{n-1}^a \\ x_{n-1}^b \end{pmatrix} = \begin{pmatrix} -1/x_{n-1}^a + 1/x_{n-1}^b \\ x_{n-1}^a \end{pmatrix} \tag{2.10}$$

If $F$ is required to be continuous then the state space $\mathsf{X}$ in this example must be taken as two dimensional Euclidean space $\mathbb{R}^2$ minus the $x$ and $y$ axes, and any other initial conditions which might result in a zero value for $x_n^a$ or $x_n^b$ for some $n$.

A typical sample path of this model is given in Figure 2.4. In this figure 40,000 consecutive sample points of $\{x_n\}$ have been indicated by points to illustrate the qualitative behavior of the model. Because of its similarity to some Australian flora, the authors call the resulting plot the *gumleaf attractor*. Ydstie in [285] also finds that such chaotic behavior can easily occur in adaptive systems.

One way that noise can enter the model (2.10) is directly through the first component $x_n^a$ to give

$$X_n = \begin{pmatrix} X_n^a \\ X_n^b \end{pmatrix} = F\begin{pmatrix} X_{n-1}^a \\ X_{n-1}^b \end{pmatrix} = \begin{pmatrix} -1/X_{n-1}^a + 1/X_{n-1}^b \\ X_{n-1}^a \end{pmatrix} + \begin{pmatrix} W_n \\ 0 \end{pmatrix} \tag{2.11}$$

where $\mathbf{W}$ is i.i.d..

The special case where for each $n$ the disturbance $W_n$ is uniformly distributed on $[-\frac{1}{2}, \frac{1}{2}]$ is illustrated in Figure 2.5. As in the previous figure, we have plotted 40,000 values of the sequence $\mathbf{X}$ which takes values in $\mathbb{R}^2$. Note that the qualitative behavior of the process remains similar to the noise-free model, although some of the detailed behavior is "smeared out" by the noise.

The analysis of general models of this type is a regular feature in what follows, and in Chapter 7 we give a detailed analysis of the path structure that might be expected under suitable assumptions on the noise and the associated deterministic model.
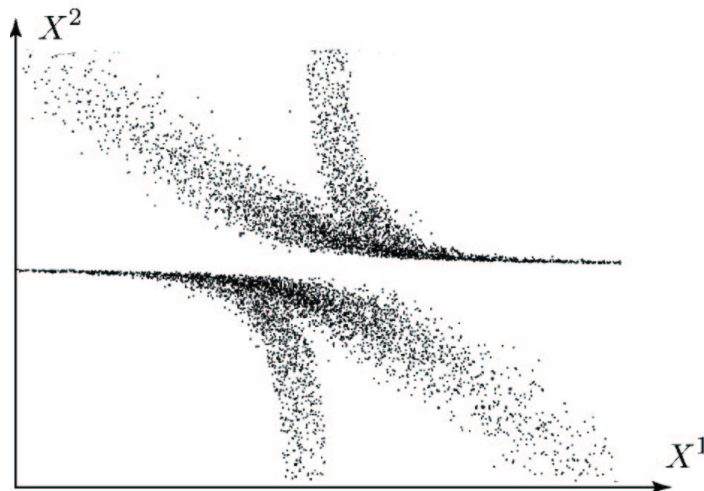
**Fig. 2.5.** The gumleaf attractor perturbed by noise

### 2.2.4 The dependent parameter bilinear model

As a simple example of a multidimensional nonlinear state space model, we will consider the following dependent parameter bilinear model, which is closely related to the simple bilinear model introduced above. To allow for dependence in the parameter process, we construct a two dimensional process so that the Markov assumption will remain valid.

---

The Dependent Parameter Bilinear Model

The process $\boldsymbol{\Phi} = \binom{\boldsymbol{\theta}}{\mathbf{Y}}$ is called the *dependent parameter bilinear model* if it satisfies

(DBL1)     For some $|\alpha| < 1$ and all $k \in \mathbb{Z}_+$,

$$
\begin{aligned}
Y_{k+1} &= \theta_k Y_k + W_{k+1} & (2.12) \\
\theta_{k+1} &= \alpha \theta_k + Z_{k+1}, & (2.13)
\end{aligned}
$$

(DBL2)     The joint process $(\mathbf{Z}, \mathbf{W})^{\top}$ is a disturbance sequence on $\mathbb{R}^2$, $\mathbf{Z}$ and $\mathbf{W}$ are mutually independent, and the distributions $\Gamma_w$ and $\Gamma_z$ of $W$, $Z$ respectively possess densities which are lower semicontinuous. It is assumed that $W$ has a finite second moment, and that $\mathsf{E}[\log(1 + |Z|)] < \infty$.
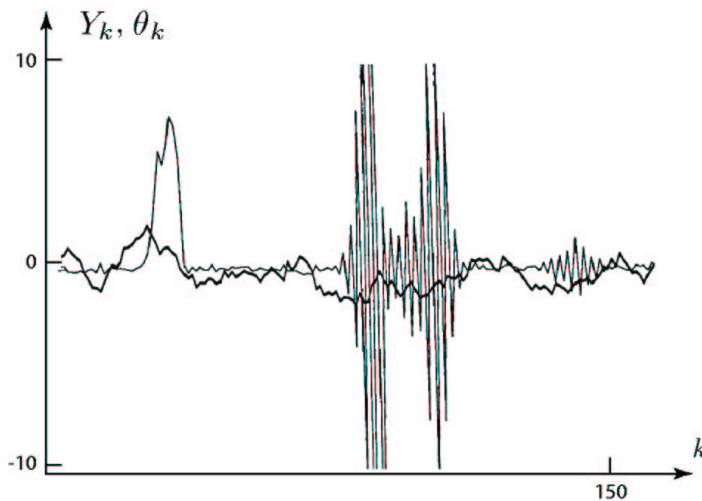
---

**Fig. 2.6.** Dependent parameter bilinear model paths with $\alpha = 0.933$, $W_k \sim N(0, 0.14)$ and $Z_k \sim N(0, 0.01)$

This is described by a two dimensional $\mathrm{NSS}(F)$ model, where the function $F$ is of the form

$$F\left(\begin{pmatrix} Y \\ \theta \end{pmatrix}, \begin{pmatrix} Z \\ W \end{pmatrix}\right) = \begin{pmatrix} \alpha\theta + Z \\ \theta Y + W \end{pmatrix} \tag{2.14}$$

As usual, the control set $O_w \subseteq \mathbb{R}^2$ depends upon the specific distribution of $W$ and $Z$.

A plot of the joint process $\begin{pmatrix} \mathbf{Y} \\ \boldsymbol{\theta} \end{pmatrix}$ is given in Figure 2.6. In this simulation we have $\alpha = 0.933$, $W_k \sim N(0, 0.14)$ and $Z_k \sim N(0, 0.01)$.

The dark line is a plot of the parameter process $\boldsymbol{\theta}$, and the lighter, more explosive path is the resulting output $\mathbf{Y}$. One feature of this model is that the output oscillates rapidly when $\theta_k$ takes on large negative values, which occurs in this simulation for time values between 80 and 100.

## 2.3 Models In Control And Systems Theory

### 2.3.1 Choosing controls

In Section 2.2, we defined deterministic control systems, such as (2.5), associated with Markovian state space models. We now begin with a general control system, which might model the dynamics of an aircraft, a cruise control in an automobile, or a controlled chemical reaction, and seek ways to choose a control to make the system attain a desired level of performance.

Such control laws typically involve feedback; that is, the input at a given time is chosen based upon present output measurements, or other features of the system

which are available at the time that the control is computed. Once such a control law has been selected, the dynamics of the controlled system can be complex. Fortunately, with most control laws, there is a representation (the "closed loop" system equations) which gives rise to a Markovian state process $\boldsymbol{\Phi}$ describing the variables of interest in the system. This additional structure can greatly simplify the analysis of control systems.

We can extend the AR models of time series to an ARX (*autoregressive with exogenous variables*) system model defined for $k \geq 1$ by

$$Y_k + \alpha_1(k)Y_{k-1} + \cdots + \alpha_{n_1}(k)Y_{k-n_1} = \beta_1(k)U_{k-1} + \cdots + \beta_{n_2}(k)U_{k-n_2} + W_k \quad (2.15)$$

where we assume for this discussion that the output process $\mathbf{Y}$, the input process (or exogenous variable sequence) $\mathbf{U}$, and the disturbance process $\mathbf{W}$ are all scalar-valued, and initial conditions are assigned at $k = 0$.

Let us also assume that we have random coefficients $\alpha_j(k), \beta_j(k)$ rather than fixed coefficients at each time point $k$. In such a case we may have to estimate the coefficients in order to choose the exogenous input $\mathbf{U}$.

The objective in the design of the control sequence $\mathbf{U}$ is specific to the particular application. However, it is often possible to set up the problem so that the goal becomes a problem of regulation: that is, to make the output as small as possible. Given the stochastic nature of systems, this is typically expressed using the concepts of sample mean square stabilizing sequences and minimum variance control laws.

We call the input sequence $\mathbf{U}$ *sample mean square stabilizing* if the input-output process satisfies

$$\limsup_{N \to \infty} \frac{1}{N} \sum_{k=1}^{N} [Y_k^2 + U_k^2] < \infty \qquad \text{a.s.}$$

for every initial condition. The control law is then said to be *minimum variance* if it is sample mean square stabilizing, and the sample path average

$$\limsup_{N \to \infty} \frac{1}{N} \sum_{k=1}^{N} Y_k^2 \qquad\qquad (2.16)$$

is minimized over all control laws with the property that, for each $k$, the input $U_k$ is a function of $Y_k, \ldots, Y_0$, and the initial conditions.

Such controls are often called "causal", and for causal controls there is some possibility of a Markovian representation. We now specialize this general framework to a situation where a Markovian analysis through state space representation is possible.

### 2.3.2 Adaptive control

In *adaptive control*, the parameters $\{\alpha_i(k), \beta_i(k)\}$ are not known *a priori*, but are partially observed through the input-output process. Typically, a parameter estimation algorithm, such as recursive least squares, is used to estimate the parameters on-line in implementations. The control law at time $k$ is computed based upon these estimates and past output measurements.

As an example, consider the system model given in equation (2.15) with all of the parameters taken to be independent of $k$, and let

$$\theta = (-\alpha_1, \cdots, -\alpha_{n_1}, \beta_1, \cdots, \beta_{n_2})$$

denote the time invariant parameter vector. Suppose for the moment that the parameter $\theta$ is known. If we set

$$\phi_{k-1}^{\top} := (Y_{k-1}, \cdots, Y_{k-n_1}, U_{k-1}, \cdots, U_{k-n_2}),$$

and if we define for each $k$ the control $U_k$ as the solution to

$$\phi_k^{\top} \theta = 0, \tag{2.17}$$

then this will result in $Y_k = W_k$ for all $k$. This control law obviously minimizes the performance criterion (2.16) and hence is a minimum variance control law if it is sample mean square stabilizing.

It is also possible to obtain a minimum variance control law, even when $\theta$ is not available directly for the computation of the control $U_k$. One such algorithm (developed in [87]) has a recursive form given by first estimating the parameters through the following stochastic gradient algorithm:

$$\begin{aligned} \hat{\theta}_k &= \hat{\theta}_{k-1} + r_{k-1}^{-1} \phi_{k-1} Y_k \\ r_k &= r_{k-1} + \|\phi_k\|^2; \end{aligned} \tag{2.18}$$

the new control $U_k$ is then defined as the solution to the equation

$$\phi_k^{\top} \hat{\theta}_k = 0.$$

With $X_k \in \mathsf{X} := \mathbb{R}_+ \times \mathbb{R}^{2(n_1 + n_2)}$ defined as

$$X_k := \begin{pmatrix} r_k^{-1} \\ \phi_k \\ \hat{\theta}_k \end{pmatrix}$$

we see that $\mathbf{X}$ is of the form $X_{k+1} = F(X_k, W_{k+1})$, where $F \colon \mathsf{X} \times \mathbb{R} \to \mathsf{X}$ is a rational function, and hence $\mathbf{X}$ is a Markov chain.

To illustrate the results in stochastic adaptive control obtainable from the theory of Markov chains, we will consider here and in subsequent chapters the following ARX(1) random parameter, or state space, model.

Simple Adaptive Control Model

The *simple adaptive control model* is a triple $\mathbf{Y}, \mathbf{U}, \boldsymbol{\theta}$ where

(SAC1)    the output sequence $\mathbf{Y}$ and parameter sequence $\boldsymbol{\theta}$ are defined inductively for any input sequence $\mathbf{U}$ by

$$
\begin{aligned}
Y_{k+1} &= \theta_k Y_k + U_k + W_{k+1} & \text{(2.19)} \\
\theta_{k+1} &= \alpha \theta_k + Z_{k+1}, \qquad k \geq 1 & \text{(2.20)}
\end{aligned}
$$

where $\alpha$ is a scalar with $|\alpha| < 1$;

(SAC2)    the bivariate disturbance process $\binom{\mathbf{Z}}{\mathbf{W}}$ is Gaussian and satisfies

$$
\begin{aligned}
\mathsf{E}[\binom{Z_n}{W_n}] &= \binom{0}{0} \\
\mathsf{E}[\binom{Z_n}{W_n}(Z_k, W_k)] &= \begin{pmatrix} \sigma_z^2 & 0 \\ 0 & \sigma_w^2 \end{pmatrix} \delta_{n-k}, \qquad n \geq 1;
\end{aligned}
$$

(SAC3)    the input process satisfies $U_k \in \mathcal{Y}_k$, $k \in \mathbb{Z}_+$, where $\mathcal{Y}_k = \sigma\{Y_0, \ldots, Y_k\}$. That is, the input $U_k$ at time $k$ is a function of past and present output values.

The time varying parameter process $\boldsymbol{\theta}$ here is not observed directly but is partially observed through the input and output processes $\mathbf{U}$ and $\mathbf{Y}$.

The ultimate goal with such a model is to find a mean square stabilizing, minimum variance control law. If the parameter sequence $\boldsymbol{\theta}$ were completely observed then this goal could be easily achieved by setting $U_k = -\theta_k Y_k$ for each $k \in \mathbb{Z}_+$, as in (2.17).

Since $\boldsymbol{\theta}$ is only partially observed, we instead obtain recursive estimates of the parameter process and choose a control law based upon these estimates. To do this we note that by viewing $\boldsymbol{\theta}$ as a state process, as defined in [39], then because of the assumptions made on $(\mathbf{W}, \mathbf{Z})$, the conditional expectation

$$
\hat{\theta}_k := \mathsf{E}[\theta_k \mid \mathcal{Y}_k]
$$

is computable using the Kalman filter (see [165, 156]) provided the initial distribution of $(U_0, Y_0, \theta_0)$ for (2.19), (2.20) is Gaussian.

In this scalar case, the Kalman filter estimates are obtained recursively by the pair of equations

$$
\hat{\theta}_{k+1} = \alpha \hat{\theta}_k + \alpha \frac{\Sigma_k (Y_{k+1} - \hat{\theta}_k Y_k - U_k) Y_k}{\Sigma_k Y_k^2 + \sigma_w^2}
$$

$$\Sigma_{k+1} \quad = \quad \sigma_z^2 + \frac{\alpha^2 \sigma_w^2 \Sigma_k}{\Sigma_k Y_k^2 + \sigma_w^2}$$

When $\alpha = 1$, $\sigma_w = 1$ and $\sigma_z = 0$, so that $\theta_k = \theta_0$ for all $k$, these equations define the recursive least squares estimates of $\theta_0$, similar to the gradient algorithm described in (2.18).

Defining the parameter estimation error at time $n$ by $\tilde{\theta}_n := \theta_n - \hat{\theta}_n$, we have that $\tilde{\theta}_k = \theta_k - \mathsf{E}[\theta_k \mid \mathcal{Y}_k]$, and $\Sigma_k = \mathsf{E}[\tilde{\theta}_k^2 \mid \mathcal{Y}_k]$ whenever $\tilde{\theta}_0$ is distributed $N(0, \Sigma_0)$ and $Y_0$ and $\Sigma_0$ are constant (see [172] for more details).

We use the resulting parameter estimates $\{\hat{\theta}_k : k \geq 0\}$ to compute the "certainty equivalence" adaptive minimum variance control $U_k = -\hat{\theta}_k Y_k$, $k \in \mathbb{Z}_+$. With this choice of control law, we can define the closed loop system equations.

---

### Closed Loop System Equations

The *closed loop system equations* are

$$
\begin{aligned}
\tilde{\theta}_{k+1} &= \alpha\tilde{\theta}_k - \alpha\Sigma_k Y_{k+1} Y_k (\Sigma_k Y_k^2 + \sigma_w^2)^{-1} + Z_{k+1} & (2.21) \\
Y_{k+1} &= \tilde{\theta}_k Y_k + W_{k+1} & (2.22) \\
\Sigma_{k+1} &= \sigma_z^2 + \alpha^2 \sigma_w^2 \Sigma_k (\Sigma_k Y_k^2 + \sigma_w^2)^{-1}, \qquad k \geq 1 & (2.23)
\end{aligned}
$$

where the triple $\Sigma_0, \tilde{\theta}_0, Y_0$ is given as an initial condition.

---

The closed loop system gives rise to a nonlinear state space model of the form (NSS1). It follows then that the triple

$$\Phi_k := (\Sigma_k, \tilde{\theta}_k, Y_k)^\top, \qquad k \in \mathbb{Z}_+, \tag{2.24}$$

is a Markov chain with state space $\mathsf{X} = [\sigma_z^2, \frac{\sigma_z^2}{1-\alpha^2}] \times \mathbb{R}^2$. Although the state space is not open, as required in (NSS1), when necessary we can restrict the chain to the interior of $\mathsf{X}$ to apply the general results which will be developed for the nonlinear state space model.

As we develop the general theory of Markov processes we will return to this example to obtain fairly detailed properties of the closed loop system described by (2.21)-(2.23).

In Chapter 16 we characterize the mean square performance (2.16): when the parameter $\sigma_z^2$ which defines the parameter variation is strictly less than unity, the limit supremum is in fact a limit in this example, and this limit is independent of the initial conditions of the system.

This limit, which is the expectation of $Y_0$ with respect to an invariant measure, cannot be calculated exactly due to the complexity of the closed loop system equations.
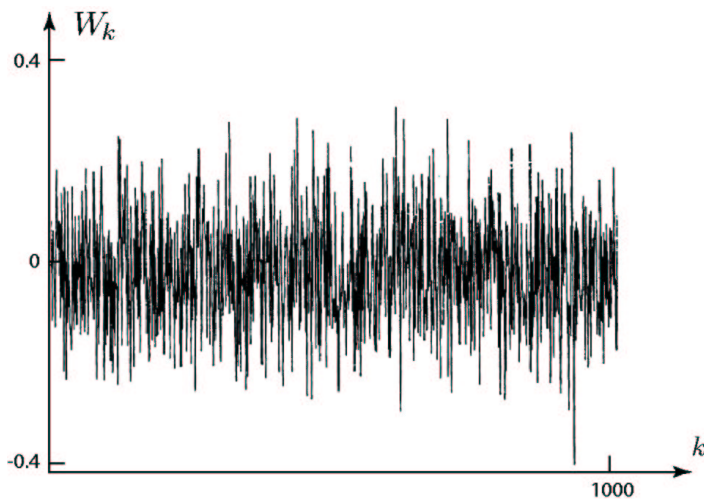
**Fig. 2.7.** Disturbance $\mathbf{W}$ for the SAC model: $N(0, 0.01)$ Gaussian white noise

Using invariance, however, we may obtain explicit bounds on the limit, and give a characterization of the performance of the closed loop system which this limit describes. Such characterizations are helpful in understanding how the performance varies as a function of the disturbance intensity $\mathbf{W}$ and the parameter estimation error $\widetilde{\boldsymbol{\theta}}$.

In Figure 2.8 and Figure 2.9 we have illustrated two typical sample paths of the output process $\mathbf{Y}$, identical but for the different values of $\sigma_z$ chosen.

The disturbance process $\mathbf{W}$ in both instances is i.i.d. $N(0, 0.01)$; that is, $\sigma_w = 0.1$. A typical sample path of $\mathbf{W}$ is given in Figure 2.7.

In both simulations we take $\alpha = 0.99$. In the "stable" case in Figure 2.8, we have $\sigma_z = 0.2$. In this case the output $\mathbf{Y}$ is barely distinguishable from the noise $\mathbf{W}$. In the second simulation, where $\sigma_z = 1.1$, we see in Figure 2.9 that the output exhibits occasional large bursts due to the more unpredictable behavior of the parameter process.

## 2.4 Markov Models With Regeneration Times

The processes in the previous section were Markovian largely through choosing a sufficiently large product space to allow augmentation by variables in the finite past.

The chains we now consider are typically Markovian using the second paradigm in Section 1.2.1, namely by choosing specific *regeneration times* at which the past is forgotten. For more details of such models see Feller [76, 77] or Asmussen [10].

### 2.4.1 The forward recurrence time chain

A chain which is a special form of the random walk chain in Section 1.2.3 is the *renewal process*. Such chains will be fundamental in our later analysis of the structure of even
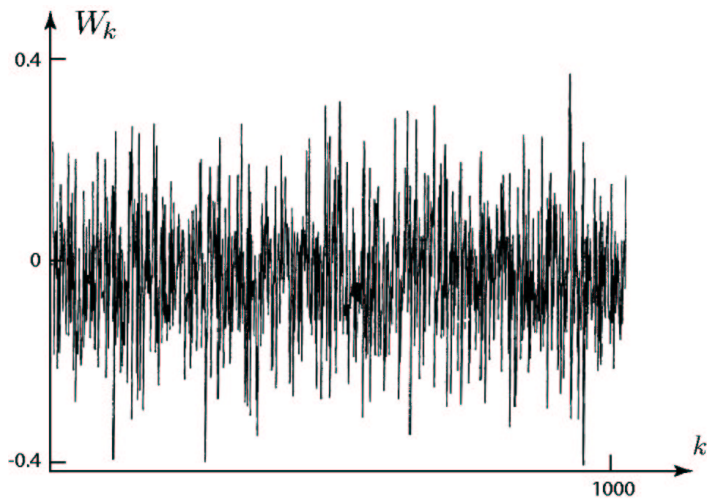
**Fig. 2.8.** Output $\mathbf{Y}$ of the SAC model with $\alpha = 0.99$, $\sigma_w = 0.1$, and $\sigma_z = 0.2$
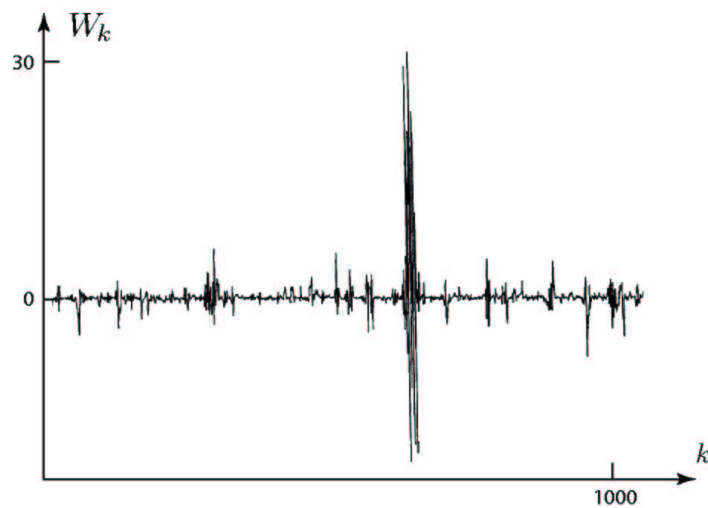


**Fig. 2.9.** Output $\mathbf{Y}$ of the SAC model with $\alpha = 0.99$, $\sigma_w = 0.1$, and $\sigma_z = 1.1$

the most general of Markov chains, and here we describe the specific case where the state space is countable.

Let $\{Y_1, Y_2, \ldots\}$ be a sequence of independent and identical random variables, with distribution function $p$ concentrated, not on the positive and negative integers, but rather on $\mathbb{Z}_+$. It is customary to assume that $p(0) = 0$. Let $Y_0$ be a further independent random variable, with the distribution of $Y_0$ being $a$, also concentrated on $\mathbb{Z}_+$. The random variables

$$Z_n := \sum_{i=0}^{n} Y_i$$

form an increasing sequence taking values in $\mathbb{Z}_+$, and are called a *delayed renewal process*, with $a$ being the delay in the first variable: if $a = p$ then the sequence $\{Z_n\}$ is merely referred to as a renewal process.

As with the two-sided random walk, $Z_n$ is a Markov chain: not a particularly interesting one in some respects, since it is evanescent in the sense of Section 1.3.1 (II), but with associated structure which we will use frequently, especially in Part III.

With this notation we have $\mathsf{P}(Z_0 = n) = a(n)$ and by considering the value of $Z_0$ and the independence of $Y_0$ and $Y_1$, we find

$$\mathsf{P}(Z_1 = n) = \sum_{j=0}^{n} a(j)p(n-j).$$

To describe the $n$-step dynamics of the process $\{Z_n\}$ we need convolution notation.

---

### Convolutions

We write $a * b$ for the *convolution* of two sequences $a$ and $b$ given by

$$a * b\,(n) := \sum_{j=0}^{n} b(j)a(n-j) = \sum_{j=0}^{n} a(j)b(n-j)$$

and $a^{k*}$ for the $k^{th}$ convolution of $a$ with itself.

---

By decomposing successively over the values of the first $n$ variables $Z_0, \ldots, Z_{n-1}$ and using the independence of the increments $Y_i$ we have that

$$\mathsf{P}(Z_k = n) = a * p^{k*}\,(n).$$

Two chains with appropriate regeneration associated with the renewal process are the *forward recurrence time chain*, sometimes called the residual lifetime process, and the *backward recurrence time chain*, sometimes called the age process.

---

**Forward and backward recurrence time chains**

If $\{Z_n\}$ is a discrete time renewal process, then the *forward recurrence time chain* $\mathbf{V}^+ = V^+(n), n \in \mathbb{Z}_+$, is given by

(RT1) $\quad V^+(n) := \inf(Z_m - n : Z_m > n), \qquad n \geq 0$

and the *backward recurrence time chain* $\mathbf{V}^- = V^-(n), n \in \mathbb{Z}_+$, is given by

(RT2) $\quad V^-(n) := \inf(n - Z_m : Z_m \leq n), \qquad n \geq 0.$

---

The dynamics of motion for $\mathbf{V}^+$ and $\mathbf{V}^-$ are particularly simple.

If $V^+(n) = k$ for $k > 1$ then, in a purely deterministic fashion, one time unit later the forward recurrence time to the next renewal has come down to $k - 1$. If $V^+(n) = 1$ then a renewal occurs at $n + 1$: therefore the time to the next renewal has the distribution $p$ of an arbitrary $Y_j$, and this is the distribution also of $V^+(n + 1)$ . For the backward chain, the motion is reversed: the chain increases by one, or ages, with the conditional probability of a renewal failing to take place, and drops to zero with the conditional probability that a renewal occurs. We define the laws of these chains formally in Section 3.3.1.

The regeneration property at each renewal epoch ensures that both $\mathbf{V}^+$ and $\mathbf{V}^-$ are Markov chains; and, unlike the renewal process itself, these chains are stable under straightforward conditions, as we shall see.

Renewal theory is traditionally of great importance in countable space Markov chain theory: the same is true in general spaces, as will become especially apparent in Part III. We only use those aspects which we require in what follows, but for a much fuller treatment of renewal and regeneration see Kingman [136] or Lindvall [155].

### 2.4.2 The GI/G/1, GI/M/1 and M/G/1 queues

The theory of queueing systems provides an explicit and widely used example of the random walk models introduced in Section 1.2.3, and we will develop the application of Markov chain and process theory to such models, and related storage and dam models, as another of the central examples of this book.

These models indicate for the first time the need, in many physical processes, to take care in choosing the timepoints at which the process is analyzed: at some "regeneration" time-points, the process may be "Markovian", whilst at others there may be a memory of the past influencing the future.

In the modeling of queues, to use a Markov chain approach we can make certain distributional assumptions (and specifically assumptions that some variables are

exponential) to generate regeneration times at which the Markovian forgetfulness property holds. We develop such models in some detail, as they are fundamental examples of the use of regeneration in utilizing the Markovian assumption.

Let us first consider a general queueing model to illustrate why such assumptions may be needed.

---

Queueing Model Assumptions

Suppose the following assumptions hold.

(Q1)   Customers arrive into a service operation at timepoints $T_0 = 0$, $T_0 + T_1$, $T_0 + T_1 + T_2, \ldots$ where the interarrival times $T_i$, $i \geq 1$, are independent and identically distributed random variables, distributed as a random variable $T$ with $G(-\infty, t] = \mathsf{P}(T \leq t)$.

(Q2)   The $n^{th}$ customer brings a job requiring service $S_n$ where the service times are independent of each other and of the interarrival times, and are distributed as a variable $S$ with distribution $H(-\infty, t] = \mathsf{P}(S \leq t)$.

(Q3)   There is one server and customers are served in order of arrival.

Then the system is called a *GI/G/1 queue*.

---

The notation and many of the techniques here were introduced by Kendall [128, 129]: GI for general independent input, G for general service time distributions, and 1 for a single server system. There are many ways of analyzing this system: see Asmussen [10] or Cohen [54] for comprehensive treatments.

Let $N(t)$ be the number of customers in the queue at time $t$, including the customers being served. This is clearly a process in continuous time. A typical sample path for $\{N(t), t \geq 0\}$, under the assumption that the first customer arrives at $t = 0$, is shown in Figure 2.10, where we denote by $T_i'$, the arrival times

$$T_i' = T_1 + \cdots + T_i, \quad i \geq 1 \tag{2.25}$$

and by $S_i'$ the sums of service times

$$S_i' = S_0 + \cdots + S_i, \quad i \geq 0. \tag{2.26}$$

Note that, in the sample path illustrated, because the queue empties at $S_2'$, due to $T_3' > S_2'$, the point $x = T_3' + S_3$ is not $S_3'$, and the point $T_4' + S_4$ is not $S_4'$, and so on.
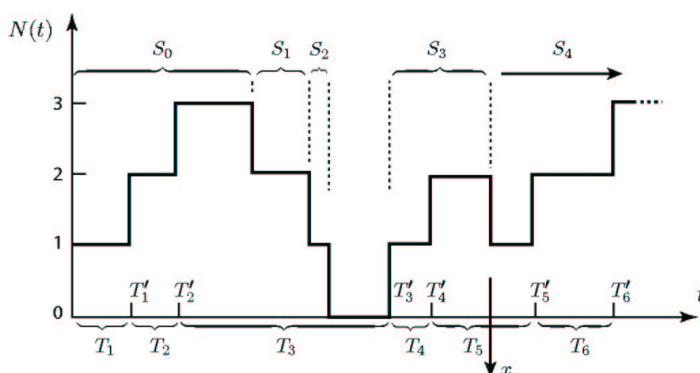
**Fig. 2.10.** A typical sample path of the single server queue

Although the process $\{N(t)\}$ occurs in continuous time, one key to its analysis through Markov chain theory is the use of *embedded Markov chains*.

Consider the random variable $N_n = N(T'_n-)$, which counts customers immediately before each arrival. By convention we will set $N_0 = 0$ unless otherwise indicated. We will show that under appropriate circumstances for $k \geq -j$

$$\mathsf{P}(N_{n+1} = j + k \mid N_n = j, N_{n-1}, N_{n-2}, \ldots, N_0) = p_k, \qquad (2.27)$$

regardless of the values of $\{N_{n-1}, \ldots, N_0\}$. This will establish the Markovian nature of the process, and indeed will indicate that it is a random walk on $\mathbb{Z}_+$.

Since we consider $N(t)$ immediately before every arrival time, $N_{n+1}$ can only increase from $N_n$ by one unit at most; hence, equation (2.27) holds trivially for $k > 1$.

For $N_{n+1}$ to increase by one unit we need there to be no departures in the time period $T'_{n+1} - T'_n$, and obviously this happens if the job in progress at $T'_n$ is still in progress at $T'_{n+1}$.

It is here that some assumption on the service times will be crucial. For it is easy to show, as we now sketch, that for a general GI/G/1 queue the probability of the remaining service of the job in progress taking any specific length of time depends, typically, on when the job began. In general, the past history $\{N_{n-1}, \ldots, N_0\}$ will provide information on when the customer began service, and this in turn provides information on how long the customer will continue to be served.

To see this, consider, for example, a trajectory such as that up to $(T'_1-)$ on Figure 2.10, where $\{N_n = 1, N_{n-1} = 0, \cdots\}$. This tells us that the current job began exactly at the arrival time $T'_{n-2}$, so that (as at $(T'_2-)$)

$$\mathsf{P}(N_{n+1} = 2 \mid N_n = 1, N_{n-1} = 0) = \mathsf{P}(S_{n-2} > T_{n+1} + T_n \mid S_{n-2} > T_n). \qquad (2.28)$$

However, a history such as $\{N_n = 1, N_{n-1} = 1, N_{n-2} = 0\}$, such as occurs up to $(T'_5-)$ on Figure 2.10, shows that the current job began within the interval $(T'_n, T'_{n-1})$, and so for some $z < T_n$ (given by $T'_5 - x$ on Figure 2.10), the behavior at $(T'_6-)$ has the probability

$$\mathsf{P}(N_{n+1} = 2 \mid N_n = 1, N_{n-1} = 1, N_{n-2} = 0) = \mathsf{P}(S_n > T_{n+1} + z \mid S_n > z). \qquad (2.29)$$

It is clear that for most distributions $H$ of the service times $S_i$, if we know $T_{n+1} = t$ and $T_n = t' > z$

$$\mathsf{P}(S_n > t + z \mid S_n > z) \neq \mathsf{P}(S_n > t + t' \mid S_n > t'); \qquad (2.30)$$

so $\mathbf{N} = \{N_n\}$ is not a Markov chain, since from equation (2.28) and equation (2.29) the different information in the events $\{N_n = 1, N_{n-1} = 0\}$ and $\{N_n = 1, N_{n-1} = 1, N_{n-2} = 0\}$ (which only differ in the past rather than the present position) leads to different probabilities of transition.

There is one case where this does not happen. If both sides of (2.30) are identical so that the time until completion of service is quite independent of the time already taken, then the extra information from the past is of no value.

This leads us to define a specific class of models for which $\mathbf{N}$ is Markovian.

---

**GI/M/1 Assumption**

**(Q4)**    If the distribution $H(-\infty, t]$ of service times is exponential with
$$H(-\infty, t] = 1 - e^{-\mu t}, \quad t \geq 0$$
then the queue is called a GI/M/1 queue.

---

Here the M stands for Markovian, as opposed to the previous "general" assumption.

If we can now make assumption (Q4) that we have a GI/M/1 queue, then the well-known "loss of memory" property of the exponential shows that, for any $t$, $z$,

$$\mathsf{P}(S_n > t + z \mid S_n > z) = e^{-\mu(t+z)}/e^{-\mu z} = e^{-\mu t}.$$

In this way, the independence and identical distribution structure of the service times show that, no matter which previous customer was being served, and when their service started, there will be some $z$ such that

$$\begin{aligned}
\mathsf{P}(N_{n+1} = j + 1 \mid N_n = j, N_{n-1}, \ldots) &= \mathsf{P}(S > T + z \mid S > z) \\
&= \int_0^\infty e^{-\mu t}\, G(dt)
\end{aligned} \qquad (2.31)$$

independent of the value of $z$ in any given realization, as claimed in equation (2.27).

This same reasoning can be used to show that, if we know $N_n = j$, then for $0 < i \leq j$, we will find $N_{n+1} = i$ provided $j - i + 1$ customers *leave* in the interarrival time $(T'_n, T'_{n+1})$. This corresponds to $(j - i + 1)$ jobs being completed in this period, and the $(j - i + 1)^{\text{th}}$ job continuing past the end of the period. The probability of this happening, using the forgetfulness of the exponential, is *independent* of the amount of time the service is in place at time $T'_n$ has already consumed, and thus $\mathbf{N}$ is Markovian.

A similar construction holds for the chain $\mathbf{N}^* = \{N_n^*\}$ defined by taking the number in the queue immediately after the $n^{th}$ service time is completed. This will be

a Markov chain provided the number of arrivals in each service time is independent of the times of the arrivals prior to the beginning of that service time. As above, we have such a property if the inter-arrival time distribution is exponential, leading us to distinguish the class of M/G/1 queues, where again the M stands for a Markovian inter-arrival assumption.

---

M/G/1 Assumption

(Q5)    If the distribution $G(-\infty, t]$ of inter-arrival times is exponential with

$$G(-\infty, t] = 1 - e^{-\lambda t}, \quad t \geq 0$$

then the queue is called an M/G/1 queue.

---

The actual probabilities governing the motion of these queueing models will be developed in Chapter 3.

### 2.4.3 The Moran dam

The theory of storage systems provides another of the central examples of this book, and is closely related to the queueing models above.

The storage process example is one where, although the time of events happening (that is, inputs occurring) is random, between those times there is a deterministic motion which leads to a Markovian representation at the input times which always form regeneration points.

A simple model for storage (the "Moran dam" [189, 10]) has the following elements. We assume there is a sequence of *input times* $T_0 = 0$, $T_0 + T_1$, $T_0 + T_1 + T_2 \ldots$, at which there is input into a storage system, and that the *inter-arrival times* $T_i$, $i \geq 1$, are independent and identically distributed random variables, distributed as a random variable $T$ with $G(-\infty, t] = \mathsf{P}(T \leq t)$.

At the $n^{th}$ input time, the amount of input $S_n$ has a distribution $H(-\infty, t] = \mathsf{P}(S_n \leq t)$; the input amounts are independent of each other and of the interarrival times. Between inputs, there is steady withdrawal from the storage system, at a rate $r$: so that in a time period $[x, x + t]$, the stored contents drop by an amount $rt$ since there is no input.

When a path of the contents process reaches zero, the process continues to take the value zero until it is replenished by a positive input.

This model is a simplified version of the way in which a dam works; it is also a model for an inventory, or for any other similar storage system.

The basic storage process operates in continuous time: to render it Markovian we analyze it at specific timepoints when it (probabilistically) regenerates, as follows.

---

**Simple Storage Models**

(SSM1)   For each $n \geq 0$ let $S_n$ and $T_n$ be independent random variables on $\mathbb{R}$ with distributions $H$ and $G$ as above.

(SSM2)   Define the random variables

$$\Phi_{n+1} = [\Phi_n + S_n - J_n]^+$$

where the variables $J_n$ are independent and identically distributed, with

$$\mathsf{P}(J_n \leq x) = G(-\infty, x/r] \tag{2.32}$$

for some $r > 0$.

Then the chain $\boldsymbol{\Phi} = \{\Phi_n\}$ represents the contents of a storage system at the times $\{T_n-\}$ immediately before each input, and is called the *simple storage model*.

---

The independence of $S_{n+1}$ from $S_{n-1}, S_{n-2}, \ldots$ and the construction rules (SSM1) and (SSM2) ensure as before that $\{\Phi_n\}$ is a Markov chain: in fact, it is a specific example of the random walk on a half line defined by (RWHL1), in the special case where

$$W_n = S_n - J_n, \quad n \in \mathbb{Z}_+.$$

It is an important observation here that, in general, the process sampled at other time points (say, at regular time points) is *not* a Markov system, since it is crucial in calculating the probabilities of the future trajectory to know how much earlier than the chosen time-point the last input point occurred: by choosing to examine the chain embedded at precisely those pre-input times, we lose the memory of the past. This was discussed in more detail in Section 2.4.2.

We define the mean input by $\alpha = \int_0^\infty x\, H(dx)$ and the mean output between inputs by $\beta = \int_0^\infty rx\, G(dx)$. In Figure 2.11 and Figure 2.12 we give two sample paths of storage models with different values of the parameter ratio $\alpha/\beta$. The behavior of the sample paths is quite different for different values of this ratio, which will turn out to be the crucial quantity in assessing the stability of these models.
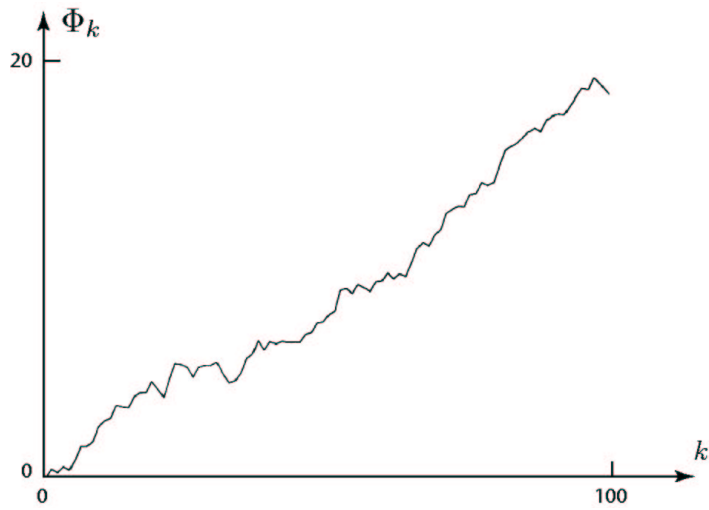
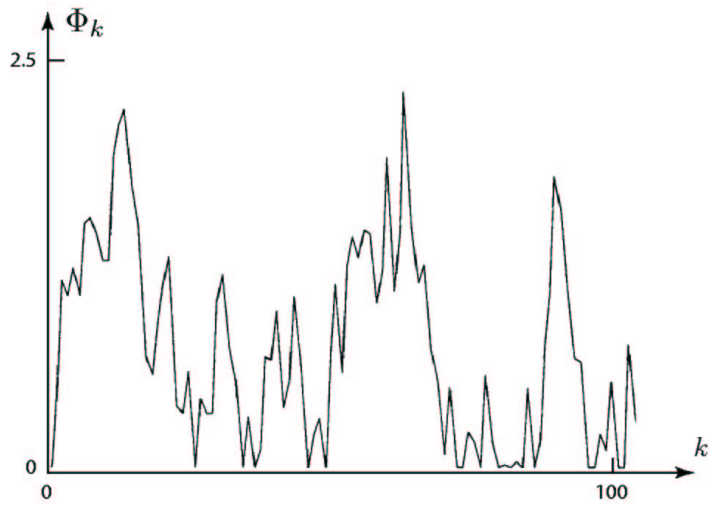**Fig. 2.11.** Storage system path with $\alpha/\beta = 2$, $r = 1$



**Fig. 2.12.** Storage system path with $\alpha/\beta = 0.5$, $r = 1$

### 2.4.4  Content-dependent release rules

As with time-series models or state space systems, the linearity in the Moran storage model is clearly a first approximation to a more sophisticated system.

There are two directions in which this can be taken without losing the Markovian nature of the model.

Again assume there is a sequence of input timepoints $T_0 = 0$, $T_0 + T_1$, $T_0 + T_1 + T_2 \ldots$, and that the interarrival times $T_i$, $i \geq 1$, are independent and identically distributed random variables, with distribution $G$.

Then one might assume that, if the contents at the $n^{th}$ input time are given by $\Phi_n = x$, the amount of input $S_n(x)$ has a distribution given by $H_x(-\infty, t] = \mathsf{P}(S_n(x) \leq t)$ dependent on $x$; the input amounts remain independent of each other and of the interarrival times.

Alternatively, one might assume that between inputs, there is withdrawal from the storage system, at a rate $r(x)$ which also depends on the level $x$ at the moment of withdrawal. This assumption leads to the conclusion that, if there are no inputs, the deterministic time to reach the empty state from a level $x$ is

$$R(x) = \int_0^x [r(y)]^{-1} dy. \tag{2.33}$$

Usually we assume $R(x)$ to be finite for all $x$. Since $R$ is strictly increasing the inverse function $R^{-1}(t)$ is well-defined for all $t$, and it follows that the drop in level in a time period $t$ with no input is given by

$$J_x(t) = x - q(x, t)$$

where

$$q(x, t) = R^{-1}(R(x) - t).$$

This enables us to use the same type of random walk calculation as for the Moran dam.

As before, when a path of this storage process reaches zero, the process continues to take the value zero until it is replenished by a positive input.

It is again necessary to analyze such a model at the times immediately before each input in order to ensure a Markovian model. The assumptions we might use for such a model are

Content-Dependent Storage Models

(CSM1)    For each $n \geq 0$ let $S_n(x)$ and $T_n$ be independent random variables on $\mathbb{R}$ with distributions $H_x$ and $G$ as above.

(CSM2)    Define the random variables

$$\Phi_{n+1} = [\Phi_n - J_n + S_n(\Phi_n - J_n)]^+$$

where the variables $J_n$ are independently distributed, with

$$\mathsf{P}(J_n \leq y \mid \Phi_n = x) = \int G(dt)\mathsf{P}(J_x(t) \leq y) \qquad (2.34)$$

Then the chain $\boldsymbol{\Phi} = \{\Phi_n\}$ represents the contents of the storage system at the times $\{T_n-\}$ immediately before each input, and is called the *content-dependent storage model*.

Such models are studied in [96, 34]. In considering the connections between queueing and storage models, it is then immediately useful to realize that this is also a model of the waiting times in a model where the service time varies with the level of demand, as studied in [38].

## 2.5 Commentary

We have skimmed the Markovian models in the areas in which we are interested, trying to tread the thin line between accessibility and triviality. The research literature abounds with variations on the models we present here, and many of them would benefit by a more thorough approach along Markovian lines.

For many more models with time series applications, the reader should see Brockwell and Davis [32], especially Chapter 12; Granger and Anderson for bilinear models [88]; and for nonlinear models see Tong [267], who considers models similar to those we have introduced from a Markovian viewpoint, and in particular discusses the bilinear and SETAR models. Linear and bilinear models are also developed by Duflo in [69], with a view towards stability similar to ours. For a development of general linear systems theory the reader is referred to Caines [39] for a controls perspective, or Aoki [6] for a view towards time series analysis.

Bilinear models have received a great deal of attention in recent years in both time series and systems theory. The dependent parameter bilinear model defined by (2.13, 2.12) is called a doubly stochastic autoregressive process of order 1, or DSAR(1), in Tjøstheim [265]. Realization theory for related models is developed in Guégan [90] and

Mittnik [186], and the papers Pourahmadi [219], Brandt [28], Meyn and Guo [177], and Karlsen [123] provide various stability conditions for bilinear models.

The idea of analyzing the nonlinear state space model by examining an associated control model goes back to Stroock and Varadhan [260] and Kunita [144, 145] in continuous time. In control and systems models, linear state space models have always played a central role, while nonlinear models have taken a much more significant role over the past decade: see Kumar and Varaiya [143], Duflo [69], and Caines [39] for a development of both linear adaptive control models, and (nonlinear) controlled Markov chains.

The embedded regeneration time approach has been enormously significant since its introduction by Kendall in [128, 129]. There are many more sophisticated variations than those we shall analyze available in the literature. A good recent reference is Asmussen [10], whilst Cohen [54] is encyclopedic.

The interested reader will find that, although we restrict ourselves to these relatively less complicated models in illustrating the value of Markov chain modeling, virtually all of our general techniques apply across more complex systems. As one example, note that the stability of models which are state-dependent, such as the content-dependent storage model of Section 2.4.4, has only recently received attention [38], but using the methods developed in later chapters it is possible to characterize it in considerable detail [178, 180, 181].

The storage models described here can also be thought of, virtually by renaming the terms, as models for state-dependent inventories, insurance models, and models of the residual service in a GI/G/1 queue. To see the last of these, consider the amount of service brought by each customer as the input to the "store" of work to be processed, and note that the server works through this store of work at a constant rate.

The residual service can be, however, a somewhat minor quantity in a queueing model, and in Section 3.5.4 below we develop a more complex model which is a better representation of the dynamics of the GI/G/1 queue.

*Added in Second Printing* In the last two years there has been a virtual explosion in the use of general state space Markov chains in simulation methods, and especially in Markov chain Monte Carlo methods which include Hastings-Metropolis and Gibbs sampling techniques, which were touched on in Chapter 1.1(f). Any future edition will need to add these to the collection of models here and examine them in more detail: the interested reader might look at the recent results [44, 191, 245, 246, 225, 166, 224], which all provide examples of the type of chains studied in this book.