TIME-SERIES ANALYSIS: A REVIEW OF TIME-DOMAIN THEORY WITH IMPLICATIONS FOR NON-LINEAR TIME-SERIES MODELLING

by

CHRISTOPHER MAH

M.A. Carleton, 1980

A THESIS SUBMITTED IN PARTIAL FULFILMENT OF

THE REQUIREMENTS FOR THE DEGREE OF

MASTER OF SCIENCE

in the Department

of

Mathematics

© Christopher Mah, 1983

SIMON FRASER UNIVERSITY

September 1983

All rights reserved. This thesis may not be reproduced in whole or in part, by photocopy or other means, without permission of the author.

APPROVAL

Name: Christopher Derrick Mah

Degree: Master of Science

Title of Thesis: Time-Series Analysis: A Review of Time-Domain Theory with Implications for Non-Linear Time-Series Modelling

Examining Committee:

Chairperson: C. Villegas

Richard Lockhart Senior Supervisor

K.L. Weldon

D.M. Eaves

T. Chang External Examiner

Date Approved:

-

September 16, 1983

PARTIAL COPYRIGHT LICENSE

I hereby grant to Simon Fraser University the right to lend my thesis or dissertation (the title of which is shown below) to users of the Simon Fraser University Library, and to make partial or single copies only for such users or in response to a request from the library of any other university, or other educational institution, on its own behalf or for one of its users. I further agree that permission for multiple copying of this thesis for scholarly purposes may be granted by me or the Dean of Graduate Studies. It is understood that copying or publication of this thesis for financial gain shall not be allowed without my written permission.

Title of Thesis/Dissertation:

Andy Cis: review o with Implicate ream Non-Line

<u>Author:</u>

(signature)

CHRISTOPHER MAH (name)

(date)

Abstract

Many phenomena of interest in the sciences can be modelled by a deterministic relation perturbed by a sequence of errors in discrete time. Such a model is called a time-series. The linear time-series is reviewed in detail and theory of imperfections of well established and recent estimation methods is suggested that a priori structural are pointed out. Ιt information is a neglected resource in this regard. To clarify meaning of structural information, representation and the prediction theory for weakly stationary processes is reviewed. seen to consist of a Weakly stationary processes are deterministic and a moving average part. Linear models are most suited to simulation of the moving average component. Recent proposals for non-linear time-series models are assessed in light of this discussion. It is argued that one recent proposal (threshold autoregression) is more suited to modelling of timeseries with complicated deterministic parts, and requires a priori structural knowledge to be fully effective. In contrast another proposal (bilinear) seems more suited to time-series which are non-gaussian moving averages, and are less dependent on a priori structure for their usefulness.

iii

Table of Contents

Approval
I. INTRODUCTION1
I. LINEAR MODELS: THEORY AND PRACTICE
6 1.2.1 STATIONARITY 1.2.2 INVERTIBILITY 1.3 BOX-JENKINS MODELS:COVARIANCE STRUCTURE AND UNIQUENESS OF THE MODEL
UNIQUENESS OF THE MODEL
II. HOW WELL DO THESE METHODS WORK?
III. METHODS OF SELECTING AR ORDER
IV. STRUCTURAL PROPERTIES: LINEAR EXAMPLES
 V. REPRESENTATION THEORY FOR WEAKLY STATIONARY PROCESSES: IMPLICATIONS FOR NON-LINEAR TIME-SERIES MODELLING 5.1 WOLD'S THEOREM AND SINGULAR PARTS 5.2 A NONLINEAR AUTOREGRESSIVE/MOVING AVERAGE DUALITY 90 5.3 PREDICTION THEORY FOR MOVING AVERAGES 92 5.4 BILINEAR MODELS
VI. APPROPRIATENESS OF FULLY NON-LINEAR MODELS: PRACTICE 6.1 VIOLATIONS OF LINEARITY
APPENDIX A - SOLVING A HOMOGENEOUS DIFFERENCE EQUATION121 BIBLIOGRAPHY

List of Tables

I.	AR order Selection
II.	AR order selection
III.	AR Order Selection
IV.	Models for Series A60
v.	Values of AIC for Various Models61

v

I. INTRODUCTION

Many phenomena of interest in the social, life and physical sciences can be usefully modelled at time-series. A time-series is discrete-time а stochastic process generated by а deterministic relation such as a difference equation, which is perturbed by errors. Theoretically, the errors are viewed either as exogenous disturbances of the deterministic dynamics of the system or alternatively, as input. Errors are modelled which are typically assumed as random variables, to be independently and identically distributed (iid).

Time-series have been of interest to statisticians because time-series data often occurs for which causal mechanisms are relatively obscure, but for which it is desirable to

- i) facilitate prediction
- ii) develop an insight into the causal structure of the underlying natural phenomenon

The distinction between of (i) (predictive the pursuit time-series analysis) and of (ii) (structural time-series analysis is due to Parzen(1974).) This thesis will focus on the close interrelation between these goals and we will maintain that a systematic recognition of this relationship is valuable for both purposes. This point will be illustrated for a well known special class of time-series, and later applied to the evaluation of some recent proposals for time-series models.

The time-series problem is to obtain, from a data set, a parametric form for the underlying stochastic process and estimate its coefficients. The solution of this problem demands more than the classical estimation problems of mathematical statistics because the appropriate parametric form is not specified in advance.

A time-series data set is a finite sequence of real numbers

(1.1)
$$\{ X(t) \}_{i=1}^{N}$$

obtained by observing a phenomenon of the type described above at the times t . When the meaning is clear from the context, a i time-series data set may also be referred to as a time-series. Note that the errors occuring in the time-series model are not mentionned in this formulation of the problem, and are assumed not to be observed. In practice therefore, the errors can only be estimated from residuals, that is, from the differences between the predictions of the model and the observed values. While it can be useful to view the errors as input this is not the perspective emphasized in this thesis.

The situation in which input is observed leads to a different statistical problem from the one we consider. This alternate problem has a large literature and is appropriately treated by different methods (ie.Spectral theory) than those to be developed here. Instead, we will discuss chiefly the timedomain approach to time-series modelling. This approach deals

with a time-series directly rather than via its spectral representation.

deals with current methods of solving the time-Chapter I series problem when only linear models admitted, are and develops the theory of linear time-series models. Chapter II assesses these methods, and compares the performance of skilled automatic procedures on the problem of versus choice of autoregressive order. Because automatic methods are of special interest in the context of non-linear time-series modelling, Chapter III develops the theory of a popular automatic method, Akaike's information criterion. We conclude that automatic methods are not a panacea, and in Chapter IV demonstrate the use of structural assumptions to facilitate linear model fittina. This is, in a sense, an alternative approach to the time-series problem, and we focus on the application of this approach to non-linear modelling. Chapter V reviews representation theory for weakly stationary processes and its implication for structural modelling. Finally, Chapter VI deals with the practical application of these conclusions.

I. LINEAR MODELS: THEORY AND PRACTICE

1.1 TIME-SERIES MODELS

In attempting the time-series problem it is necessary to consider only narrow special classes of candidate models, for otherwise the problem would be intractible in practical applications. Specific hypotheses are choosen for computational convenience and so that they mirror a priori knowledge about the real processes being studied. We will enumerate the general types of times series model in common use, and discuss the solution of the time-series problem in a special case: that of the well known Box-Jenkins (linear) models. This treatment will raise general issues relevant to the structural/predictive dichotomy in non-linear time-series analysis.

Two basic types of hypothesis may be distinguished, which we will refer to loosely as models. First, models of moving average type. The most elementary form of such a model is, called simply a moving average and has the form

(1.1) $X = q \epsilon + q \epsilon + \dots + q \epsilon$ t 0 t 1 t-1 n t-n where t is integer, q are real constants, and ϵ is an i t uncorrelated mean-zero sequence of random variables.

In the notation of Box & Jenkins such a model is an MA(n)

(moving average process of order n) model and is written as

(1.2)
$$\begin{array}{c} m & i \\ X = \Sigma & q & B & \epsilon \\ t & i & i & t \end{array}$$

where the q are real constants, B is the backward shift i operator defined by BX = X and Q(B) is a polynomial defined t t-1 by the second equality. If Q(B) instead of being a linear operator is some more general operator defined on the space of sequences , then the time-series is of moving average type. In considering this model it is important to recall again that the errors are not observable.

The second class of models are those of autoregressive type. In the most elementary form of such a model X satisfies t

(1.3)
$$X = pX + pX + \dots + pX + \epsilon$$

t 1t-1 2t-2 nt-n t
where t is integer and the p are real constants.

In the notation of Box& Jenkins such a model is an AR(n) (autoregressive model of order n) process and it is often written

(1.4)
$$\sum_{i=1}^{n} \sum_{j=1}^{i} P(B)X = \epsilon$$

where P(B) is a polynomial defined by the first equality. We will refer to the term

$$\begin{array}{c} \text{(1.5)} \\ \text{B} \\ \text{X} \\ \text{t} \\ \text{t} \\ \text{t} \end{array}$$

as the term corresponding to the lag of order m.

If P(B) instead of being a linear operator is some more general operator defined on the space of sequences, then the time-series is of autoregressive type. Time-series which satisfy relations of the form

(1.6) P(B)X = Q(B)et t

where the symbols are as above, are said to be of mixed type.

The elementary linear versions of these models were introduced by Yule(1927) and are supported by a large body of theory (Box&Jenkins,1970:Koopmans,1974; Deutsh,1965), as well as a set of relatively refined estimation techniques. For example, many commonly available time-series analysis programmes (such as the IMSL routine FTCMP) are devoted to the estimation of parameters in Box-Jenkins models.

1.2 BOX-JENKINS MODELS: INVERTIBILITY AND STATIONARITY

We now begin a discussion of the solution of the timeseries problem in the linear case. To do this it is necessary to describe the theory of Box-Jenkins autoregressive/moving average (ARMA) models. We begin with the definition of the natural sample statistics to compute if one is interested in the linear predictibility of the observations. Consideration of these definitions leads to assumptions on the stochastic such procedures make processes so that sense; namely weak

stationarity and invertibility. The meaning and mathematical interpretation of these conditions is discussed. Finally, it will be shown that under these assumptions the covariance structure of a process leads to a nearly unique determination of the coefficients of an ARMA model.

In examining an empirical data series and assessing its predictibility it is natural to compute the autocorrelation function of the observations. This is defined by

(2.1)
$$r = \frac{N-j}{N-j} \sum_{i=1}^{N-j} (x - \bar{x}) (x - \bar{x}) / \sum_{i=1}^{N} (x - \bar{x}) (x - \bar{x})$$

where N is the length of the series and \bar{x} is the series mean. This statistic summarizes the extent to which we can (within the sample) linearly predict X from X. For this t+j t procedure to give a result generalizable to another sample from the same time-series, it must be assumed that

(2.2) $E[X] \equiv \mu = \int X(\omega)dP = \lim_{T \to \infty} \frac{1}{2T} \sum_{T \to T}^{T} X < \infty$

independent of t and

3)
$$\rho = \int (X (\omega) - \mu) (X (\omega) - \mu) dP$$
$$j \qquad t+j \qquad t$$
$$= \lim_{T \to \infty} \frac{1}{2T} \sum_{T} (X - \mu) (X - \mu) < \infty$$

(2.

The first equalities for all t comprise the assumption of weak stationarity while the second set of equalities (limits being taken in the L sense) characterize ergodicity. The assumption of weak stationarity may be contrasted with (strong) stationarity, for which the joint probability distribution of any finite collection rather than only the first two moments, must be independent of Thus the process X is strongly t. stationary if and only if the joint distribution function of the collection

does not depend on t.

Much of the treatment of Box-Jenkins models assumes only weak stationarity of the error series. In particular, it is assumed that E[e] = 0 and E[e e] is constant in t for each t t-k zero for k=0. k, and In general this implies only weak stationarity of the resulting series X However, should be noted that when the error series is normally distributed, a Box-Jenkins series Х has multivariate normal joint distributions t because of linearity. Thus, in this case weak stationarity of the errors implies strong stationarity of the model.

moving average is always stationary. Α If the errors become independent sufficiently rapidly, it is also ergodic. Suppose for example that ϵ is independent of ϵ for k>K , all moments of the sequence exist, and the variance is joint constant. Now if K=1 (the sequence is independent) then the ε t sequence is mean-ergodic because, the variance of the sum 2 Thus the variance of the average is order Σε is N σ . 1 t t N and so, tends to zero as $N^{-} > \infty$. Suppose now K > 1. Computing the variance of the sum using the assumptions we find that the number of non-zero covariances in the computation grows as NK. 2 Thus, since each covariance is bounded by σ the variance of the average is again order 1 and 50 tends to zero. Covariance ergodicity is obtained similarly on computing the variance of This shows that the error sequence the sum Σεε is t t t-l ergodic in the sense of (2.2-3), and from this it follows that finite moving averages are ergodic. A further easy argument extends this result to infinite moving averages which converge geometrically, as the variance of the tail of the sum tends to zero.

It is a comparatively simple matter to compute the the theoretical autocovariance function γ for an MA(n) process.

It is computed as

$$\gamma = \sigma \sum_{i=1}^{2} q_{i}$$

 $\gamma = \sigma \sum_{m \in i} 2 q q$

for m≤n

 $\gamma = 0$ for m>n m

where σ is the variance of the ϵ .

It is a little more difficult to compute the autocovariance function of a (stationary) autogressive process and this will be postponed until later. For now we wish to deduce its behavior for large lags in order to contrast it with an moving average process. To do this notice that because X satisfies the t difference equation (1.4), the autocovariances $Cov(X, X) = \gamma$ satisfy the difference equation t t-n t, n

(2.6)
$$\gamma - q \gamma + \dots + q \gamma = 0$$

t, m 1 t, m-1 m t, m-n

To see this multiply the original equation by X and take t-nexpectations. Because of stationarity, the t subscript may be omitted and it is seen that the autocovariance sequence $\{\gamma\}$ m satisfies the homogeneous difference equation. The difference equation (2.6) can be solved to yield a sequence which is a linear combination of sequences of the form

(2.7)
$$t kt$$

 $\gamma \text{ or } t \gamma$

For notice that the operator P(B) may be resolved into-

(2.8)
$$(I - r B) (I - r B) \dots (I - r B)$$

1 2 n

the identity operator and B is the backward shift where Ι is operator, on using the fundamental theorem of algebra, and the resulting factors commute. Then substitute terms of the form above to see that they are annihilated by the resulting The conclusion is operator. that autocorrelations of an AR process must grow or decay exponentially with time. This is explained more fully in appendix A.

1.2.1 STATIONARITY

We have seen that a moving average process is automatically stationary. In contrast there are non-trivial conditions for an autoregressive process to be stationary. For consider the timeseries

$$(2.9) X = pX + \epsilon$$

where $p \neq 1$, the ϵ are a stationary zero mean process and ϵ is t t independent of X . The variance of X can be found by t-1 t squaring both sides of this equation and taking expectations. Doing this, bearing in mind the independence of X and ϵ we t-1 t

find

(2.10)
$$\begin{array}{c} 2 & 2 & 2 \\ E[X] = p & E[X] + \sigma \\ t & t^{-1} & \epsilon \end{array}$$

where σ is the variance of the ϵ . If {X } were stationary ϵ t t

 $E[X] = \gamma$ a constant, and the equation becomes t o

(2.11)
$$\gamma = p \gamma + \sigma$$

Inspecting equation (2.11) it is seen that there is no constant positive solution for γ when |p|>1. Thus no process satisfying the difference equation (2.9) with |p| > 1 and ϵ independent of t

t-1

If |p| < 1, then the sum $\sum p \epsilon$ defines a stationary k=0 t-k ergodic random process with a finite variance. It is easy to satisfies equation (2.9).Is this solution see that this To answer this question note $E[X] = \gamma$ that is unique? t 0 uniquely determined by (2.11). Moreover, because the covariances satisfy the homogeneous difference equation (2.12) $\gamma = p\gamma$ t+1 t

whenever Var[X] = E[X] is known, Cov[X,X] for k>0 are
 t t t
uniquely defined. Thus, there is a weakly stationary process

with unique covariance structure satisfying the difference equation with |p|<1.

When the errors are iid this is the stochastic process which results from starting the Markov chain X in its t stationary distribution. Henceforth, this situation will be summarized as; an AR(1) process is stationary iff |p|<1. Note that we now say stationary, and understand, weakly stationary.

Now this is not the whole story on stationarity of linear processes, for if the above procedure is repeated under the assumption that ϵ is independent of X in the defining t difference equation, it is then found that the resulting process is stationary iff |p|>1. This point will be clarified in section III but first we will extend the above reasoning to more complicated AR processes.

Observe that if X is a stationary process so is P(B)X for t any P. To see this just write out the operator, multiply by any time translate and take expectations. This fact, together, with the discussion of the above process implies that

(2.13)
$$P(B)X = e t t$$

can define a weakly stationary process only when the roots of P(x)=0 lie outside the unit circle (note that the roots of P(x)=0 are the reciprocals of the numbers r in the resolution i of the operator P(B)). For resolve P(B) as (2.14) $P(B) = (I-r B) (I - r B) \dots (I - r B)$

2

We procede by contradiction. Suppose that some $|r_i| \ge 1$, say for i=1 and that X is stationary. Define $T(x) = (I - r x)^{-1} P(x)$ (2.15)and set (2.14) $Z_{t} = T(B)X_{t}$ Then Z satisfies t (2.16)= (I - r B)Zby the discussion Z is non-stationary, so but being a 'polynomial' in lags of X it is stationary by the hypothesis. gives a contradiction, and the assertion is proved. The This implications of roots on the unit circle will be discussed in а later section dealing with non-stationary behavior.

Finally note that if |r| < 1 for i=1,n then the operator i P(B) may be inverted, factor by factor using the formula

(2.17) $(I - r B)^{-1} = I + rB + r B + \cdots$

to obtain an infinite moving average representation for the series which converges in probability and solves the difference equation (2.13). In section 1.3 it will be seen that the covariance structure of any solution must be the same as the covariance structure of this sum.

1.2.2 INVERTIBILITY

We are now prepared to consider the notion of invertibility of a process. The invertibility of a process is an algebraically dual condition to stationarity which also has a practical meaning. Consider the mixed process

(2.18) P(B)X = Q(B)e

As shown above, the condition for stationarity is that P(x)=0 have roots lying outside the unit circle. The algebraic condition for invertibility of a process is that Q(x)=0 have roots lying outside the unit circle. One consequence of this will clearly be that Q(B) may be inverted to yield a purely autoregressive representation

(2.19)
$$e = Q$$
 (B) P(B)X
t

for the mixed process.

The invertibility condition may also be interpreted as follows: Suppose that X(t) is a series of interest and that at time t=n a forecast f(n,h) is required of the future value X(n+h). If a least square criterion is adopted then the optimal forecast is

(2.20)
$$f(n,h) = E\{X(n+h) | past\}$$

or in other words

(2.21)
$$f(n,h) = \sum_{i=1}^{x} p_i f(n,h-i) + \sum_{i=1}^{e} q_i f(n,h-j)$$

where

$$f_{n,h-j}^{x}(n,h-j) = \begin{bmatrix} X(n+h-j) & \text{for } j \ge h \\ f(n,h-j) & \text{for } j < h \\ f(n,h-j) &= \begin{bmatrix} 0 & \text{for } j < h \\ e(n+h-j) & \text{for } j \ge h \end{bmatrix}$$
Thus, consider the MA process
(2.22) $X = e + qe = (I + qB) e \\ t = t + t = t + t = 1 & t \end{bmatrix}$
where $e = is a$ stationary zero mean process. On applying the above principle conclude
(2.23) $f(n,1) = qe = n$
however, the $e = were$ not observed and we know only $X = t = t$
must be estimated from previous values of $X = t$
We now investigate the feasibility of this as the data series becomes longer. First note that if there was a good guess (or any guess)
 $\hat{e} = for e = 0 = 0$

available then the recursion could be used to get a sequence of estimates

(2.24) $\hat{e} = X - q \hat{e}$ t t t-1

since X(t), t=1,N were observed. It turns out that this is

only sensible when |q| < 1. To see this define the inaccuracy h

as

(2.25)
$$h = \hat{e} - e_{t}$$

then

$$h = \{ [X] - q\hat{e} \} - e \\
 t t t^{-1} t$$

 $= [e + qe] - q\hat{e} - e = - qh$ t t - 1 t t - 1 t t - 1

Thus, using this method, on a long data series, the inaccuracy in an initial guess dies out when |q|<1 and fails to die out when $|q|\geq 1$. Because the errors are not observed, this is the only possible method of forecasting. Therefore the practical interpretation of invertibility is that the optimal forecast may be approached arbitrarily closely with sufficiently long data series. The above example (2.22) is invertible iff |q|<1.

For a more complicated MA process

(2.26)
$$X = Q(B) e_{t}$$

The optimal forecast (as noted) still involves e . But for any t

initial guess ê we have

(2.27)
$$Q(B) h = 0$$

and the solutions of this difference equation are stable only if the roots of Q(x)=0 lie outside the unit circle. When this condition is satisfied it is again true that the inaccuracies h t in the choice of ê die out, and the optimal forecast may be t approached arbitrarily closely.

1.3 BOX-JENKINS MODELS:COVARIANCE STRUCTURE AND UNIQUENESS OF THE MODEL

To discuss covariance structure it is most convenient to have the model in moving average form. We claim that any mixed ARMA model P(B)X = Q(B)e for which $|r| \neq 1$ in the resolution of t t P, can be written

(3.1)
$$X = \Sigma \quad v \quad e \quad \equiv \quad V(B) \quad e \\ t \quad -\infty \quad j \quad t-j \qquad t$$

where V(x) is a formal power series. For certainly the factors I - rB in the resolution of P with |r| < 1 may be inverted as noted above, while for each factor with |r| > 1 rewrite as

(3.2)
$$(I - rB) = BFr (I r - B)$$

$$=-B r (I - r F)$$

where B is the backward shift operator and F is the forward shift operator, so that

(3.3)
$$(I - rB)^{-1} = -(I - rF)^{-1} rF$$

In this manner it becomes meaningful to write

(3.4)
$$X = P(B)Q(B)e = V(B)e$$

t t t

The autocovariances in this model can be computed as

(3.5)
$$E[X X] = \Sigma \Sigma V V E[e e] t+jt k l k l t+j-k t-1$$

$$= \sigma \Sigma \nabla \nabla \nabla = \gamma$$

$$\epsilon 1 1 + j 1 j$$

Next, define the autocovariance generating function

(3.6)
$$\Lambda(B) = \sum \gamma B$$

Substituting the previous expression for the autocovariances obtain

(3.7)
$$\Lambda(B) = \sigma \sum_{\epsilon} \sum_{j=1}^{2} v v B = \sigma \sum_{i=1}^{2} \sum_{j=1}^{j-k} v B = \sigma \sum_{i=1}^{j-k} v B = \sigma \sum_{i=1}^{$$

We can now show that for a given autocovariance generating function, the resolutions of the operators P(B) and Q(B) are defined up to reciprocals of the roots. For the autocovariance generating function corresponding to (3.4) is

(3.8)
$$2 - 1 - 1$$

 $\sigma \{ P(B) Q(B) \} \{ P(F) Q(F) \}$

using (3.7).

Thus, from the two operator identities(3.9)

(i)
$$(1 - \phi_{B})(1 - \phi_{F}) = \phi_{i}^{2} (1 - \phi_{B})(1 - \phi_{F})$$

i i i i i

(ii)
$$(1 - \psi B) \begin{pmatrix} -1 & -1 & -2 & -1 & -1 & -1 & -1 \\ (1 - \psi B) & (1 - \psi F) &= & \psi & (1 - \psi B) & (1 - \psi F) \\ i & i & i & i & i \\ & & & & i & & i \end{pmatrix}$$

here ψ and ϕ are real constants, it is seen that all the i

processes

i

W

where δ is chosen to scale the variance appropriately have the same autocovariance generating function and hence the same autocovariances.

Thus, the restrictions of invertibility and stationarity allow the choice of exactly one of the reciprocal pairs of each factor (provided the roots are off the unit circle) and ensure that the model corresponding to a given covariance structure is unique up to cancellation of factors in the equation.

It is important to note that the non-uniqueness discussed above is not a non-uniqueness in distribution. In fact when the errors are normally distributed, knowledge of the covariance function is equivalent to knowledge of the joint distributions of

(x, x, ..., x) t t t (3.11)

where is any finite collection of times. Rather, the nont i uniqueness arises from the fact that the Box-Jenkins representations contain information about the relation of past,

present and future. No restrictions follow from this when the defining relationship involves a finite number of random variables. For, if the present is time t=n then a stationary process which satisfies

(3.12)
$$X = p X + p X + \dots + p X$$
$$n+m \quad o \quad n+m-1 \quad 1 \quad n+m-2 \qquad n \quad n$$
for each n must also satisfy

(3.13) X = p X + p X + ... + p X $n \quad o \quad n-1 \quad 1 \quad n-2 \quad m \quad n-m$

But the situation is different for one-sidedly infinite relations. Consider the stationary AR(1) process

(3.14) X = qX + e t t - 1 t

Recall that this can be written as

(3.15) (I - qB) X = et t

and this implies (for $|q| \le 1$)

(3.16) X = e + qe + qe + ...t t t t-1 t-2

Hence, by a straightforward calculation we see that when the autocovariance generating function makes sense (ie. when X t is stationary) it depends only on q. There is no non-uniqueness here. However it was seen above that for each value of q there is a Box-Jenkins time-series equation giving the same autocovariance generating function which depends on 1/q. Does this equation generate a distinct stationary process? Recall in the discussion of stationarity conditions it was stated that if in an AR(1) difference equation е is that assumed t independent of X, the resulting process is stationary iff |q|>1. Suppose then that we assume |q|>1 and e independent of Rewrite the difference equation for this stationary process х.

(3.17) (I - F/q) X = w t t

where F is the forward shift operator and

(3.18)
$$w = -e / q$$

t t

then since |q| > 1 this process may be inverted as a function of future errors:

(3.19)
$$\begin{array}{cccc} & & -1 & -2 \\ X &= w &+ q & w &+ q & w &+ \cdots \\ t & t & t + 1 & t + 2 \end{array}$$

Thus there are two distinct processes with the same autocorrelation function . While the two processes above will have the same autocovariances and possibly the same joint distributions, one describes a process independent from past e 's and the other a process independent of future e 's. t

Fortunately the distinction is merely philosophical since the errors are not observed, so that there is no way to distinguish these models from data. It makes good sense however, to choose the representation in which a convergent expansion in terms of the past is stationary.

In a parallel fashion consider the MA(1) process

$$(3.20) X = (I - pB) e t t t$$

this may be inverted as (for |p|<1)

(3.21)
$$e = X + pX + pX + ... t t t t-1 t-2$$

There is another Box-Jenkins model which gives the same autocovariance generating function, which is constructed by replacing B with F and p by 1/p and multiplying by an appropriate scale factor. To verify that a distinct stationary process is generated note that if |p|>1, the original equation may be rewritten as

(3.22)
$$Z = (I - p F) e$$

t t-1

where F is the forward shift operator and Z = -X / pt t

and may be 'inverted' as

(3.23)
$$e = Z + p Z + p Z + ...$$

t-1 t t+1 t+2

Thus in the first case (AR(1)) the resulting process X t depends on its past and a current error e while in the second t (MA(1)) X depends on its future and a current error e. This ambiguity is again resolved by appealing to physical intuition. 1.4 NON-STATIONARY MODELS

Box-Jenkins models provide naturally for non-stationary behavior. The most useful type of non-stationarity occurs when |r| = 1 in one or more of the factors I-rB in the polynomials P,Q. First note that if r = 1 that the factor I-rB becomes the difference operator D defined by

(4.1)
$$DX = X - X$$

t t t-1

This suggests a natural model for data which locally exhibits a stochastic polynomial trend over time. Differencing once transforms a sequence which grows linearly to a sequence of constants, and differencing twice takes a sequence growing quadratically to a sequence of constants. Therefore sequences which contain polynomial trends of degree d but which become stationary when the trends are removed, may be modelled as the solutions of

$$\begin{array}{c} d \\ P(B) D X = Q(B) e \\ t \end{array}$$

If d=1 this model is a generalization of a random walk model in which the increments, instead of being independent, are a Box-Jenkins process. These models $d\geq1$ are known as integrated autoregressive moving average models and are identified by the notational convention ARIMA(p,d,q) where p and q are respectively the autoregressive and moving average orders and d

is the order of differencing. On writing down the associated difference equations for the autocovariances as we did in section 1.2, it is seen that the autocorrelation function for an ARIMA model with $d \ge 1$ will tend to grow rather than die out, and is the failure of the sample autocorrelation function to die it out that is taken as indication of non-stationarity in an practice. On the other hand if $r = \cos[\theta] \pm i \sin[\theta]$, $r \neq 1$ then this implies periodic behavior of the series, for example if |r| = 1, $r \neq 1$, i=1,2 and i

(4.3) $P(B)X = (I - r B) (I - r B) X = e_{t}$

then the null space of the difference operator includes functions such as $\cos[\theta t]$ and $\sin[\theta t]$, as the real and imaginary

parts of r must separately solve the homogeneous equation. This means the difference equation admits non-stationary solutions such as

$$(4.4) \qquad \cos[\theta t] + w \\ t$$

where w is any solution of the defining relation (4.3).

In this manner it is possible to model data which contains periodic components with random changes in amplitude and phase. Because the autocovariances must satisfy the homogeneous difference equation, it may be seen as before that the autocorrelation function of such a process will fail to die out.

1.5 SOLVING THE TIME SERIES PROBLEM

Box&Jenkins (1970) distinguish the following stages of model building:

- (1) Identification, in which we attempt to determine a class of models which may be sensibly entertained in the light of a priori considerations and a rough data analysis, but is small enough that reasonably efficient parameter estimates are available
- (2) Estimation, in which, having decided upon a class of models we form estimates of the parameters of the models which is supposedly being observed
- (3) Diagnostic checking, in which by residual analysis it is tested whether the model is consistent with the observations post hoc.

We now discuss each of these stages in turn, sketching the general procedures involved. The aim of this section is to demonstrate that the linear time-series problem is a difficult statistical problem which leads to many unresolved issues. Ιt will become clear that the fitting of a Box-Jenkins model, even basic assumptions are correct may leave much room, when the (relative to other statistical procedures) for judgement. It will be suggested that a somewhat neglected method of coping with such difficulties is to make more deliberate use of information. structural A discussion of some structural properties of Box-Jenkins models follows which will later be contrasted with the properties of non-linear time-series models.

1.5.1 IDENTIFICATION

When it is known in advance what the correct form of the linear model is, estimating the coefficients of a time-series is relatively straightforward. Unfortunately such knowledge is rarely available in convenient form. Thus the goal in the early stages of data analysis is to decide whether a time-series is purely a moving average, purely autoregressive or mixed and to determine the orders of the components. To decide initially on the correct class of linear models it is usual to begin by computing the autocorrelation function and the partial autocorrelation function of the observations.

The partial autocorrelation function is defined as follows. Let H for t,s t > s+1 and let \hat{x} and \hat{x} denote respectively the predictions of S and X from the true regressions of X and X on H $\,$. Х Then t t,s th the t-s partial autocorrelation of a stationary process is, defined as

(5.1)
$$\phi = \text{Correlation}[X - \hat{X}, X - \hat{X}]$$
$$t-s \qquad t \quad t \quad s \quad s$$

where we agree that when both the arguments of the correlation are constant we will set ϕ = 1. When both the arguments of t-s

the latter correlation are non-constant random variables the n

th

partial autocorrelation is

(5.2)
$$\hat{p}^{''}$$
 = the last coefficient in the true multiple

We have seen that the theoretical autocorrelations for an MA process have a cutoff, while the autocorrelations of an AR process are non-zero for large lags. Now the latter statement holds with autocorrelations replaced by partial autocorrelations and AR and MA interchanged. Thus, the computation of the autocorrelation function and partial autocorrelation function is a preliminary procedure for distinguishing MA and AR processes. The claim is proved as follows. Consider the autoregression

(5.3)
$$e = \sum_{i=0}^{n} p X_{i-i}$$

where p = 1, the series converges in probability and possibly
o
n=∞. If the object is to minimize

(5.4)
$$E[(X - \hat{X})^2]$$

where

 $\hat{p} X + \hat{p} X + \dots + \hat{p} X$ for $m \ge n$ then we $1 t-1 \qquad 2 t-2 \qquad m t-m$

must minimize

Ŷ

(5.5)
$$Var\{(\hat{p} - p) X + (\hat{p} - p) X + \dots \\ 1 \quad 1 \quad t-1 \quad 2 \quad 2 \quad t-2 \quad \dots$$

+
$$(\hat{p} - p) X + \hat{p} X + \dots + \hat{p} X + e$$

n n t-n n+1 t-n-1 m t-m t

Since the X are jointly non-degenerate random variables, and e is uncorrelated with X for k>0 conclude $\hat{p}_{i} = p_{i}$ for i=1, m(5.6)and in particular $\hat{p} = 0$ for i > n. Thus for i > n, where n is the order of the autoregressive process ϕ = 0. When m < n the latter variance (5.5) becomes Var{ $(\hat{p} - p) X + (\hat{p} - p) X + \dots$ 1 1 t-1 2 2 t-2 (5.7) + $(\hat{p} - p) X + p X + ... + p X + e$ m m t-n m+1 t-m-1 n t-n t and the first argument fails. In fact, some $\hat{p} \neq p$ unless x + ... + p X + e is uncorrelated with X ,...,X m+1 t-m-1 n n t t-1 t-m but these quantities are known to be correlated. Therefore we note that the 'true' coefficients of the multiple regression of X , ... , X for m < n, are not the coefficients of, t-1 t-m on X the underlying Box-Jenkins process.

The preceding argument shows that when a process is autoregressive of finite order, the partial autocorrelation function cuts off, and when the process is a moving average, n is infinite and so the partial autocorrelation function is generally non-zero for arbitrarily large lags as claimed. Ramsey(1974) has characterized the partial autocorrelation

function of a stationary process by the following striking result. Define the set of sequences S as

(5.8) $S = \{ [\phi_{i}] \in \mathbb{R}^{\infty} : |\phi_{i}| \leq 1 \text{ and } |\phi_{i}| = 1 \Rightarrow \phi_{i+1} = 1 \}$ then a function ϕ_{i} belongs to S iff it is the partial autocorrelation function of a stationary process. A time-series whose partial autocorrelation function is eventually unity belongs to the class of singular processes, which are processes perfectly predictable in mean square.

Because the partial autocorrelation function arises in the section on AR order selection some remarks on its computation are in order. This discussion will also gives some additional insight into the nature of the partial autocorrelation function. For this purpose we introduce the Yule-Walker equations.

When the process concerned is purely a moving average or purely autoregressive there is a simple relationship between the autocovariances and the parameters of the Box-Jenkins process. If the process is autoregressive then the parameters may be obtained uniquely from the Yule-Walker equations which are obtained as follows. Suppose the process is of order m. Then multiplying both sides of the defining difference equation by X for k=1,m and taking expectations obtain

which may be solved directly for p____i=1,m

This leads to a natural estimate for the p's (provided that order of the autoregression is known) obtained by replacing the theoretical autocorrelations their above by sample the estimate (called the method of This moments estimates. estimate) is used in practice as an initial guess for maximum likelihood computations (cf. Box& Jenkins, 1970; Appendix 6.2).

In a similar fashion, it was shown in section 1.2 that for a moving average process of order m the autocovariances are

(5.10)
$$\gamma = (q q + q q + ... + q q) \sigma$$

k 1 1+k 2 2+k m-k m e

where o is the variance of the errors. These non-linear e equations may be solved for q i=1,m. While the non-linear i equations will have multiple solutions it follows from the discussion of section 1.3 that there is a unique invertible Box Jenkins model corresponding to a given autocorrelation function.

We now return to the partial autocorrelation function . It is not hard to see that if the Yule-Walker system has an order than the true one, then the resulting coefficients are the less p̂ in the multiple regression coefficients of X ,..., X . This is because to t-1 t-m minimize the Х on t variance Var[X - X] we compute this variance as a function of t t the p's (a quadratic form with coefficients γ and variates \hat{p}) set its gradient with respect to the p's to zero. This and results in a set of formulas for the partial autocorrelation in terms of determinants. The linearity of this function problem enables an easy computation of the residual variance .. $H = span \{X, \dots, X\}, \quad x = \hat{X} \text{ and } y = X.$ Define n t-1 t-n t t let For the inner product [z,v] = E{ zv } . With these definitions, the problem is equivalent to the Hilbert space problem of finding S = Min || y - x || where X is a subspace of H of dimension m. n хεХ the solution x to this problem is well

It is well known that the solution x to this problem is * * *characterized by [x - y, x] = 0 and that the norm of x is the norm of the projection of y on H. The latter projection is m

equal to

(5.11)
$$E[X(pX,...,pX)] = \sum_{k=1}^{m} p\gamma$$

t 1 t-1 m t-m k=1 k k

where p are the true coefficients of the autoregressive k representation of the process. Thus, applying the pythagorean theorem we obtain

(5.12)
$$S = \gamma - \Sigma p \gamma$$

m o k=1 k k

as the residual of the autoregression of order n.

We can now state the recursive formulae due to Durbin(1960) for computation of the partial autocorrelation function . As noted earlier, \hat{p} is the m partial autocorrelation of X. For m t ease of notation we now drop the ^ on \hat{p} , since the distinction

is clear from the superscript. Durbin's algorithm is a

recursive joint computation S and p for k=1,m. The recursion m k

0

begins

(5.13)

$$S = \gamma \\ 0 \qquad 0$$

$$p_{1}^{1} = \frac{\gamma_{1}}{\gamma_{0}}$$

$$S = \{ 1 - (p_{1}^{1})^{2} \} S$$

and the general steps are

(5.14)
$$p = \{ \gamma - \Sigma p \gamma \} / S$$

 $n n k=1 n n-k n-1$

$$n \quad n-1 \quad n \quad n-1$$

$$p = p \quad -p \quad p \quad for \quad k=1, n-1$$

$$k \quad k \quad n \quad n-k$$

$$S_{n} = \{ 1 - (p) \} S_{n}$$

From the last equation the necessity of the conditions on the set S above are evident. For clearly $|p| \le 1$ and if |p|= 1 this implies S = 0 . When S = 0 this means that the n+1 n+1 process satisfies a homogeneous difference equation of order n+1 or less. Thus the correlation in the definition of ϕ reduces n+1

to a correlation between identically zero variates and so, by definition, is unity. Ramsey's result arises by observing that these equations can be used to generate an admissible sequence of autocovariances from the partial autocorrelation function .

His result is not completely trivial, for in contrast to the weak conditions on a sequence of numbers which make it admissible as a partial autocorrelation function , many sequences cannot be the autocorrelation function of a stationary process. Thus, for example, there are non-vaccuous conditions on a set of m autocorrelations generated by a moving average process in addition to exhibiting cut off. For example given an MA(1) process

> X = e + qe t t t-

compute

(5.15)

(5.16)
$$\gamma = E[X X] = \{1 + q\}\sigma$$

0 tt e

$$\gamma = E[X X] = E\{e + qe \}\{e + qe \}$$

1 t-1 t t-1 t-2 t t-1
$$= q\sigma$$

so that

(5.17)
$$\rho = \frac{\gamma_1}{\gamma_0} = \frac{q}{(1+q^2)}$$

and solving for q in terms of ρ we find q is real iff $|\rho| \leq .5$. Since $|\rho| = .5$ implies q = 1, the resulting process is an invertible moving average iff the inequality is strict.

A general necessary and sufficient condition for a sequence of numbers ρ , i=1,..m to be the autocorrelation function of an i invertible moving average process was given by Wold (1953). Define

(5.18)
$$U(x) = 1 + \sum_{j=1}^{m} \rho(x + x)$$

now define

(5.19) V(z) = U(x) if z = x + x

Then the condition is: V(z) has no root in the interval [-2,2]. We sketch the proof. To begin, we recall that from the

th discussion of section 1.2, any m degree polynomial O corresponds to some moving average model $X[\Theta]$. Such a moving average model is always stationary, but conceivably not invertible. For invertibility of the associated moving average model, the roots of the operator resolution of Θ must lie in the interior of the unit disk of the complex plane. However, the discussion of section 1.3 shows that if one or more roots lie in the interior of the complement of the disk, a new (unique) invertible moving average process may be defined with r replaced by $\frac{1}{\Theta}$ possessing the same autocovariances. Thus any choice of Θ

such that the roots lie off the unit circle corresponds to a unique invertible moving average (and thus a set of m moving average autocorrelations). On the other hand it can be shown th

that any sequence of m autocorrelations corresponds to an m degree Box-Jenkins (moving average) operator Θ . The preceding remarks show that a given autocorrelation sequence thus corresponds to an invertible moving average process iff roots of the associated operator Θ do not lie on the unit circle. Our object is to translate this condition on Θ into a condition on the autocorrelation function (or equivalently, the autocovariance function) of the process.

To do this first recall that the autocovariances of a moving average process are $\gamma = \sigma \Sigma q q$. Inspection of this t jt+jj formula shows that we can represent the autocovariance sequence alternatively as follows. Define the fourier transform of a

sequence q , i=0,m-1

Set
$$x = e$$
 and consider (abusing notation slightly)

(5.21)
$$\Gamma(\phi) = \Theta(\mathbf{x}) \cdot \Theta(\mathbf{x})$$

where the * denotes complex conjugation. The result of this computation is certainly real and can be shown to equal twice the cosine transform of the autocovariance sequence. Thus the autocovariances could be recovered from Γ . This representation useful because on inspecting (5.21) we notice $\Gamma(\phi) = 0^{\circ}$ iff is Thus, as $\mathfrak{O}(\mathbf{x})$ is just the Box-Jenkins operator $\tilde{\Theta}(\mathbf{x}) = 0$. the unit circle, we can detect the presence of evaluated on roots of Θ on the unit circle by inspecting Γ , which is twice the cosine transform of the autocorrelation sequence of length Therefore our aim of translating the condition on the roots m. of Θ has been acheived. To complete the proof, note that the

transformation z = x + x maps x on the unit circle to the real interval [-2,2].

For example, in the model (5.15) we obtain for U(x), on substituting q for q

(5.22)
$$[(1 + q^{2}) + 2q \cos(\phi)]$$
$$= [1 + \rho (x + x^{-1})]$$

In this example it is clear that that to avoid a root on the unit circle in (5.22) $|\rho| < .5$. This is the same condition 1 obtained earlier in this section.

1.5.2 MODEL ESTIMATION

Once the order of the ARMA processes have been specified it becomes feasible to estimate the parameters of the model. The usual method is the method of maximum likelihood or some variation of it. The method of maximum likelihood has several peculiarities in its application to time-series analysis which are treated in this section. In particular we describe the conditional and unconditional likelihoods and sum of squares functions. We then derive the exact likelihood function (Gaussian case) of a moving average which in practice suffices for a general ARMA process.

Suppose we are given an ARMA model

(5.23) $e = X - p X \dots -p X + q e + q e$ t t 1't-1 nt-n 1t-1 nt-m Imagine for now that X and e the values of the observations, -i -i -i and the errors prior to the start of the series are known. Then using the known values for X the recursion can be solved i to give the unobserved values of e. Thus, e is a function of P,Q,X and e which we denote by e [P,Q|X e]. -i -i -i t -i -i Since these values are iid normally distributed (5.24) $p(e, \dots, e) = a|\sigma|^{-N} EXP[-\frac{1}{2\sigma^2} (\sum_{i=1}^{N} e_i))]$ where a is a real constant.

Thus we may write the likelihood conditional on the choice and e as above. The term in the exponential of X

 $S^{-}(P,Q) = \sum_{t=1}^{N} \sum_{t=1}^{2} [P,Q|X,e]$ (5.25)

where
$$P(B) = \sum_{i=1}^{n} B and Q(B) = \sum_{i=1}^{m} A and Q(B) = \sum_{i=1$$

is called the conditional sum of squares function. The conditional likelihood is not yet a true likelihood as it contains unobserved values. In a long series however, this will make little difference because the inaccuracy in the choice of initial values will die out. However, since the likelihood is a useful practical tool we give the derivation of an exact likelihood function. This derivation for a moving average model will also furnish approximate likelihoods for mixed and autoregressive models on inversion and truncation of the operator series at convenient orders.

Thus consider the MA(m) model

nx1

X = e + q e + ... + q et t 1 t-1 n t-m (5.26)where the model is assumed to be invertible is and e X for k>0. We derive the likelihood for a independent of t: series of length n. Because e and hence X are normally distributed we have (5.27) $p([X] | Q, \sigma) = (2\pi\sigma) | M | EXP[-[X] M [X]]$ (5.27) $p([X] | Q, \sigma) = (2\pi\sigma) | M | EXP[-[X] M [X]]$

where
$$[X]_{nx1} = (X, ..., X)_{n}^{T}$$

and $M \sigma$ is the covariance matrix of $[X]$
Thus, to make this expression explicit one must evaluate
 $\begin{bmatrix} X \end{bmatrix}_{nx1}^{T} M \begin{bmatrix} X \end{bmatrix}_{nx1}^{T}$. To do this, first write the n+m equations
 $\begin{bmatrix} nx_{1} & nx_{1} \\ nx_{1} & nx_{1} \end{bmatrix}$. To do this, first write the n+m equations
 $\begin{bmatrix} x \\ nx_{1} & nx_{1} \\ nx_{1} & nx_{1} \end{bmatrix}$.
(5.28) $e = e$
 $e = e$
 $e = x + q e + \dots + q e$
 $e = x + q e + \dots + q e$
 $1 = 1 + 1 = 0$
 $e = x + q e + \dots + q e$
 $n = 1 + 1 + 1 = 0$
 $e = x + q = 1 + \dots + q = 1$
 $e = x + q = 1 + \dots + q = 1$
Now set
(5.29) $[e]_{m+nx_{1}} = (e_{n}, \dots, e_{n})^{T}, [\tilde{e}] = (e_{n}, \dots, e_{n})^{T}$
Dropping the dimension subscripts, partition the resulting matrices as
(5.30) $\begin{bmatrix} 1 & 0 & 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix}$

	100	. 0 0 0	[~]
	0 1 0	. 0 0 0	[ẽ]
[e] =	001.	. 0 0 0	
	•••••	••••	••••
	qq.q 12 n	1 0 0	
	••	t 1 0	
	v v 21 2n	. t t 1	[X]
	vv - n1 nn	. t t 1	

$$= \begin{bmatrix} I & O \\ \cdots & \cdots \\ V & T \end{bmatrix} \begin{bmatrix} \widetilde{e} \\ \cdots \\ [X] \end{bmatrix}$$

where I is the identity, O is the zero matrix, T is lower triangular, and the coefficients of T and V are complicated functions of q obtained on expressing the e as functions of and [X]. Ιt is seen from the first form that e the - i transformation has unit jacobian, so the joint distribution of [X] and [ee] is

> $p\{[X], [\tilde{e}] | Q, \sigma\} = 2 - (n+m)/2$ (2\pi \sigma) EXP [- 1 S(Q, [\tilde{e}])]

where, defining

(5.31)

$$\mathbf{K} = \begin{bmatrix} \mathbf{I} \\ - \\ \mathbf{V} \end{bmatrix} \qquad \mathbf{L} = \begin{bmatrix} \mathbf{O} \\ - \\ \mathbf{T} \end{bmatrix}$$

(5.32) $S(Q, [\tilde{e}]) = \{ L [X] + K [\tilde{e}] \}^{T} \{ L [X] + K [\tilde{e}] \}$

It is now sought to minimize $S(Q, [\tilde{e}])$ as a function of $[\tilde{e}]$. To do this write (5.33) $L[X] + K[\tilde{e}] = \{L[X] + K[\tilde{e}]\} + \{K([\tilde{e}] - [\tilde{e}])\}$ so that if

(5.34)
$$\begin{array}{c} T & T \\ K & [\bar{e}] = K & (L & [X]) \end{array}$$

the quantities in braces are orthogonal vectors and so

(5.35) $S(Q, [\tilde{e}]) = S(Q) + \{[\tilde{e}] - [\bar{e}]\}^T \kappa \kappa \{ [\tilde{e}] - [\bar{e}]\}$ and it follows [\bar{e}] defined by the penultimate equation (5.34) is the minimizing value of [\tilde{e}]. Now by the above it follows that (5.36) $p\{[X], [\tilde{e}]|Q, \sigma\} =$

$$\begin{array}{c} 2 - (m+n)/2 & T \\ (2\pi\sigma) & EXP \left[-\frac{1}{2\sigma^2} (S(Q) + \{ [\tilde{e}] - [\bar{e}] \} K K \{ [\tilde{e}] - [\bar{e}] \}) \right] \end{array}$$

And since

(5.37) $p\{[X], [\tilde{e}] | Q, \sigma\} = p\{[X] | Q, \sigma\} p\{[\tilde{e}] | [X], Q, \sigma\}$ Conclude that (5.38) $p\{[\tilde{e}] | [X], Q, \sigma\} =$

$$(2\sigma)$$
 $|KK|$ $EXP[-\frac{1}{2\sigma^2}([\tilde{e}]-[\bar{e}])KK([\tilde{e}]-[\bar{e}])]$

Therefore, on inspecting the exponential conclude (5.39) $[\bar{e}] = E\{ [\tilde{e}] | [X], Q \}$

This means that $[\bar{e}]$, the minimizing estimate of $[\tilde{e}]$ may be computed from the difference equation by a technique known as backcasting. This is a trick which utilizes the fact that the distributions of X for the model written in forward form are

the same as those for the model written in backward form, though as noted in section (1.4) the assumptions on the error sequence are different. The strategy is to compute the expected preliminary errors e from the expected preliminary values of

X . Once the conditional expectations of the X are known we -i -i

can combine this with knowledge that for a moving average of order m the conditional expectations $E\{e | X \} = 0$ for k > m. t-k Using these facts in (5.26) we can obtain the conditional expectations for e k=0,1-m given the data . Thus for example t-k that E[e | data] = E[X | data]. If we -mwe compute (via 5.26) denote conditional expectation given the data by square brackets this yields [e] = [X], [e] = [X] - q[e], etc.m -m -m+1 -m+1 1 -m The trick is necessary because it will not do to compute the X (or e) by running the original recursion (5.26) backwards with estimated errors ê for i>0. Because of invertibility, in this model inaccuracies die out with increasing time, but they must blow up as the time index decreases. A reversed version of this statement is true of the backward model $\begin{array}{rcl} X &= \epsilon &+ q \ \epsilon &+ \hdots + q \ \epsilon \\ t & t & 1 \ t+1 & n \ t+m \end{array}$ (5.40)is independent of X for k>0. That is, the where ε t-k for ϵ dies out as the time inaccuracy in a guess index decréases. Thus, using (5.40), accurate estimates of the ϵ ~- k may be computed by guessing values near the end of the series and proceding backwards. These values of ϵ are used to t construct the conditional expectations for X i>0 in the

obvious way, and hence the expectations for e $\mbox{ may be found as } -k$

described above. In this manner we arrive at

(5.41) $p\{[X]|Q,\sigma\} = (2\pi\sigma) |KK| EXP[-\frac{1}{2\sigma^2}S(Q)]$

where
$$S(Q) = \sum_{t=1-m}^{n} t$$

(a being conditional expectations of e given the data). t

The necessary complicated function of the q may be found by

iteration. Note however, that maximum likelihood estimates for time-series may not always be found by setting the derivatives of the loq likelihood to zero, since the parameters corresponding to a given covariance structure are not unique. is necessary to constrain the maximization so that the Thus it resulting estimates satisfy the stationarity and invertibility conditions. We skip over the large literature on numerical computation of maximum likelihood and least squares estimates, because it is not relevant to the argument.

purpose here is to make clear the conventional The assumptions involved in the derivation and computation of maximum likelihood estimates for time-series. When the model is correct the likelihood function contains all useful information estimation of parameters (cf. relevant to the It appears that addition likelihood Box&Jenkins,1970). in functions may have a deeper significance than is evident from this derivation. For even if false but close to correct (in some sense) we will see that the (false) likelihood contains valuable information. Thus, although the model is misspecified and the above theory is incorrect, a likelihood function may still be useful.

1.5.3 REVISING THE MODEL

It is possible and indeed likely that the identification stage of the time-series problem will have produced an incorrect choice of model. In other words, the application of the operator

(5.42)
$$Q(B) P(B)$$

(where Q and P are estimated), to X may fail to transform

the series to second order white noise (a sequence of uncorrelated random variables with zero mean). This outcome is likely first because the class of Box-Jenkins models is sufficiently rich that serious errors in model selection are probable. In addition, the requirement that a time-series be generated by a Box-Jenkins model is restrictive, and can be a misspecification in itself. When the parametrization is seriously wrong the maximum likelihood estimates will be invalid and we are in a situation which is at present, imperfectly understood.

There are two possibilities in this regard. The simplest occurs when it is possible to establish that the model is

incorrect through some diagnostic test. To cope with this situation, it is necessary to detect it, to elucidate the nature of the problem and modify the model appropriately. The other possibility is that several models may fit the data equally well, so that while identification has in a sense, failed, no obvious course of action is available. In practice it is usual to choose the simplest model which is acceptable.

To detect a failure in identification, two criteria have been used extensively in non-automatic procedures, namely, Quenouille's(1947) test and Box&Jenkins(1970) Portemanteau test.

If the k+i partial autocorrelation is zero for i=0,m-k then (Anderson,1971) if $\hat{\phi}$ is the 1 estimated partial 1 autocorrelation for the series, $\sqrt{N\phi}$ i=0,m-k have a limiting

joint distributions which are independent N(0,1) under the null hypothesis. Thus, the statistic

$$(5.43) \qquad \qquad \begin{array}{c} m & 2 \\ Q = N \sum \phi \\ k+1 & i \end{array}$$

is distributed asymptotically as a chi-square random variable with m-k degrees of freedom. For an autoregressive order k gaussian model, the hypothesis that the model is correct is equivalent to the hypothesis that all partial autocorrelations after order k are zero. For white noise, this is true for k=0. Thus, this statistic may be applied to the residuals to detect failures of model fitting. The Portemanteau test is an analogue of Quenouille's test with partial autocorrelations replaced by the usual autocorrelations Thus, if ê are the residuals of the m order t

 $\begin{array}{ccc} & & & & & \\ \text{fitted model and } \rho & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ &$

(5.44)
$$B = N \sum_{i=1}^{m} [\rho]^{2}$$

has an asymptotic chi-square distribution with m-p degrees of freedom.

These procedures for model selection terminate, in principle, when the residuals do not fail a test of randomness, and no simpler model has this property.

Several automatic and semi-automatic procedures for model determination have been proposed. The most well known of these is Akaike's (1973) criterion. This criterion is based on information theory and will be discussed in detail in Chapter. III. Applied to the problem of choosing the correct autoregressive order, Akaike's criterion which would select an autoregressive order by minimizing

(5.45) AIC(p) =
$$\ln[\hat{\sigma}] + \frac{2p}{N}$$

where N is the sample size, $\hat{\sigma}$ is the estimate of the residual p variance for the fitted model of order p.

When the hypothesis that the residuals are white noise is rejected, or if the model is unsatisfactory for some other reason, several kinds of situation may obtain. By linear model underfitting we refer to the fitting of a model with an insufficient number of parameters, that is models which set truely non-zero parameters to zero. Since it is usual to begin by fitting the simplest models this sort of problem should show up early in the fitting procedure, and be manifested in autocorrelated residuals and a large estimated error variance. Linear model overfitting refers to a situation in which a model of the correct type is fitted but the fitted model makes provision for non-zero values of parameters which are actually The effect of this situation may be to reduce or increase zero. the estimated prediction error variance and in a non-automatic procedure the only general means of testing for its occurence is fit all simpler models to see whether they are acceptable. to One case of overfitting, called parameter redundancy may however result in a certain pattern in the contours of the sum of, in the parameter space. squares (likelihood) function For example consider the model

(5.46) (1 - aB)(1 - rB)X = (1 - sB)et t

When |r-s| is small the factors on each side of the difference equation will nearly cancel so that the model is nearly (5.47) (1 - aB) X = e

As this is an operator equation, some discussion is required to establish that cancellation is meaningful. To see this, rewrite the model as

(5.48) (1 - aB)((1 - sB) + (s - r)B)X = (1 - sB)et

and invert explicitly to obtain

$$(5.49) \{1 + (s-r)B\{1 + sB + sB + sB + ...\} \} (1 - aB)X = e_{t}$$

This means that the behavior of the series with nearly equal values of s-r will ressemble one another. In particular, changes in (s,r) in the neighbourhood of the line s=r will have little effect on the goodness of fit of the model. Thus if the true parameters r and s are nearly one it will be difficult to discriminate a stationary from a non-stationary model when the more complex one is fitted. However examination of the likelihood will reveal a ridge in the sum of squares function along the line s=r, indicating parameter redundancy. At the expense of a small increase in lack of fit, the offending factors may be dropped from the model. This then, is a practical situation in which the nature of misspecification is evident from the likelihood function. We will see later that other, more theoretical, approaches to the detection of misspecification are also based upon likelihood functions.

1.6 USE OF RESIDUALS TO MODIFY THE MODEL

Suppose that the residuals $\hat{\epsilon}$ from the fitted model t, (5.50) $P(B)X = Q(B)\epsilon$ t t appear to be non-random in the sense of section 1.5.3. Suppose further that one could correctly fit the model (5.51) $R(B)\epsilon = S(B)e$ t t where e are iid random variables with zero mean and the model t defined by R and S is stationary and invertible.

Then, the model could be corrected to

(5.52)
$$P(B)X = Q(B) \{R(B)S(B)\}e$$

or, avoiding infinite series

$$P(B)R(B)X = Q(B)S(B)e$$
t

1.7 SUMMARY

The time-domain theory for linear time-series was reviewed, including the relationship between covariance structure and Box-Jenkins coefficients, and the interpretation of the invertibility and stationarity conditions. Current methods of attacking the linear time-series problem were reviewed. It was seen that these procedures may require judgement in detecting misspecification. There is some indication, however that the likelihood function contains much of the relevant information for this judgement. Thus, it is natural to ask whether automatic procedures based on the likelihood function are a practical way to address the problem of misspecification. If such automatic procedures were marked improvements on the skilled procedures, they would be helpful in the (harder) nonlinear case. Succeeding sections will evaluate this possibility.

II. HOW WELL DO THESE METHODS WORK?

Bora-Senta& Kounias(1979) systematically The paper of compares order determination schemes using Ouenouille's test, Portemanteau test and Akaike's criterion for various the estimation methods in a Monte Carlo simulation. We present some of the results of the study in Tables I-III. The method of estimation used was the method of unconditional least squares. Significance levels for the portemanteau test were determined from the theoretical Chi-squared variate distributed with m-p degrees of freedom where m=10 and p is the order of the fitted Similarly, significance levels for Quenouille's test model. were determined from the Chi-squared variate distributed with mp degrees of freedom with m=6 and p the order of the fitted model. Thus the comparison between Portemanteau and Quenouilles test may not be fair, though neither test can be valid owing to the small sample size (n=20). The model adopted was the model lowest order which passed the respective tests at the .05 of level of significance.

All simulations were involved one of the 6 autoregressive models quoted in Tables I-III. For a large sample size (n=200) fifty simulation runs were performed and all criteria performed well. Quenouille's test was the best performer, failing in only 1/6 cases to select the correct AR order on all fifty runs. The other two criteria, each failed in 2/6 cases to have a perfect record, misclassifying at least 10 runs when they failed.

For the small sample size, performance is less impressive.

Examination of the modal choice of order shows only Quenouilles test has the wrong mode (AR(5)) for the first model. Both portemanteau and Quenouille are wrong for the second AR(1) model, though only Portemanteau is wrong if the choice AR(5) is disallowed. All criteria fail to have correct modal choice of order for the first AR(2) model, each criterion instead selecting the AR(1) model most frequently. For the second AR(2) model correct modes resulted with each criterion.

Performance was further summarized by the use of the index W

(1.1)
$$W(p) = \sum_{\hat{p}=0}^{5} (p - \hat{p}) f / N$$

where p is the true order, and f is the frequency of selection \hat{p} of order \hat{p} .

By this index, Portemanteau wins for the first and second AR(1) models with Akaike a close second. In the other 4/6 cases Akaike wins, decisively in 3/6 cases. Mean values of W are respectively 1.398, 5.607, 2.147 for Akaike, Quenouille and Portemanteau respectively. The relatively poor performance of Portemanteau and Quenouille might be expected here since as noted they are large sample tests, however, it should also be pointed out that the order selection problem considered by Bora-Senta et. al. was a simplified one in which the alternatives were just six autoregressive models (order=0,...,5) which were assumed a priori to be stationary.

When models are allowed to be non-stationary by introducing

powers of difference operators into the defining equations and moving average terms, then the model-discriminating powers of these procedures begins to be taxed even at large sample sizes. Ozaki(1977) compared the models of data series A-F (data given in Box&Jenkins book) obtained by Box&Jenkins with those adopted by Akaike's criterion. We discuss his results for series A, which are the most completely presented.

Series Α from Box&Jenkins(1970) consists of 197 observations concentration readings for a chemical process of taken every two hours. From inspection of the sample autocorrelation function and partial autocorrelation function Box&Jenkins suggest that this time-series might be described by IMA(0, 1, 1) or ARMA(1,0,1) model. Ozaki concluded however an that (7,0,0) and (6,1,0) and (1,1,1) were also candidate models, by the same procedure.

The results of fitting these models to the data are given in Table IV. All these models passed the Portemanteau test.

Akaikes criterion was useful in automatically weeding out a large number of other models (which likely would have failed the eyeball test) but somewhat less helpful in discriminating amongst models which were good candidates. Values of Akaike's criterion for the various models considered are given in Table V. While Box&Jenkins concluded that each of (1,0,1) and (0,1,1) fitted equally well Akaikes criterion gives (1,1,1) as the model of choice amongst

(1.2) (i) { $(p,d,q), d=0,1 \quad 0 \le p,q \le 2$ }

and (3,0,3) as the model of choice amongst

(ii) { $(p,d,q), d=0,1 \quad 0 \le p,q \le 9$ }

However, the (3,0,3) model failed the Portemanteau test.

It can be seen that the models adopted by Box&Jenkins exhibit AIC values near the minimum. On the other hand, while Ozaki claims that the automatic procedure's results are similar Box &Jenkins, examination of the numbers seems to to those of show that Akaike's procedure recomends 9 of a possible 18 models (i) as well as recommending the higher order models of Table IV. In particular 9 models had AIC values between -449 and -451 whereas, for comparison, the models selected by Box&Jenkins had AIC(1,0,1) = -450.2 and AIC(0,1,1) = -448.16. See Table V. Now it is very well to say that one should apply Ockham's razor but in the first place, 3 of these models have the same number of parameters, and in the second, it is not sensible to do this without regard for the explanatory power of the model Akaike's criterion, which is supposed to formalize the tradeoff, between explanatory power and parsimony is what gives us these results in the first place. Thus the data seem ambiguous. Examination of the other results in Ozaki 's paper suggest that situation is not untypical. This problem is this less disastrous than it appears, since in several cases the different models differ only by the addition of nearly cancelling factors on each side of the equation, which implies that very similar forecast would be obtained from these apparently different

models.

The conclusion appears to be that while practice shows that useful prediction can be obtained, optimal prediction may be difficult, for presumably not all of the models the data admits will be correct. Automatic procedures give a clear reduction in labour, though not necessarily an improvement in performance. Ozaki remarks that a classical Box-Jenkins identification by inspection of the autocorrelations etc, can produce visual reasonable results, but only by 'skilful analysis which requires expertise perhaps not commonly possessed by ordinary practicing statisticians.' Both automatic and skilled procedures invariably in the acceptance of some useful model. However, the result results of either procedure are usually far from unique when the class of admissible models is large.

Table I - AR order Selection

1

Model		X(t) = .3 X(t-1) + e(t)	
order selecte	Akaike	Quenouille	Portemanteau
0 1 2 3 4 5	360 60 30 40 10 0	140 10 0 30 0 320	480 10 10 0 0 0
W	1.28	10.76	.98
Model		X(t) = .8 X(t-1) + e(t)	
0 1 2 3 4 5	70 300 70 50 0 10	40 120 10 10 30 290	310 180 10 0 0
W	1.0	10.6	.64

Method of maximum likelihood was used for estimation of coefficients. Sample series were of length 20.

Table II - AR order selection

Model	X(t) = .3 X(t-1) + .5 X(t-2) + e(t)			
order selected	Akaike	Quenouille	Portmanteau	
0 1 2 3 4 5	240 105 110 35 5 5	145 55 90 45 45 210	410 90 0 0 0 0	
W	2.33	5.5	3.28	
Model	X(t)	= .8 X(t-1)9 X(t-2)	+ e(t)	
0	0	0	90	
1	0	Ō	0	
2	400	300	410	
2 3 4 5	70	40	0	
4	20	30	0	
5	10	130	0	
W	.48	2.6	.72	

Table III - AR Order Selection

Model	X(t)= 1.6	X(t-1)79 X(t-2) ·	+ .12 X(t-3) + e(t)
_	Akaike	Quenouille	Portemanteau
order			
select			
0	0	0	40
1	110	30	380
2	320	230	80
1 2 3 4 5	60	20	0
4	10	10	0
5	0	210	0
W	1.54	2.4	3.76
Model	X(t) = 1.8 X	X(t-1) - 1.14 X(t-2)	+ .272 X(t-3) + e(t)
0	0	0	40
	130	20	310
1 2 3 4 5	320	255	150
3	25	60	0
4	20	35	0
5	5	130	0
W	1.76	1.78	3.50

r

Table IV - Models for Series A

Model	N-d	Fitted Model		Portemanteau Statistic	đ
(7,0,0)	197	1355B186B ² 023B ³ 029B ⁴ 021B ⁵ 085B ⁶ 203B ⁷	.0925	13.61	18
(6,1,0)	196	1+.624B +.423B ² +.386B ³ +.339B ⁴ +.346B ⁵ +.238B ⁶	.0945	15.05	19
(1,1,1)	196	(1218B)(1-B)X(t) =(1825B) e(t)	.0985	23.37	
(1,0,1)	197	(1915B)X(t) =(1583B)e(t)	.0977	26.45	
(0,1,1)	196	(1-B)X(t) =(1705B)e(t)	.1007	29.87	

All the above models passed the portemanteau test with p=.05, k=25 (d.f.=k-p-q for model (p,d,q)).

Table V - Values of AIC for Various Models

Model	AIC	Model	AIC
(1,0,1)	-450.26	(1,1,1)	-450.55
(2,0,1)	-450.26	(2,1,1)	-449.72
(1,0,2)	-449.73	(0,1,2)	-449.53
(2,0,2)	-448.67	(1,1,2)	-449.23
(2,0,0)	-445.73	(0,1,1)	-448.16
(1,0,0)	-434.58	(2,1,2)	-447.47
(0,0,2)	-423.76	(2,1,0)	-430.56
(0,0,1)	-403.08	(1,1,0)	-425.57
(0,0,0)	-358.76	(0,1,0)	-390.40
*(3,0;3) *(7,0,0)	-456.16 -450.89		

*(6,1,0) -450.89

III. METHODS OF SELECTING AR ORDER

Statisticians are often faced with the problem of choosing appropriate dimensionality of a model that will fit a given the set of observations. The problem is particularly acute in the analysis of time-series data. This is because, in practice, the time-series problem is often posed to the statistician in the form we have posed it, without much knowledge of the underlying this situation, the problem must be attacked structure. Ιn frontally despite its formidable nature. Thus, whether the aim structural characterization or optimal prediction, the is time-series parametrization (misspecification) problem in basic. It is the purpose of this section to analysis is describe some recent work on this problem.

It is part of statistical folklore (cf. Schwartz,1978) that the method of maximum likelihood performs badly in the selection of over-parametrized models. In particular there is a bias toward the selection of models with too high a dimensionality. With the method of least squares for example, the minimization of prediction error on the data set leads to non-zero estimates of parameters which are zero and biases the estimate of mean square error.

Quenouilles test and the Portemanteau test were derived with no particular alternative in mind and in practice are applied repeatedly (which is not a recommended approach to hypothesis testing) however they represent a considerable improvement on simple maximum likelihood. In particular these

tests may be understood as likelihood ratio tests against particular alternatives so that they are maximum power tests under the usual conditions(Hosking, 1978). Thus, the use of the Portemanteau test statistic

(1.1)
$$m \qquad 2 \qquad 2$$
$$N \qquad \Sigma \qquad \hat{\rho} \qquad \chi$$
$$k=1 \qquad k \qquad m-p$$

is equivalent to a likelihood ratio test of

(1.2) H: P(B)Y = a

versus H: T(B)P(B)Y = a

where a and å are white noise processes possibly with t t different variances. Thus the alternative is an autoregression fitted to the residuals or from another point of view, an autoregression of order m+p with p roots restricted to the same values as the null hypothesis.

Quenouille's test statistic (Whittle, 1952)

(1.3)
$$N \sum_{k+1}^{m} \sum_{i=1}^{2} \sum_{k+1}^{2} \sum_{i=1}^{2} \sum_{i=1}^$$

is equivalent asymptotically to a likelihood ratio test of
(1.4)
H: P(B)Y = a
0 t t

versus H: S(B)Y = a

where deg P = p and deg S = p+m for small variations about the null hypothesis (the size of the variation affects the approximations used in deriving the results).

As we noted, Quenouille's test and the Portemanteau test

t

are usually applied repeatedly. Therefore this approach is not completely satisfactory. Another set of approaches may be thought of as based on the following empirical observation, well known to statisticians . Namely, that if spurious explanatory variables are added to a regression model, estimated mean square error first exhibits a drop as the first few are added, then after a certain number have been added, mean square error begins to rise. This suggests that some function of estimated mean square error might be a sensible order selection criterion. that the Gaussian case Note in the mean square error is an likelihood. estimate of loa There have been a number of to formalize such procedures of which we will discuss attempts the most systematic.

To illustrate the phenomenon we mention the simulation study of Jones(1976). In this study, autoregressive models of various orders were fitted to white noise for samples of size 20 and 40 and the results averaged over 100 realizations. Possibly because white noise was being fitted, no initial drop in mean square error was evident as the order of the autoregression was increased. However, estimated mean square error rose drastically as the order, p of the fitted autoregression approached the length of the sample. Thus, estimated error doubled as p approached 15/20 and 25/40 respectively.

Akaike's(1969) final prediction error criterion was based directly on this phenomenon. It is

(1.5)
$$FPE(p) = \hat{\sigma}^2 (1 + (p+1))$$

 $p = N$

where

$$\hat{\sigma}_{p} = \sum_{t=t}^{N} (\hat{e}_{t})^{2} / (N - 1 - p)$$

[p] and ê denote the residuals from the fitted autoregression of

order p. The factor multiplying the mean square error is calculated to compensate for the effect of errors in estimating that the parameters so FPE is an estimate of the average prediction error when a fitted model of order p is used. Jones(1975) and Gersh&Sharp(1973) presented favourable simulation and practical results for the use of this criterion in the selection of autoregressive order. Minimization of final prediction error turns out to be a special case of minimization of Akaike's(1973) information criterion which, because of its widespread use and interesting theoretical rationale, we discuss in some detail.

Akaike's information criterion is

(1.6) -2log likelihood + 2 (no. of parameters adjusted) Heuristically, the idea is to measure the amount of structure in the estimated model (ie. in its probability density) and to choose the fitted model of a class which exhibits the least structure. Such a procedure, made precise, is conservative in the sense that it attributes the least possible structure to the data, consistent with membership in a certain class of models. It will turn out that the definition of structure employed will assign less structure to models with small error variance, and

models to with small numbers of parameters. While the notion that more more structure should parameters mean cause no difficulty, that a small variance should correspond to less structure may seem problematic. It is useful to think of the case of a discrete random variable in which a large variance might mean that more different, widely scattered values occur in the distribution. Alternatively Akaike's criterion can be thought of as one possible choice of loss function in the tradeoff between reduced (sample) error variance and number of a formalization of the principle parameters, that is, of parsimony. In the case of a one parameter family of distributions, Akaike's method obviously reduces to the method of maximum likelihood.

We present a series of heuristic arguments justifying now Akaike's criterion, beginning from the desire to develop a model selection procedure which is conservative in the above sense. begin, we argue that a parameter in a model represents То structure only because it indexes significant deviations from some reference model. So it is reasonable to think of small deviations from reference as implying less structure, and larger deviations as implying more structure. It is natural to take model as reference, so that the attribution of least the true structure leads to a correct choice when the model is known. These considerations mean that we need to look for an index of the deviations of a putative model from the true one.

The well known fact that for large N the maximum likelihood

estimator approaches that theoretical minimum variance estimator suggests that the likelihood function is a sensible index of deviation of the model parameters from the true values. This motivates the assumption that the information (structure) function I is proportional to

$$(1.7) \quad -I\{f,\Phi\} = \int \Phi\{ f(x|\theta) / f(x|\theta) \} f(x|\theta) dx$$
$$= E\{ \Phi\{ f(\theta) / f(\theta) \} \}$$

where f is the likelihood function of the observations. The choice of Φ is equivalent to the choice of a loss function for decision making, for if the loss function were chosen as $\lambda\{\theta, \theta\}$ o and it were sought to minimize $E[\lambda[\theta, \theta]]$ then this expectation o could be defined as the previous integral for some Φ provided that

$$\frac{\partial f}{\partial \theta} \neq 0$$

using the inverse function theorem. The choice of the, appropriate Φ follows from the hypothesis that, if X and Y are independent random variables then

(1.9) I(X + Y) = I(X) + I(Y)

where I(.) is the information function

Since the joint density (X,Y) is the product of the densities of X and Y, it immediately follows that the appropriate choice of Φ is (1.10) $\Phi\{r\}=c \log(r)$ 68

iff f=g a.e. since this is true when f and g are simple functions, using a lagrange multiplier argument.

In practice, the true density will be unknown, so that we will work with the natural estimate of $\Phi(\theta)$ namely,

$$(1.12) \qquad \qquad \frac{1}{N} \sum_{i} \log f(x \mid \hat{\theta})$$

It is now possible to formulate the problem of order estimation precisely. Suppose then that we want to choose a density $f\{x | \theta\}$ which best approximates the density $f\{x | \theta\}$

for $\theta \in \Theta$ and $\theta \in \Theta$ where Θ is a k dimensional subspace of Θ . Thus, the problem of order choice is replaced with a problem in density estimation. We define θ and θ to be the maximum likelihood estimates of θ on restriction to the respective linear spaces. Note the notational distinction between θ which, k is the maximum likelihood estimate of heta projected onto Θ , and k θ which is the maximum likelihood estimate in the restricted Θ. Now the ultimate object is to obtain the parameter space ŀ

estimate $\underline{\theta}$ which minimizes $I(f(\theta))$. However this is not completely straightforward. If k is fixed and $\theta \in \Theta$ then the o k result of maximizing the function

(1.13) $F(\theta) = -\frac{2}{N} \sum_{i} \log f\{x \mid \theta\}$ will be just the ordinary maximum likelihood estimator. However
when $\theta \notin \Theta$ the statistic $F(\theta)$ cannot estimate $I(f(\theta))$ or e^{A} or e^{A} and some adjustment becomes
necessary. This is done as follows: define the matrix
(1.14) $J = E[\frac{\partial \log f}{\partial \theta[i]} \{x \mid \theta\} \cdot \frac{\partial \log f}{\partial \theta[j]} \{x \mid \theta\}]|_{\theta}$ $= -E[\frac{\partial}{\partial \theta[1]\partial \theta[m]} f\{x \mid \theta\}]|_{\theta}$

for i, j=1,...L where L is the dimension of Θ . This is the Fisher information matrix.

It is ordinarily positive definite, being the variance covariance matrix of the random variable(s)

 $(1.15) \qquad \qquad \frac{\partial f\{\theta\}}{\partial \theta[i]} \frac{1}{f\{\theta\}}$

Let $\frac{\theta}{2}$ be the maximizing argument of $I(\theta, \theta)$ restricted to Θ .

Finally define the J inner product as $\{u,v\} = u J v$ then (i) it can be shown that approximately

(1.16)
$$I(\theta, v) = -\frac{1}{2} | \theta - v |$$

and

(ii) $\underline{\theta}$ is approximately the projection of θ on Θ in the J inner o k

product. Thus (approximately) by the pythagorean theorem $|| \theta - \theta||_{=}^{k} || \theta - \theta ||_{+}^{2} + || \theta - \theta ||_{+}^{2}$ (1.17) and thus $(1.18) \quad -I(\theta, \theta) = -I(\theta, \theta) + N||\theta - \theta||$ We want to estimate I(θ , $\underline{\theta}$), but can only hope to observe а function equal to a constant (relative to k) plus N I (θ, θ) (1.19) $= -\left[\sum \log f\{x \mid \theta\} - \log f\{x \mid \theta\} \right]$ on restriction to each subspace. However, it can be shown (iii) that asymptotically $N | | \theta^{k} - \theta | |_{\tau}^{2} \chi_{t}^{2}$ (1.20)the statistic (1.19) can be made unbiased (up to a, So that constant in k) for $I(\theta, \theta)$ by adding E χ . Since k is just the number of independently adjusted parameters this reasoning leads to Akaike's statistic. Following Akaike(1973), we now sketch the proof of the

necessary facts. Approximately here means ignoring terms of order $|| \theta - v ||^3$. Fact (i) follows from a Taylor expansion

of $I(\theta, v)$ assuming that

$$(1.21) \qquad \qquad \frac{\partial I(\theta_0, \theta)}{\partial \theta} \qquad = 0$$

Fact (ii) follows from the first fact and the variational characterization of projections. To get the third fact procede as follows. Assume that Θ contains θ and recall θ denotes the o

unrestricted maximum likelihood estimate of θ . Assume

(1.22) $\Sigma \frac{\partial \log f(x[i]|\theta)}{\partial \theta[m]} = 0$ for m=1,...,L

(1.23) $\Sigma_{i} \frac{\partial \log f(x[i]|\theta)}{\partial \theta[m]}^{k} = 0$ for m=1,...,k

Then using the mean value theorem

We assume that the matrices in the second terms converge in probability as N $\rightarrow \infty$ to the corresponding Fisher Information matrices.

On differentiating the above expressions with respect to $\underline{\theta}$ and dividing by N it is seen that for l=1,k

$$= \sum_{m=1}^{L} \left(\frac{\theta[m]}{\mu} - \theta[m] \right) \frac{1}{N} \sum_{i=1}^{N} \frac{\partial^2 \log f(x)}{\partial \theta[m] \partial \theta[1]} \left(\frac{\theta}{\mu} + \rho(\theta - \theta) \right)$$

We need the first member in this sequence of equations to tend to zero as N -> ∞ . This is reasonable given some mild conditions on f since $\underline{\theta}$ maximizes the structure function I by hypothesis. It would hold for example if $\underline{\theta}$ is a normal mean or variance. If so then the last equations imply that asymptotically (for N large)

(1.26)
$$J \left(\begin{array}{c} k \\ \theta - \theta \end{array} \right) = J \left(\begin{array}{c} \theta - \theta \end{array} \right)$$

kxk Lxk

or in other words

(1.27)
$$\Pi J \Pi \left(\begin{array}{c} \theta - \theta \end{array} \right) = J \Pi \left(\begin{array}{c} \theta - \theta \end{array} \right)$$
LxL LxL LxL

where Π denotes the J-orthogonal projection onto Θ , J is the k Fisher Information matrix.

Since the equation (1.27) characterizes the projection we conclude that asymptotically

(1.28)
$$\Pi(\underline{\theta}-\overline{\theta}) = \underline{\theta}-\overline{\theta}^{K}$$

which may be written

(1.29)
$$(\underline{\theta} - \overline{\theta}^{k}) = (\underline{\theta} - \overline{\theta}^{k})$$

We can now use this equality to derive the distribution of $N \mid | = \frac{k}{\theta} - \frac{2}{\theta} \mid \int_{J}^{2} from the assumption that N = \frac{1/2}{(\theta - \theta)} = \frac{D}{--->}$ $N(\underline{0}, J)$. First note that for any Y in Θ , if Y = MVN(0, K) (multivariate normal) this implies there exists a matrix M such that that MY is standard normal of dimension 1 and M = K. Since $(1.30) \qquad ||Y||_{K}^{2} = Y K Y = Y M M Y$

it follows the norm (1.30) is χ^2 . Thus, we need only show that in a space of dimension L, if Π is the J-orthogonal projection (for J positive definite) onto a subspace Θ of

dimension k < L, and $X \sim MVN(0, J)$ then L

(1.31) $\Pi X \text{ is } MVN(0, J|).$ $k |\Theta$

This is certainly true if J is diagonal, since the density of an independent MVN distribution will factor. Let P diagonalize J. Then in P coordinates the same result holds. This is equivalent to the projection result (5.31). Finally because J is positive definite the L-dimensional distribution, and thus the k-dimensional distribution is non-degenerate. Thus k 2

we conclude that $||\underline{\theta} - \overline{\theta}||$ is distributed as χ as required. J

Another technique which was formulated to solve the parametrization problem is known as cross validation. It is based on the notion of optimizing final prediction error estimated in a slightly different manner. The procedure is as follows for each candidate parametrization. For each data point { x,y} compute the maximum likelihood estimate for the parameters of interest with this data point omitted. Thus ϑ for the omitted ith point. Measure the compute prediction error by the squared difference from the observed value, or more generally log f((x,y)| ϑ). Repeat for i i -i

each sample point and add prediction errors. Chose the parametrization giving the least summed error. Stone(1977) shows that this technique is asymptotically equivalent to Akaike's criterion under assumptions similar to those above. Shibata (1976) showed that Akaike's technique is not consistent, and that it over-estimates autoregressive order with non-zero probability as N $\rightarrow \infty$. This implies that cross-validation also

fails to be a consistent technique. Hannan(1979) has modified Akaike's criterion so that the result is consistent.

IV. STRUCTURAL PROPERTIES: LINEAR EXAMPLES

We have seen that data may be ambiguous with respect to the correct choice of parametrization. This suggests that it will advantageous to try to employ any available structural be information even though it may not be in convenient form. It suggested that recent proposals for non-linear models will be ought to be examined with this emphasis rather than as general data prescriptions. That is, in contrast to the widely followed emphasis of Box&Jenkins on modelling data for which little or no priori structure is postulated, we suggest that time-series а data analytic models must be chosen in accord with known general facts about the natural process. Succeeding sections will deal with how (and why) non-linear models may be utilized for this purpose.

We first show that within the class of linear models it is possible to augment the useful information available by an appeal to relatively vague structural knowledge. The fact that these examples due to Granger& Morris(1976) are not better known, illustrates the neglect of structural data by time-series analysts. Suppose it is reasonable to assume that an underlying time-series X is Markov. In the linear case this means that partial autocorrelation function cuts off and $\phi = 0$ for the i>1, that is no additional useful information for prediction is known. This might be the case for instance exists once X t-1

in a physical model for which all theoretical models are differential equations, and so are all Markov. Suppose however that an independent source of white noise is added to the observations.

The time-series described above will be identified by Box-Jenkins procedures as an ARMA(1,1) process. To see this define (1.1) (1 + aB) X = w

where

a is a real constant such that $|a| \le 1$, B is the backward shift operator, and w and v are uncorrelated white noise processes. t t

Suppose that we are observing

then:

(1.3) (1 + aB) Z = (1 + aB) v + w (= m)t t t t

Computing the autocorrelations for the process on the right hand side (=m) of the latter equation we get:

(1.4)

$$u = E(m) = (1 + a) \quad o + o = 0$$

$$u = E(m m) = a \sigma$$

$$1 \quad t \quad t = 1 \quad v$$

and

We will construct a putative representation in the form (1.5) (1 + c B) Z = (1 + d B) e (=n)t t t t by choosing the parameters c,d and σ^2 to produce the values of u computed above . To this end we let c=a, and define the i other parameters via the equations

(1.6)
$$u = d \sigma$$

and

$$u = (1 + d) \sigma$$

to complete the construction. However, observe that these equations can always be solved, and that the resulting processes (n) have the same covariance structure as (m). Thus the t construction works and generates a process with the same covariances as (1.2).

In contrast, it may not be possible to represent ARMA(1,1) as AR(1)+ white noise . The realizability (necessary) conditions may be stated as follows.

Let the u be given by the ARMA(1,1) process (1.5), as computed above. Define (1.7) p = u / u1 1 0

Then

(1.8)

$$p = \frac{22}{a} \frac{2}{\sigma} / ((1 + \frac{2}{a}) \frac{2}{\sigma} + \frac{2}{\sigma})$$

$$= \frac{2}{a} \frac{2}{(1 + \frac{2}{a})}$$

$$= \frac{2}{c} \frac{2}{(1 + c^{2})}$$

long as these inequalities are satisfied, some choice As 2 2 the parameters σ and σ will allow the construction of а candidate AR(1)+noise model. Therefore (1.8)is а realizability condition the ARMA process for on an AR(1)+noise representation.

Thus, under the original assumptions, this condition should be satisfied by parameter estimates for AR(1)+noise and its application will simplify computation of maximum likelihood estimates. It is possible moreover, to reduce the number of parameters in some cases. For example, using methods like these it can be shown that AR(2)+noise identifies as ARMA(2,2). Writing out the respective models it may be seen that ARMA(2,2), involves five parameters and AR(2)+noise only four. For if

$$\begin{array}{cccc} (1.9) & & & Z = X + Y \\ & & & t & t & t \end{array}$$

 $(1 + a B + a B) X = \epsilon$ and $Y = \eta$ 1 22 t t t t

then

(1.10) $(1 + a B + a B) Z = (1 + a B + a B) \eta + \epsilon$ 1 2 t 1 2 t t

where ϵ and η are independent white noise processes, while the t t

ARMA model is

(1.11) $(1 + c_B + c_B) Z_t = (1 + d_B + d_B) \zeta_t$ where ζ_t is a white noise series. This fact is reflected in the realizability conditions for representing ARMA(2,2) as AR(2)+noise . As before let the symbols u denote the successive autocorrelations of the process i on the right hand side of the above equations. Define (1.12) p = u / u1 = 1 = 0

and

p = u /u 2 2 0

It may be shown that a given ARMA(2,2) process is expressible as AR(2)+noise provided

(1.13)
$$c/(1+c+c) > p \ge 0$$

2 1 2 2

and

c/c(1 + c) = p/p2 1 2 2 1

Because these conditions involve an equality, the dimension of the parameter space has been reduced. This is a gain in simplicity which, if justified, will reduce the final error of prediction.

V. <u>REPRESENTATION THEORY FOR WEAKLY STATIONARY</u> PROCESSES: IMPLICATIONS FOR NON-LINEAR TIME-SERIES MODELLING

5.1 WOLD'S THEOREM AND SINGULAR PARTS

We contend that to make full use of non-linear time-series models the use of available structural information is essential ¹. The insistence on structural restrictions is closely related to prediction since enlarging the class of admissible models means that the errors made in choice of parametrization become more important. These errors affect the final prediction error. It seems clear from our discussion of order choice that this phenomenon leads to difficulties in AR the case of linear identification, and that in the non-linear case the problem will be worse.

This observation is of little use without some clarification however, since the notion of structure is multifaceted. To the scientific investigator, structural models are those which express interpretable relationships between a

Some non-linear time-series modelling might be possible with minimal resort to structural information if a useful nonparametric approach could be formulated. There is already a non-parametric regression Stone, 1977; literature on (eg. Friedman& Stuetzle, 1981). To make this work for time-series one might start with an extremely general autoregressive model. For example, a time-series might be modelled as the solution of an autonomous continuous time stochastic differential equation Ludwig, 1975) and the parameters (cf. Jones,1981; (SDE) estimated locally using techniques of non-parametric regression. The special case when the underlying deterministic flow is a gradient field was investigated and partially solved, in work related to this thesis. Unfortunately this investigation had to be abandonned due to time constraints.

priori state variables in the system of interest. То the statistician structure means the specification of а data prescription. In order to reconcile these views we need to formulate the idea of the structure of a process in a more precise fashion.

Recall that our definition of а time-series explicitly mentions a deterministic relation forced by errors, so that conceivably, not all weakly stationary processes may correspond time-series. This section is devoted to clarifying the to probabilistic structure of weakly stationary processes and their representations as time-series . this From discussion conclusions are drawn about the connection between interpretable relationships and data prescriptions, well as the as appropriateness of different types of model.

For stationary Box-Jenkins models it was shown that a moving average representation can always be computed. However this representation is not completely general. A weakly stationary time-series consists of the sum of a moving average part and a deterministic part. Specifically, Wold's theorem guarantees that if x is weakly stationary, then it can be

where x is deterministic, in the sense that it may be forecast 1t linearly from previous values with zero mean square error, and

$$\begin{array}{ccc}
 & & & j \\
 x & = \Sigma & b & \eta \\
 2t & j=0 & t-j
\end{array}$$

where η is a zero mean uncorrelated sequence of random t variables.

The deterministic part is also known as the singular component. Its appearance in a weakly stationary process may seem puzzling because it seems to imply that the level of the series depends on time. This is indeed the case for fixed ω . However, the definition of mean-stationarity refers only to an ω average. While the presence of a non-zero singular part generally implies that the time-series is not ergodic, ω averages may still remain stationary. In fact it is possible to construct non-ergodic stationary time-series whose singular parts may be arbitrary finite fourier series.

To show this it suffices to construct a stationary timeseries which has singular part $A \cos(\pi bt) + B \sin(\pi bt)$ because any more complicated one may be constructed by summing independent series of this form. Let C and C be independent

normal random variables with variance σ with mean zero. Consider the time-series

> $X = C \cos(\pi bt) + C \sin(\pi bt)$ t 1 2

for all t , only C being random i Obviously this series has mean zero. Computing E[X X] we t t-k

find the latter (1.2)

 $= \sigma \left[\sin(\pi bt) \sin(\pi b(t-k)) + \cos(\pi bt) \cos(\pi b(t-k)) \right]$

$= \cos(\pi bk)$

so that the process has stationary covariances. Since the magnitudes of the covariances computed using a time average depend on ω in this example, this process is not ergodic. Since for each realization this time-series is periodic it could not be discriminated from a non-stationary one in practice. For this reason it is usually assumed that the overall stochastic process being sampled has no singular component.

This does not mean that in practice time-series do not have singular parts. Ιt is important to point out that Wold's theorem is a theorem about stochastic processes, and not about time-series. Specifically, the process constructed by starting discrete time-series in its stationary distribution, is а different from the process obtained when the initial value is fixed. This distinction is usually not made in discussing timeseries but it can make the difference between а non-zero singular component and none.

A time-series conditioned on a fixed initial value usually contains a singular component which must die out more or less rapidly as the stationary solution is approached. For example, in the linear case the solution may contain a sine-wave which damps more slowly as a root of the equation approaches a point on the unit circle. Thus, as the root approaches a point on

the unit circle, the time-series looks more and more singular (non-stationary). However, no matter how slowly the singular part may die out, if it tends to zero then the recursion admits stationary solutions, and the series is stationary in principle. In speaking of time-series then, 'near non-stationarity' means that a slowly dying singular component may appear in the solution of the defining equation, for fixed initial values. For non-linear time-series the 'almost non-stationary ' case seems the most interesting.

The generic non-linear time-series is different in character for the following reason. To construct a nonstationary solution in the linear case, one need only find a non-zero deterministic solution for the difference equation

$$(1.3) \qquad \Theta(X) = 0$$

which specifies the time-series. The non-stationary solution is constructed by adding the latter to any stationary solution and applying linearity. The result is truely non-stationary. In the non-linear case this construction generally fails. It seems, reasonable to speculate that most (in some sense) useful nonlinear recursions will admit true non-stationary solutions only with zero probability. For input errors will interact with a deterministic solution and cause the system to drift away from it.

In fact it is possible to give sufficient conditions for stationarity and ergodicity of a non-linear time-series via a result of Tweedie(1975). The relevant result is as follows.

Let {X } be a Markov chain taking values in a normed space S n

(such as R[°]) with temporally homogeneous transition probabilities

(1.4)
$$P(x,A) = P[X \in A | X = x]$$

where $x \in S$, and A is a Borel set. Assume that there is a σ -finite measure ϕ such that that whenever $\phi(A) > 0$

(1.5)
$$\sum_{n}^{-n} \sum_{n}^{n} \Sigma_{x,A} > 0 \quad \text{for every } x \in S$$

That is almost all of S is reachable from any $x \in S$. Further assume that for every Borel set A of S, P(x,A) is a continuous function of x in the topology of the norm. The result gives sufficient conditions for the existence of a finite invariant measure μ , that is a measure such that for each Borel set A

(1.6)
$$\mu(A) = \int \mu(dy)P(y,A)$$

Specifically,(i) μ with this property exists if (1.7) E[||X|| | |X| = x] is bounded for all x n+1 n

and

(ii) there exists a compact set K with $\mu(K) > 0$ and $\epsilon > 0$ such that (1.8) $E\{||X||| - ||X||| |X| = x \} < -\epsilon$ n+1 n n

whenever x & K .

Whenever such a measure μ exists it can be shown (Tweedie ,1974) that except for $y \in N$ where N is a null set of S (ie. $\phi(N) = 0$)

(1.9)
$$\frac{1}{n} \sum_{i=1}^{n} P(y,A) \longrightarrow \mu(A) \text{ as } n \longrightarrow \infty$$

So that the Markov process settles down eventually into its invariant distribution. This is equivalent to our definition of stationarity.

Most time-series are not Markov because X depends on N X ...X . Moreover, time-series do not always possess n-1 n-k temporally homogenous transition probabilities, as was shown by the example of the second order autoregression with roots on the unit circle in section 1.4 . However, when the time-series defined by a non-linear recursion has temporally homogeneous transition probabilities it is possible to apply Tweedie's result. Consider an autoregression satisfying

(1.10) $X = F(X, ..., X) + \epsilon$ n+1 n-1 n-k n

where the error series ϵ is normal. If we define n

(1.11) y = [x, ..., x]n = [x, ..., x]

then the processes $\{Y \}$ are Markov whenever p > k. Let np n=1 || . || be defined by the summed absolute values of the p-tuple

and let $S = R^{p}$. We claim that the following conditions are sufficient for stationarity.

(1.12)
$$E(|X| | |Y|)$$

n+1 n

are bounded on R .

(iii) There exists a ball B of radius R such that

(1.13)
$$E[|X| | |Y|] < R/k$$
 outside B
n+1 n

As $X = F(Y) + \epsilon$, the condition (1.13) depends on the n+1 n t

distribution of the errors as well as on F.

We prove this result by first verifying Tweedie's condition

for the process {Y } under these assumptions. Without loss nk n

of generality we may take R such that

(1.14) $E\{|X| | |Y| \} < R/k \text{ for all } Y$.

To do this take $R' = Max(R, k \sup E\{|X | |Y \})$ using n+1 nk

boundedness of the conditional expectation. Finally, define

 $\Delta_{nk} = || \Upsilon_{nk}^{k-1} ||.$

We now show that for each p, Δ > R implies

nk

E[|X ||Y] < R . Using the definition of the norm, this nk+p nk

k – 1 will establish Tweedie's condition for the process {Y } since nk k-1 consists of k coordinates of this form. Y (n+1)kif exceeds То begin observe that Δ R then nk Y] < R/k by (1.13). Note also that (1.13) means that nk E[X kn+1 for each p=1, k-1, Δ > R implies kn+p-1k-1E[|X ||Y] < R/k. (1.15)Using the Markov property, for p≥1 we can write E[|X| | |Y|]kn+p+1 nk (1.16) $= E[E[|X_{kn+p+1}||Y_{kn+p}||X_{kn+p}]]$

Applying (1.13) to the inner expectation we see that it is always strictly less than R/k as required.

We now argue that because the transition probabilities are not time-dependent the stationary distribution we obtain does not depend on how the X were grouped to produce a Markov n process and thus that there is a stationary distribution for X.

While such series are stationary, there are many examples of stable non-linear recursions with complicated slowly dying solutions. Because non-linear recursions admit a rich variety of singular parts in the associated time-series, non-linear autoregressive models are appropriate when knowledge of a complicated singular component exists.

5.2 A NONLINEAR AUTOREGRESSIVE/MOVING AVERAGE DUALITY

It will be shown that linear time-series models compare well with non-linear models where series of moving average type are concerned. However it is difficult to make precise the claim that non-linear autoregressive models are generically more singular than linear ones (and thus more useful in modelling singular parts). The computation of singular and moving average parts is not easy in the general case and no procedure seems to known. However, a type of duality between series of inputbe output type and series of autoregressive type mav be demonstrated formally for certain non-linear processes. This is acheived by a novel expansion technique due to Jones(1978) who employed it to do moment calculations. The result suggests that large non-linearities correspond to large singular parts. Note however that the computed processes are not, strictly speaking, moving averages.

Suppose that a non-linear autoregression satisfies the equation

(2.1) $X = \lambda(X) + e_{n+1}$ n t

where λ is infinitely differentiable. For example $\lambda(x) = \frac{2}{2} \exp(-\frac{1}{2}x)$. Consider the family of processes indexed by ξ which

satisfy equations of the form

(2.2)
$$W_{n+1} = \xi \lambda(W) + e_n$$

Suppose for now that when t=0, W = w for each ξ . but that t o each process has a stationary distribution (as a Markov process). Then for each ξ by substitution (2.3) $W = e + \xi \lambda(w)$

$$W = e + \xi \lambda (e + \xi \lambda (w))$$

2 2 1 0

$$= e + \xi \lambda(e) + \xi \lambda(w) \lambda (e) + \dots$$
2
1
0
1

On expanding λ whenever necessary obtain (2.4)

 $W_{n}(\xi) = e + \xi \lambda(e) + \xi \lambda(e) \lambda (e) + higher terms$ n n n n-1 n-2 n-1

To approximate the original series set $\xi = 1$.

Since the processes have a stationary distribution, we expect that the effect of the starting value w will die out. Thus the influence of e must diminish as j grows, and the n-j of e for small j must have greater influence. Such values . n-j heuristic reasoning suggests that the expansion obtained above may be a valid decomposition of X into stationary processes which are explicit functions of the That is, е. а

representation in input-output form. While this is not a true moving average, examination of this result seems to indicate that in this type of representation is more useful when ξ (which determines the size of the singular part) is small. This suggests that a generic non-linear autoregressive time-series will have a slowly dying singular part and will not be usefully expressible in input-output form. In the linear case а timeseries with singular component (non-stationary model) а similarly has no moving average representation.

5.3 PREDICTION THEORY FOR MOVING AVERAGES

Many theoretical results exist for the moving average part of a weakly stationary process. See for example Anderson (1971) for general linear theory and Schetzen (1980) for an exposition of some general non-linear theory. Such results seem to show that the emphasis on linear models is often entirely proper. In particular, certain results in prediction theory show that the improvement in prediction attainable through non-linear predictive schemes is limited.

Specifically, when the constituent variables η of a linear

moving average are independent identically distributed (iid) random variables, classical prediction theory gives the error of the best linear forecast (eg. Deutsch, 1965) as follows. Suppose X has such a representation

 $\begin{array}{ccc} (3.1) & X = \Sigma & Y \\ & n & j & j & n-j \end{array}$

Now define

(3.2) $\Phi(s) = \begin{vmatrix} \Sigma & a e \\ j = -\infty & j \end{vmatrix}$

and set

(3.3)
$$\Delta^{2} = \exp(\frac{1}{2\pi}) \int \log \Phi(s) ds$$

If E(Y)=0 and $Var(Y)=\sigma < \infty$ then $\Delta \sigma$ is the mean square error of optimal one step prediction, that is (3.4) $\Delta \sigma = \inf E((X - \Sigma b X))^2)$ n+1 j=0 j n-j

where the inf is taken over all finite sequences b of real i

It is well known that for Gaussian time-series the linear forecast (obtained by a Box-Jenkins model) is the optimal one in mean square. This is not true for the Non-Gaussian case. Therefore some gain in forecasting power may conceivably be obtained using some non-linear function of preceding, observations. However, Kanter(1979) has proved the following result for linear moving averages which assumes only some mild conditions on the distribution of Y (which must be iid). Let Ω

denote the set of all Borel measurable functions f from R into R. Then, letting

(3.5) $x^{n} = (\dots, x_{n-1}, \dots, x_{n-1}, x_{n-1})$

we have

(3.6)
$$\inf_{\substack{f \in \Omega \\ f \in \Omega}} E((X - f(X))) \ge Q(Y) \Delta$$

is defined be variance of where the constant Q(Y) to the Gaussian random variable whose entropy equals the entropy of Y . entropy of a random variable Y is simply -E{ log f(Y) } The where f is the density of Y . It is thus just the information function defined earlier evaluated at f(.) . For Gaussian 2 sequences Q(Y) is just σ . The result means that for sequences 0 Y with Q(Y) > 0 and $\sigma' < \infty$ then non-linear prediction can improve prediction by at most a factor Q(Y)/ σ . In particular,

if perfect non-linear prediction is possible then perfect linear prediction is possible for moving average processes with finite variance and Q(Y) > 0.

Thus, the use of Box-Jenkins models to compute predictions for weakly stationary processes with iid moving average representations is sensible provided the random variables in the moving average are not very far from Gaussian. Since however, the random variables of Wold's theorem may not be iid, in principle there may still be some value in other non-linear representations for weakly stationary processes with small singular parts.

For example, Granger and Newbold(1976) considered series of the form

(3.7)
$$y = g(x)$$

where y is an instantaneous transform of x. If a linear model t t is identified for y then this implies a non-linear model for t x. Let the inverse transformation of g be denoted by h. If h t is known and is in some sense well behaved then an analytical expression for the optimal non-linear h-step forecast f x, of n,h x can be obtained. It may be shown that if the n+h

forecastability of the series is measured by

(3.8)
$$R = 1 - var(f) / var(x)$$

h,x n,h n+h

2 then R

2 2 R > R . That is, the non-Gaussian series constructed h,y h,x

in this manner is always less forecastable than the 'original', Gaussian series. Since in addition the construction of nonlinear forecasts is more difficult than the construction of linear forecasts, this observation ought to motivate the search for transformations assuming such non-Gaussian series occur in practice. Unfortunately this search is not easy. It must be born in mind that for a series to be Gaussian, all the joint distributions must be multivariate normal, an extremely stringent condition. The most well known set of candidate transformations is the Box-Cox (1964) class which are of the form

(3.9)
$$y = \{ (x + m) - 1 \} / \theta$$

and include logarithmic (θ =0) and linear transformations (θ =1) as special cases. Nelson (1976) found however, that when this transformation was applied to a number of economic time-series the resulting optimal forecasts were often no better than the linear ones, suggesting that in fact the transformations did not give Gaussian series. Thus, without a priori knowledge the search for such a transformation seems problematic.

non-linear Other proposals for non-singular models eg.Nelson&Van Ness(1973) are less well known and difficult to interpret. Because the potential improvement in prediction is limited and identifying non-linear models is difficult, attempts to formulate general non-linear models which are based upon modelling the moving average portion of a stationary time-series will not likely produce great rewards.

5.4 BILINEAR MODELS

Despite the preceding remarks it may sometimes be worthwhile to develop non-linear representations for processes with small singular parts. For one thing, even small increments in predictive power may be highly desirable. In addition, the fitting of such models may be useful as a preliminary procedure in a more detailled structural model fitting protocol . The most promising class of non-linear weakly stationary non-singular models to appear are the bilinear models of Granger& Andersen(1978) and Subba Rao (1980). While experience with these models is still lacking, it appears that these models are sucessful in formulating non-linear forecasts which acheive some of the gain in predictive power available without resort to structural information. Moreover consideration of these models may lead to useful tools for the empirical detection of nonlinearity. These models may be derived as follows. First observe that if a time-series follows the non-linear difference equation

then one can derive an AR(1) model by approximating g using the first term of a Taylor expansion. Suppose now that we iterate once before doing the expansion. It is found (again taking first terms only) that

(4.2) $X = g(g(X) + \epsilon) + \epsilon$ t+1 t-1 t-1 t

$$= F(X, \epsilon) + \epsilon$$

t-1 t-1 t

which may be approximated by the model (4.3) X = aX + bX + c + de + et+1 + t-1 + t-

both in X and ϵ . Under some circumstances it has a useful t-1 t

structural interpretation. The terms in the expansion may be thought of as representing respectively, the dependence of X on X , the interaction of X and ϵ , and the dependence of X t-1 an input disturbance ϵ . Such a model might also arise if on (say) an AR(1) process were observed with multiplicative white noise observation error. For if (4.4)and Y = X (1 + u)where us is a white noise. then (4.5) $Y = \beta Y$ ((1 + u)/(1 + u)) + ϵ (1 + u) t t-1 t t-1 t t t which is approximately $Y = \beta Y \quad (1 - u + u) + \eta$ t t-1 t t (4.6)These models may be generalized to (4.7) (note the absence of a term in X ϵ). This is the model which t t would arise on approximating X = F(X, X, ..., X; e, e, ..., e)t t-1 t-2 t-n t-1 t-2 t-n (4.8 by the first terms of the Taylor series for F, where F is

derived by iterating once in an n-1 order non-linear AR model. To compute effectively with this model is it necessary to assume that the errors e are independent random variables, in contrast t to the linear theory, in which errors need only be uncorrelated.

theory, in which errors need only be uncorrelated.

Although there is a large literature on engineering and control applications of such models in which the error series is replaced by a control variable (Bruni,Dipillo&Koch,1974). a complete theory does not exist. For example necessary and sufficient conditions for stationarity and invertibility are not known. Therefore we confine ourselves to a subclass of models (the so-called diagonal models) in the sequel, and actually discuss only one example. This class of models is known to generate useful weakly stationary processes. However, a true bilinear model may be discriminated from a Box-Jenkins model (with iid errors) by examining the covariance structure of the

squared series X. Following Granger&Anderson (1978) we t demonstrate these assertions for the particular diagonal bilinear model

(4.9) $X = \beta X \quad \epsilon \quad + \quad \epsilon$ t t-1 t-1

Define $\lambda = \beta \sigma$. Using the standard operator techniques it ϵ is found

(4.10) $X = \epsilon + \Sigma \quad (\epsilon \quad \beta B) \quad \epsilon$ t t k=1 t-1 t

99

st

Multiplying by ϵ and taking expectations gives

$$(4.11) \qquad \qquad E[X \epsilon] = \sigma \\ tt \epsilon$$

Thus the solution given makes the latter product meanstationary. Returning to the defining relation this implies

(4.12)
$$E[X] = \beta \sigma = \lambda \sigma$$

t $\epsilon \epsilon \epsilon$

To obtain the variance of X first square the defining relation t

to obtain

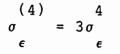
Now take expectations

To compute the expectation on the right hand side of the latter equation, use the squared defining relation, multiply by ϵ and t take expectations. From the independence of the error series we thus get

This difference equation has a positive solution for

independent of t iff $|\lambda| < 1$. Thus this model is variance stationary iff $|\lambda| < 1$. To obtain the first order autocovariance ρ procede as follows. Using independence note that $(4.17) \quad E[XX] = E[(\beta X \epsilon + \epsilon) X]$ $t t-1 \qquad t-1 \quad t-1 \qquad t \quad t-1$ $= \beta \mathbb{E} \begin{bmatrix} X & \epsilon \\ t-1 & t-1 \end{bmatrix}$ 2 Now substitute for X using the difference equation and take t-1 expectations to obtain (4.18) $E[X X] = 2 \beta E[\epsilon] E[X \epsilon]$ t t-1 t-2 t-2 2 = 2 (E[X]) t Since (4.19) $\rho = \operatorname{cov}(\mathbf{X}, \mathbf{X}) / \operatorname{var}(\mathbf{X})$ a formula for ρ may be computed.

For Gaussian errors, the formula can be made a little more explicit . In particular,



so that

(4.19)
$$E[X] = \sigma \frac{2\lambda + 1}{\epsilon 1 - \lambda^2}$$

Thus, recalling the definition of $\boldsymbol{\lambda}$

(4.20)
$$\operatorname{var} X = E[X] - (E[X])^{2}$$
$$= \sigma_{\epsilon}^{2} \frac{(1+\lambda^{2}+\lambda^{4})}{1-\lambda^{2}}$$

Since the expression for cov[X, X] reduces to $\sigma \lambda$ we can t t-1 ϵ

compute

(4.21)
$$\rho_{1} = \frac{\lambda^{2} (1-\lambda^{2})}{1+\lambda^{2}+\lambda^{4}}$$

Using this expression it is easily calculated that ρ attains a 1 maximum of .1547 at approximately $\lambda = \pm 0.605$.

Now consider

 $(4.22) E[XX] = \beta E[X \epsilon X] \qquad k>1$ tt-k t-1 t-1 t-k

On substituting for X obtain t-1

$$= 0 + \beta \sigma \sum_{\epsilon}^{2} E[X]$$
$$= (E[X])$$
$$= (E[X])$$
t

Hence

(4.24) $\operatorname{cov}(X, X) = \mathbb{E}[X X] - (\mathbb{E}[X])^{2}$ $t t^{-k} t t^{-k} t t^{-k} t^{-k}$ $= 0 \qquad k>1$

Thus the correlation structure of this model is that of an MA(1) Box-Jenkins model.

Consider now the series Z = Y where Y is generated by an t t t

MA(1) Box-Jenkins model with iid errors e then t

(4.25) E[Z Z] = t t-k

= 0 for k>1

Thus Z again has the correlation structure of an MA(1) series.

(4.26)
$$\begin{array}{c} w & 2 & w \\ \rho &= \lambda & \rho \\ \tau & \tau^{-1} \end{array}$$
 for $\tau > 1$

while ρ is a complicated function of λ . Thus the correlation 1 structure of the squared series is that of an ARMA(1,1) process. In this manner, the linear and bilinear processes may be

distinguished.

It is possible to compare in principle the performances of predictors based on the two models if the process is actually bilinear. If a linear moving average model is used and the covariance structure is known then b in the model

$$\begin{array}{cccc} X &= b \ \theta &+ \ \theta \\ t & t-1 & t \end{array}$$

where $\{\theta\}$ is white noise, may be calculated as the solution of t

(4.27)
$$\rho_1 = \frac{b}{1+b^2}$$

|b|<1 for invertibility

Then as shown earlier

(4.28)
$$\operatorname{var} \theta = \operatorname{var}(X) / (1 + b)$$

t t

Recall that variances and covariances were calculated as functions of λ and σ . Thus, the relative performance of the ϵ two models can be summarized by the ratio (4.20) $P(\lambda, \sigma, \lambda) = mpr(0, \lambda)$ (mark)

$$(4.29) P(\lambda,\sigma) = var(\theta) / var(\epsilon)$$

$$\epsilon t t$$

For $\sigma = 1$, P increases from 1.000 to 4.20 as as λ increases

from .01 to .75 and P increases from 4.20 to 27.838 as λ increases from .75 to .95 . Thus, for larger values of λ impressive gains in forecasting power are possible. We note that for $|\lambda|>1$ the model becomes non-stationary.

However, whether such gains are commonly realized in

practice depends on the frequency with which bilinear models occur in nature. As Tong&Lim(1980) remark bilinear models do seem to simulate a large fraction of the types of nonnot bilinearity is linearity possible, thus the specification of relatively restrictive. Bilinear series do not have interesting the However, the distributions of random singular parts. variables of a bilinear process differ from those of linear a For example, Granger&Anderson show the diagonal process. processes have non-zero third central moments in general , so the distributions may be skewed, and some higher moments . that need not exist, assuming the input errors are normal. These observations suggest that bilinear models are best suited to the simulation of time-series which are non-Gaussian moving averages.

Bilinear models have been fitted to the classical dataseries such as the Wolfer sunspot numbers and Canadian Lynx data as well as other sample series. It is difficult to compare the results investigators (e.q. Subba-Rao between (1979), Granger&Anderson (1978), Tong&Lim (1980)) owing to the tendency for each to employ individual fitting procedures and to fit their models to different intervals of the series concerned. addition the performance of these models is typically Ιn assessed by omitting a few (usually 10-20) points from the end of the data series and examining the ability of the model to predict these omitted points. While this method cannot be criticized in view of the shortness of these data series, it

means that no one instance is particularly conclusive. Invariably however, an improvement in power of prediction over the linear model between 5 and 25% of mean square error is reported. This suggests that bilinear modelling procedures do provide non-linear predictive power in the absence of structural data. VI. APPROPRIATENESS OF FULLY NON-LINEAR MODELS: PRACTICE 6.1 VIOLATIONS OF LINEARITY

We have indicated that small non-linearities can be modelled via a general model such as the bilinear class. When models with large non-linearities are admitted potentially greater rewards are possible. However, the richness of structure which such models admit can lead to difficulties in model selection. We have maintained that it is necessary to use structural information in this case. In the following section we discuss the application of this suggestion in practice.

It seems best to think of serious non-linearity as one does non-stationarity - as an assumption forced upon the analyst either a priori, or by the failure of the data to exhibit the pleasant features of linear model data. We will discuss examples. A non-stationary Box-Jenkins model, as discussed earlier is of the form

is the difference operator. As noted, this model form where D is designed to simulate data which is homogeneous in time except for a deterministic or stochastic polynomial trend. Α nonsimulate significant contrast, must linear model, in deterministic or stochastic inhomogeneities over space, where by most promising space we mean the values taken on by х. The classes of models for this purpose in our view are the so called

piecewise linear autoregressive models.

Non-linear models are required when it is known that a linear model leads to a failure to simulate behavior of the system under study. To clarify this criterion, note that linear structure leads to a number of properties which might loosely be called independence properties, because their statements assert that some output property is independent of some input property. Just as different types of non-stationarity justify different choices of differencing order d, violations of different independence properties will indicate that certain non-linear models are appropriate.

We now recall some independence properties of linear models. Most importantly, a linear system is amplitude independent in that if c is a constant

(1.2) C L(e) = L(c e)

That is, inputs of differing amplitudes are output without distortion. This assumption is an unrealistic one for many models since most physical, biological and economic systems exhibit some form of saturation, as well as a linear range. The simplest conceptual example (not the simplest mathematical example) is to define the operator F such that

(1.3) F(e)(s) = et s

provided |e | < c for all t < s and

F(e)(s) = 0 otherwise

This is the behavior exhibited by a fuse. Models of flooding and river overflow (Sugawara, 1962) as well as earthquakes and stock market prices are natural examples of such phenomena.

Tong&Lim (1980) propose, and fit piecewise linear models which are similar to the fuse in that they contain amplitude dependent thresholds, though they do not have infinite memory. By experimenting with parameters in such models it is possible to illustrate several types of violation of the independence properties possessed by linear systems.

non-linearity arise Other manifestations of from amplitude dependence in the frequency domain. interpreting an Such interpretation is sensible structurally because many natural phenomena are known to be periodic (often because they satisfy linear differential equations) and thus it can happen that the input to a system is a superposition of a finite number frequency components. When response depends on input of amplitude, clear empirical relations may appear in the data particular between frequency and amplitude. In the data which large vibrations of hiqh (low) contains epochs in frequency are interspersed with small vibrations of low (high) frequency. This is observed for example, in records of EEG during sleep. Tong&Lim give the example

1.4)
$$X = 1.6734 - 0.8295X + 0.1309X - n^{-1} n^{-2}$$

(1)
 $0.0276X + \epsilon n^{-3} n^{-3} n^{-3}$
if $X > 0.5 n^{-1}$
(2)

$$= 1.2270 + 1.0516X - 0.5901X - 0.2149X + \epsilon$$

n-1 n-2 n-3 n

otherwise

where $var(\epsilon) = 0.003$, i=1,2 n

which simulates the first mentioned behavior, and another similarly complex example which simulates the opposite case. An important question here seems to be determining the simplest process which gives such behavior.

The superposition principle is an obvious consequence of linearity, which also has consequences in the frequency domain. It may be seen that if a sinusoidal input is applied to a linear system defined by

then setting

the output is given by

(1.6)
$$\begin{array}{ccc} & -\infty & \operatorname{im}(n-1) \\ X &= \Sigma & a & e \\ n & 1 = 0 & m \end{array}$$

Thus, the output is a sum of attenuated (amplified) and delayed (advanced) sinusoids of the same frequency. This is true of a general linear system, and unsurprisingly may fail for nonlinear systems. For example the output of a non-linear system may contain harmonics of the input frequency. Tong&Lim give the example

$$\begin{array}{ccc} (1.7) & -(2 + \sqrt{2}) & y - (1 + \sqrt{2}) & \text{if } -1 < y & \leq -1/\sqrt{2} \\ n & n & n \end{array}$$

$$-\sqrt{2} \ Y - 1 \quad \text{if } -1/\sqrt{2} < Y \leq 0$$

n n

X =

$$\sqrt{2}$$
 Y -1 if 0 < Y $\leq 1/\sqrt{2}$
n n

$$(2 + \sqrt{2})Y - (1 + \sqrt{2})$$
 if $1/\sqrt{2} < Y \le 1$

which responds to some input sinusoids Y with a waveform having ,

frequency double that of the input. This phenomenon is difficult to detect by inspection of the data though it must be born in mind when a system exhibits systematic periodicity of a higher (lower) frequency than can be justified by a linear theory.

From the point of view of time-series analysis the most restrictive feature of linear models is that a stationary linear

system converges to zero in the absence of input, independent of the initial condition. Recall that to see this one puts the system in autoregressive form and a difference equation is solved for the output. It was seen earlier that it is possible to generalize the linear model to simulate processes containing sinusoidal components but that the other types of nonstationarity encompassed by this model lead to processes with unbounded variance and are thus of limited interest.

In contrast non-linear recursions may have deterministic solutions such that for various initial conditions { X ϵ I } o there are distinct limit sets A(I) with the property (1.8) x ϵ A iff for all ϵ >0 |x-X | < ϵ n

infinitely often

will call such a set an attractor. In particular if the We state space is partitioned into a finite number of initial sets I such that for each i>1 A(I) is a single point then we say i the system possesses multiple equilibria. If for some initial, A(I) contains a closed trajectory of the system we say set I the system admits a cycle. If there exists an I containing an open neighbourhood of A(I) and A(I) contains only finite closed trajectories then the system admits stable limit cycles. Finally, trajectories of non-linear recursions need not be closed even though they are bounded. In this case it is said the system admits chaotic solutions. (Li&Yorke,1975). When limit cycles or chaotic solutions exist deterministic non-linear

phenomena may account for most of the variance of the series, although it may present a random appearance to some conventional analyses (May, 1974; Bunow&Weiss, 1979).

The paper of Bunow &Weiss illustrates the complexities which ensue when simple deterministic recursions are analysed as linear time-series via their autocovariance functions. These examples are non-linear time-series which are entirely singular. Several recursions were studied including the discrete logistic recursion (LI&Yorke, 1975)

(1.9)
$$X = r X (1 - X)$$

n+1 n n

and the triangle recurrence (Guckenheimer,Oster&Ipaktchi(1977) (1.10) X = 1 - 2 | X - .5 |n+1 n

by systematically varying parameters and initial values.

For the logistic recursion, the most common case (eg. r=3.64, X=.878) is an almost uniform comblike autocorrelation 0 function of high amplitude extending to very large lags, in excess of 100 points. Less common (eg. r=3.75, X=.374) was a 0 rapidly decaying oscillatory autocorrelation function with small spindles at large lags. Sometimes (eg. r=4.00, X) the 0 autocorrelation function was indistinguishable from that of a sequence of independent random variables, or at the other extreme, was periodic. The triangle recursion graph is similar to that of the logistic recursion with r=4.00 so it is perhaps not surprizing that its autocorrelation function was usually identically zero.

Limit cycles, chaos, and other non-linear phenomena characterize the rich variety of singular parts in non-linear time-series models. This fact is practically important because of examples in which a priori models which predict non-linear phenomena and yet investigators have persistently fitted linear models to data. This is true particularly of the ecological literature (cf. Campbell&Walker,1977; Tong,1977) as well as certain phenomena in economics (cf. Granger&Anderson,1978).

6.2 THE THRESHOLD MODELS OF TONG&LIM

Tong&Lim(1980) proposed a variety of non-linear autoregressive models, with the property that different linear difference equations were satisfied for different values of X.

As noted, such models can simulate many of the non-linear phenomena described above. However, because this class of models is considerably richer than the linear class there is some difficulty in conducting the identification procedure. The fitting procedure for one type of model is described as follows.

Let P be a partition of the real line corresponding to the 1+1 points {r ,r ,...r). arranged in increasing order. The r o 1 1 i will be referred to as thresholds. Define R = (r ,r]. j i-1 i Then, a self-exciting threshold autoregressive model of order (1;k,...,k) or SETAR(1;k,...,k) (k is repeated 1 times and gives the order of the linear autoregression for each of 1

regimes) has the form

(2.1)
$$X = a + \Sigma a X + \epsilon$$
,
n o i=1 i n-1 n

where we write

$$\underline{x} = (x, x, \dots, x)^{T}$$

and J is a random variable such that J = j if (k) $\underline{X} \in R$ for j=1,1. It is assumed that the errors are n-1 j Gaussian (the series is not) independent in different regimes, and that R is a cylinder set of R, of the form j $R \times R \times ... R \ldots \times R \times R$ depending on some fixed lag d. Because

these models are autoregressions they are always invertible.

To fit such a model, it is necessary to estimate the thresholds, the critical lag d and the respective orders of the linear autoregressive regimes. Tong&Lim describe the procedure followed in fitting a SETAR(2,k,k). To begin, Q candidate, 1 2 thresholds t and D critical lags are selected. In this case

q each threshold distinguishes two regimes. Method of model comparison used is Akaike's criterion which in this case is proportional to

(2.2) AIC(k) = N ln(RSS / N) + 2k

where N is the number of observations, k the number of fitted parameters and RSS denotes the residual sum of squares. On each

data set determined by each fixed t and d the autoregressive qorders are determined to minimize AIC(k) for $0 \le k \le K$ where K is some fixed maximum order. Because the errors on each regime are independent one may next write

(2.3)
$$AIC(t) = AIC(\hat{k}) + AIC(\hat{k})$$

This criterion is minimized in turn to estimate t . The latter \mathbf{q}

procedures establish a minimum AIC model for each value of d.

A difficulty arises in comparing models with different values of d because of the asymptotic nature of Akaike's statistic. In particular, the selection of a different lag results in a different effective number of observations. Tong&Lim apply a normalization and thus compute

(2.4) AIC(d) = AIC(
$$\hat{r}$$
) / (n -max(d, K))
i

This procedure produces semi-automatic estimates for all the necessary parameters. As in the linear case residuals may be plotted and analysed and the model possibly rejected.

It is difficult to evaluate the sucess of Tong&Lim's procedures at this relatively early stage. In fitting classical data series such as the Canadian Lynx data and the Wolfer Sun spot numbers they report 'disappointing' results. Nonetheless, for the fitting periods chosen by Tong&Lim, modest improvements in prediction power over the linear (and bilinear) model were reported. For example for the Wolfer sunspot numbers from 1770-1869 the mean square error of one step prediction over the next

20 years was 346.6 for an AR(2) model(Box&Jenkins), (2.5)X = 14.70 + 1.425X0.731X +-2 t-1 t 293.4 for the linear model in combination with a bilinear model fitted by Granger&Anderson(1978) to the residuals ϵ $= -.0222\epsilon \eta t-2 t-1$ (2.6)+ 0.202¢ + - 1t and 267.6 for the SETAR(2;3,4) model fitted by Tong&Lim (2.7) X = 5.269 + 1.889 X - 1.5289 X+ 0.3039X n-2 n-3 n-1 0.3387X + n-4 if X <36.6 n-3 2 0.3900+1.1366X -0.3645X + 0.0524Xn-1 n-2 n-3 n if X >36.6 n-3

This seems an impressive gain in predictive power, however, there appears to be some non-stationarity in this data which is difficult to model. One step predictions over a longer period deteriorate. The predictions of a SETAR(2;4,2) fitted to the 1837 to 1920 data deteriorate after 1944.

For the (logarithmically transformed) Lynx data (1821-1920) a linear AR(12) model, selected using Akaike's criterion Tong(1977) had a mean square error of prediction of .018, while a SETAR(2; 8,3) had mean square error .0144 and a bilinear model fitted by Subba-Rao (AR order 11 with cross terms up to X ϵ) had the same mean square error as the linear model t-8 t-10

over 14 subsequent years. The most striking feature of these results is that the mean-square errors are so similar. Certainly no clear margin of superiority is evident for piecewise linear over bilinear models, or even over linear ones. It seems consistent with the evidence to speculate that there is а practical limit on the expected improvement over a linear model without resort to structural data, and both bilinear and piecewise linear models are just rich enough to be close to the optimum in these cases.

In contrast a fitted model for flow of the Kanna river in Japan as a function of rainfall and past flow gave a reduction of 18% in mean square error against an (unspecified) competing linear model over 86 time points. This is their most clearly sucessful model in this respect. It is a complicated model. There is one (rainfall) threshold distinguishing two regimes. One involves log riverflow X to lag 5 and rainfall Y to lag 4.

The other regime involves each variable to lag 2. While the model may be criticized for its complexity, before presenting the fitted model Tong&Lim note that seasonal variations of Japanese rivers are quite regular due to the rather well defined rainy season there. In addition, the ground soil is rarely dry so that variations in the water table are small. Since rainfall is not dissipated instantly, it is reasonable that present riverflow is a function of past river flow. They conclude that it is reasonable to model river flow data as a function of rain fall which may contain thresholds. This model thus has some structural justification, unlike those for the classical series. In view of what we have said so far the success of this model is predictable. Akaike notes in the discussion that in his (implementing control strategies for industrial experience processes) success had been acheived by characterizing the system physically in order to chose the particular nonlinearitities to be modelled. He does not suggest that а coherent theory need exist, only that the appropriate conditioning variables should be chosen from physical considerations, or simple examination of the data.

In evaluating Tong&Lim's approach it is useful to focus on Chatfield's guestions. 'How can we tell if a given time-series is non-linear?' and 'How can we decide if it is worth trying to fit a non-linear model?'. In our view, a serious attempt to these questions even in particular analyses is a large answer step towards a successful model. Ιf so there would be much future research to determine the simplest threshold value in models which produce particular violations of linearity, and to catalogue examples where threshold models are scientifically justified.

6.3 SUMMARY

The problems inherent in identifying time-series models were discussed and it was shown that structural information can be judiciously utilized when a parametrization is to be selected. Recent proposals for non-linear time-series models were examined in this light. Weakly stationary time-series were seen to consist of deterministic and moving average components. The potential for non-linear modelling of the moving average shown to be limited, but it was suggested that component was worthwhile improvements could result from modelling this component with a weakly non-linear (bilinear) data prescription. Conversely, we suggest that modelling of the deterministic component might best be accomplished via the strongly non-linear (piecewise linear) data prescriptions, but that for this purpose some structural information may be required.

APPENDIX A - SOLVING A HOMOGENEOUS DIFFERENCE EQUATION

Let Z(t) be a real function of an integer time index satisfying

(A.1) $P(B)Z(t) = Z(t) - p Z(t-1) - p Z(t-2) - \dots - p Z(t-n) = 0$

We will discuss the method of solution of such a recursion. It turns out that the solutions are completely analogous to the solutions of homogeneous linear differential equations . Thus the solutions can be obtained by first substituting for Z(t) functions of the form

(A.2) (i) A r or (ii) A t r

Next, determine r from the resulting algebraic equation, and count linearly independent solutions to see that the general solution has been obtained. Another way to express this procedure is as follows. Note that because P(B) is formally a polynomial in B, it can be factored using the fundamental theorem of algebra as

Because IB=BI, each of the bracketed factors commute, any factor we choose may be written last. Also note that any operator polynomial T(B) operating on the zero sequence gives zero by linearity. These facts together reduce the solution of the difference equation to the solution of

(A.4)
$$(I - rB)^{m}Z = 0$$

Thus for each root of multiplicity m there corresponds an arrangement of the factors so that the operator in the latter equation appears last. If this equation is solved, we have a solution of the original equation, for the effect of the last m factors is to annihilate Z(t). To solve this equation use an inductive approach. First note that for n=1 the obvious t solution is Ar

We claim that the equation (A.4) has m independent solutions of the form

Since the result is true for m=1 it suffices to show

(A.6)
$$(I - rB) \{ tr \} = K t r + K t r + ... + K j-1 j-2 0$$

but

So that the degree of the putative solution is reduced by 1 by the operator I-rB. Thus, the result follows because applying each of the m factors one after another must eventually annihilate the original form (A.5).

BIBLIOGRAPHY

- AKAIKE,H. (1969) Fitting autoregressions for prediction. Annals of the Institute of Statistical Mathematics, 21 ,No.2,243-247.
- AKAIKE,H. (1973) Information theory and an extension of the maximum likelihood principle. In <u>Second International</u> <u>Symposium on Information Theory</u> ,(B.N.Petrov&F.Csaki,ed),pp.267-281,Budapest Akademai Kiado.
- ANDERSON, T.W. (1971) The Statistical Analysis of Time-Series , New York, Wiley.
- BORA-SENTA, E. & KOUNIAS, S. (1979) Parametric estimation and order determination of autoregressive models. In <u>Analysing Time-</u> <u>Series</u>, (O.D. Anderson, ed.), pp 93-103, Amsterdam, North Holland.
- BOX,G.E.P.& COX,D.R. (1964) An analysis of transformations. Journal of the Royal Statistical Society, Series B , 26 , 211-252.
- BOX,G.E.P.&JENKINS,G.M. (1970) <u>Time-Series Analysis</u>, Forecasting and Control ,San Francisco, Holden Day.
- BRUNI,C., DIPILLO,G.& KOCH,G.(1974) Bilinear systems: an appealing class of 'nearly linear' systems in theory and applications. IEEE Transactions in Automatic Control , TAC-74 , 334-348.
- BUNOW, B. &WEISS, G.H. (1979) How chaotic is chaos? Chaotic and other 'noisy' dynamics in the frequency domain. Mathematical Biosciences, 47,221-257.
- CAMPBELL, M.J. & WALKER, A.M. (1977) A survey of statistical work on the Mackenzie river series of annual Canadian lynx trappings for the years 1821-1934 and a new analysis. Journal of the Royal Statistical Society, Series A, 140,411-431.
- DEUTSCH, R. (1965) Estimation Theory, Englewood Cliffs, N.J. Prentice Hall.

- DURBIN, J. (1960) The fitting of time-series models. Review of the International Institute of Statistics, 28,233-289.
- FRIEDMAN, J.H. & STUETZLE W. (1981) Projection Pursuit Regression. Journal of the American Statistical Association , 76,817-823.
- GERSH, W. & SHARPE, D.R. Estimation of power spectra in finite order autoregressive models. <u>IEEE Transactions in Automatic</u> <u>Control</u>, <u>AC-18</u>, 367-379.
- GRANGER, C.W.J. & ANDERSEN, A.P. (1978) <u>Introduction to</u> <u>Bilinear Time-Series Models</u>. Gottingen, Vandenhoeck & Ruprecht.
- GRANGER,C.W.J. & MORRIS, M. Time-series modelling and interpretation . Journal of the Royal Statistical Society, Series A , 139 , 246-257.
- GRANGER, C.W.J. & NEWBOLD, P. Forecasting transformed series. Journal of the Royal Statistical Society, Series B , 38 , 189-203.
- GUCKENHEIER, J., OSTER, G. & IPAKTCHI, A. (1977) Dynamics of density dependent population models, <u>Theoretical Population</u> <u>Biology</u>, <u>4</u>, 401-415.
- HANNAN, E.J. & QUINN, B.G. (1979) The determination of the order of an autoregression . Journal of the Royal Statistical Society, Series B. 49, 190-195.
- HOSKING, J.R.M. (1978) A Unified derivation of the asymptotic distributions of goodness of fit statistics for autoregressive time-series models. <u>Journal of the Royal Statistical Society</u>, <u>Series B</u>, <u>40</u>, 341-349.
- JONES, D.A. (1978) Non-linear autoregressive processes. Proceedings of the Royal Society of London, A , 360 ,71-95.
- JONES, R.H. (1975) Fitting autoregressions. Journal of the American Statistical Association , 70 , 590-592.

- JONES, R.H. (1976) Autoregression order selection. <u>Geophysics</u>, <u>41</u>, 771-773.
- JONES, R.H. (1981) Fitting a continuous time autoregression to discrete data. In <u>Applied Time-Series Analysis II</u>, (D.F. Findley, ed.) Academic 1981, pp 651-674.
- KANTER, M. (1979) Lower bounds for non-linear prediction error in moving average processes. <u>Annals of Probability</u>, <u>7</u>,128-138.
- KOOPMANS, L.H. (1974) The Spectral Analysis of Time-Series . New York, Academic Press
- LI, T.Y. & YORKE, J.A. (1975) Period three implies chaos. American Mathematical Monthly, 82,988-992.
- LUDWIG, D. (1975) Persistence of dynamical systems under random perturbations. SIAM Review, 17,605-640.
- MAY, R.M. (1974) Biological populations with non-overlapping generations : stable points, stable cycles & chaos. <u>Science</u>, <u>186</u>,645-647.
- NELSON, H. (1976) The Use of Box-Cox Transformations in <u>Economic Time-Series, An Empirical Study</u>, Unpublished Ph.D. thesis, Economics Department, University of Southern California at San Diego.
- NELSON, J.Z. & VAN NESS, J.W. (1973) Formulation of a nonlinear predictor <u>Technometrics</u>, <u>15</u>, 1-12.
- OZAKI, T. (1977) On the order determination of autoregressivemoving average models. Applied Statistics, 26, 290-301.
- PARZEN, E. (1974) Some recent advances in time-series modelling. <u>IEEE Transactions in Automatic Control</u>, <u>AC-19</u> ,713-715.
- QUENOUILLE, M.H. (1947) A large sample test for goodness of fit of autoregressive schemes. Journal of the Royal Statistical

Society, Series A , 110 ,123-129.

- RAMSEY, F.L. (1974) Characterization of the partial autocorrelation function. <u>The Annals of Statistics</u>, <u>2</u>,1296-1301.
- SCHETZEN, M. (1980) The Volterra and Wiener Theories of Nonlinear Systems, New York, Wiley.
- SCHWARTZ, G. (1978) Estimating the dimension of a model. Annals of Statistics , <u>6</u>,461-464.
- SHIBATA, R. (1976) Selection of the order of an autoregression by Akaike's information criterion. <u>Biometrika</u>, <u>63</u>,117-126.
- STONE, M. (1977) An asymptotic equivalence of choice of model by cross-validation and Akaike's criterion. <u>Journal of the</u> Royal Statistical Society, Series B, <u>39</u>,44-47.
- STONE, C.J. (1977) Consistent non-parametric regression. Annals of Statistics , 5 ,595-645.
- SUBBA-RAO, T. (1980) Discussion on paper by Tong&Lim (1980). Journal of the Royal Statistical Society, <u>42</u>,278-280.
- SUGAWARA, M. (1962) On the analysis of run-off structure of several Japanese rivers. Japanese Journal of Geophysics, 2 ,1-76.
- TONG, H. (1977) Some comments on the Canadian lynx data. Journal of the Royal Statistical Society, Series A , <u>140</u>,432-436, 448-468.
- TONG, H. & LIM, K.S. (1980) Threshold autoregression, limit cycles and cyclical data. <u>Journal of the Royal Statistical</u> <u>Society</u>, <u>Series B</u>, <u>42</u>,245-292.
- TWEEDIE, R.L. (1974) R-theory for Markov chains on a general state space I: solidarity properties and R-recurrent chains, Annals of Probability, 2, 840-864

TWEEDIE, R.L. (1975) Sufficient conditions for ergodicity and recurrence of Markov chains on a general state space. <u>Stochastic Processes and their Applications</u>, <u>3</u>,383-403.

WOLD, H. (1953) <u>A Study in the Analysis of Stationary Time-</u> Series Second edition. , Stockholm, Arquist and Wiksell.

WHITTLE, P. (1952) Tests of fit in time-series. <u>Biometrika</u>, <u>39</u>, 309-318.

YULE, G.U. (1927) On the method of investigating periodicities in distributions of series with special reference to Wolfer's sunspot numbers. <u>Philosophical Transactions of the Royal</u> <u>Society of London, Series A</u>, <u>226</u>, 267-298.