REPRINTED FROM THE BOOK:

# STATISTICAL ANALYSIS OF RAINFALL AND RUNOFF

A PART OF THE
Proceedings of the International Symposium on Rainfall-Runoff Modeling held May 18-21, 1981 at Mississippi State University, Mississippi State, Mississippi, U.S.A.

**WRP**

# STOCHASTIC MODELS OF RAINFALL-RUNOFF RELATIONSHIP

**Vit Klemes**
Research Hydrologist
National Hydrology Research Institute
Department of the Environment
Ottawa, Ontario, K1A 0E7, Canada

## ABSTRACT

Problems of operational and physically based stochastic modelling of rainfall-runoff relations are discussed, dangers of the reliance on stochastic methods of model building are illustrated, and the need for developing stochastic rainfall-runoff models on the basis of sound hydrological concepts is emphasized.

## INTRODUCTION

The title of this symposium suggests that the subject matter to be discussed is not hydrology. For if it were we would not be talking about rainfall-runoff modelling but about the modelling of the land phase of the hydrologic cycle. The domain of "rainfall-runoff relations" (hereafter usually RRR) originated of course from a problem caused by a hydrological event, namely floods. The problem as such, however, was not hydrological – it was an engineering problem or, to put it in our modern jargon, a decision problem: how to build a bridge across a river so that it would not be washed away or, at any rate, seriously damaged, by a flood. The following quotation from Sokolovskii (1971, pp. 322, 323) sets the scene:

> "In the following we shall review computation formulas for maximum discharge developed in Czarist Russia...
> The first computation formula developed at the Ministry of Railroad Transportation in Czarist Russia was introduced in connection with directives to determine storm runoff norms, and was developed by the Austrian engineer Koestlin in 1868. As a result of the incorrect design of a bridge and conduits for the passage of storm water, this structure collapsed and resulted in a train crash on the Moscow-Kura railroad near Kukuevo station. After this catastrophe the following formula was put into use in 1882:
>
> $$Q = k_d a \alpha F,$$
>
> where a is the theoretical storm intensity,... F is the drainage-basin area; $k_d$ is a dimensional factor, equal for measurement in the metric system to 16.67...; $\alpha$ is a factor allowing for losses of storm water by seepage and for the non-uniformity in their time lag to the outlet...

A subsequent train accident in 1900 as a result of a storm, on the Kharkov-Balashov railroad, made it necessary to undertake special studies on the usefulness of the Koestlin formulas. Chairman Nikolai of a special investigation commission received expert opinion from the railroad administration and concluded that for such a vast territory as Russia the Koestlin norms were in general unsatisfactory, since 40 supported the norms in their present state while 60 were in favor of increasing them..."

Similar histories are behind all "flood formulae", (sometimes they are alluded to even in their names, as for instance in the cases of the Bavarian Railways formula or the U.S. Bureau of Public Roads formula) and, by implication, behind the whole domain of RRR. To put it briefly, the decision maker needed a "maximum", "mean", "minimum", discharge or water level and ordered his engineers to get it, and get it fast, because ... some Governor was ordered by some High Commissar to build a bridge at some Kukuevo within two years..., etc., where one can freely substitute, say, Glen Canyon Dam for Kukuevo bridge, Senator for High Commissar, promised to his constituency for was ordered to, etc., etc., as the case may be.

And, since every schoolboy knows that floods come from rains and droughts from no rains, it is quite obvious where to look for quick answers. Given the additional fact that hydrology still is basically in the hands of engineers (as a rule having its academic home in Civil Engineering departments of universities and obtaining research funds from water resource management sources), it is quite understandable that, instead of improving hydrologic knowledge and understanding, we have been devoting most of our efforts to the thankless task of improving operational models for RRR.

For, from the hydrological point of view, these relations are of necessity very loose. As D.W. Mead put it more than 60 years ago:

"It is evident... that there exists no simple relation between rainfall and runoff from which either monthly or annual stream discharges can be calculated with any great degree of accuracy... Runoff should be regarded as the overflow or residual remaining after various other demands are supplied and not as a proportion of rainfall" (Mead, 1950, p. 568).

All one can add is that the same applies to the computation of maximum and minimum flows, in the latter case with the salutary exception that there is no runoff where there is no precipitation, which makes the modelling of RRR a rather promising and exact science in the conditions of Central Australia and Saudi Arabia.

The mainstream of the current practice in the development of stochastic models of RRR makes me feel uncomfortable by evoking a familiar parallel. Consider a trial. The prosecution (Water Management) charges the accused relatives, Rainfall and Runoff, that a close relationship exists between them. The prosecution has a vested interest in the conviction but not enough evidence to prove the case. It thus provides incentives (Research Grants, Prestige) to its investigator (Modeller) to obtain a confession from the accused. The investigator lets them "speak for themselves" but they have little to say beyond admitting that they may be relatives of some kind but know no details. The investigator records the result on a standard form (Regression) and reports back. The prosecutor is not satisfied and encourages the investigator to try harder. The investigator has two options. He can

recommend that external evidence be gathered because he satisfied himself that the accused said all they knew and nothing new could be learned from them. This does not happen too often because he knows it is exactly what the prosecutor does not want to hear and may fire him. Thus he usually takes the second option, goes back, extorts the confession by torture (High Calibre Mathematics) and/or compiles it on the basis of his own imagination and the prosecutor's wishes (Arbitary Assumptions). The case is successfully tried on the basis of such fabricated evidence and the facade of scientific respectability of stochastic RRR modelling gets a fresh coat of a brighter paint – another impressive-looking operational model of no hydrologic consequence.

## LIMITATIONS OF OPERATIONAL AND CAUSAL MODELLING OF RRR

An operational model such as a regression relation, a polynomial fit, a transfer function model, etc. is perfectly in order if it formalizes an empirically established relationship or pattern (e.g. Clarke, 1980) and the modeller or user does not transgress its inherent limitations. In particular, it should be clear that

(1) An operational model is nothing more than "geometrical" interpolation formula and will give best results in the range where the amount of empirical evidence is largest;

(2) There is no reason inherent in the mathematical structure of the model why it should be valid outside the range of data to which it has been fitted.

(3) Improvement of an operational model requires either more empirical data or more causal input, i.e. more information on the phenomenon being modelled. It cannot be achieved merely by more mathematics. The degree of mathematical sophistication has no point beyond the "carrying capacity" of the amount of the physical information available.

(4) Even if the empirical relationship being modelled by an operational model is very close, it does not necessarily imply that it is causal in the sense that a change of one variable would necessarily lead to a change in the other variable as "predicted" by the model.

As a consequence, it is obvious that operational models of RRR cannot provide reliable information where it is most needed, i.e. beyond the range of observed runoff and for other conditions than those to which the data correspond. Their development has reached a dead end because the only way to progress, other than turning to physically based modelling of the hydrologic cycle, is to wait for more empirical data. And in the context of stochastic modelling, "more" means "more by an order of magnitude" which unpleasant fact "cannot be overcome by any mathematical sleight-of-hand" (Moran, 1957). The presently fashionable polishing of operational stochastic RRR models can hardly be considered as more than a mathematical pastime to fill the intervening 500 years or so before the data base catches up.

Compared to operational modelling of RRR, the physically based stochastic modelling of the hydrologic cycle represents a jump over perhaps several orders of magnitude in terms of difficulty. An indication of this difficulty can be obtained from Eegleson's study (1978) which is the only work to date that can be properly classified as an attempt at physically based stochastic modelling of the hydrologic cycle. The difficulty has several aspects. One is the inclusion of all

the important mechanisms involved in the transformation of precipitation into runoff. So far only the simplest hydromechanical mechanisms have been considered while the thermodynamical ones which dominate the mass balance (70 to 100% of the precipitation does not reach the runoff stage) have been drastically simplified as well as the electro-chemical and other mechanisms in the soil-vegetation interphase which control the distribution of precipitation into the various phases of surface and subsurface runoff. Another aspect is an adequate description of the soil matrix and of the boundary conditions of the water-carrying medium. Yet another aspect is the modelling of the inputs into the system, namely of the processes of precipitation and the various energy fluxes. These can be represented only by operational models since they are outside of the necessary "free body cut" defining the hydrologic model (Klemeš, 1981).

This last aspect seems to make a physically-based approach to stochastic modelling of RRR self-defeating since one may say that there is no point in replacing an operational model of RRR with a physically based model if the latter is to be fed with operationally modelled inputs. There are at least three arguments to the contrary. The first is that these inputs should have a simpler stochastic structure than that of the runoff because they are largely unaffected by the complex processes operating in the basin. Since, at the same time, they are usually documented by longer records than runoff, their operational models should be more reliable than those of the runoff itself. The second is that, eventually, operational models of these inputs are likely to be replaced by causal models as a result of physically based climatic modelling. The third is that even if the two former assumptions did not materialize, a physically based hydrologic model would at least make it possible to assess the consequences of various more-or-less plausible input "scenarios" and consider them in water resource planning and design.

In summary, the situation in stochastic modelling of RRR can be described as follows:

> There are, generally speaking, two circumstances in which it is now difficult to develop a reliable stochastic model of the rainfall-runoff relationship: the first is when we treat it as empirical; the second is when we treat it as causal*.

However, the second alternative has the distinct advantage that it involves more than merely improving the pavement on a dead-end street.

## THE STARTING POINTS

For the purpose of physically based modelling of the land phase of the hydrologic cycle the system to be modelled can be divided into the two following broad categories: (1) Inputs and boundary conditions represented by operational models and (2) hydrodynamical, thermodynamical, electro-chemical, and other processes within the basin to be described by the respective physical laws.

In the first category, the interest has concentrated on the modelling of the precipitation process. The evidence indicates that, with the exception of the smallest basins, the classical assumption of

---

\*) A paraphrase of B.L. Hendrics as quoted in Berlinski (1976, p. 112): "There are generally speaking, two circumstances in which it is difficult to analyze mathematically a social system: the first is when the system is not linear; the second is when it is."

uniform basin rainfall modelled as a one-dimensional process in discrete time (hourly, daily, not to say monthly, intervals) is too crude and invalidates the result of even the most adequate basin response model. Attention has therefore turned to event-based space-time models (e.g. Sorman, 1975; Duckstein et al., 1979; Gupta and Waymire, 1979; Waymire and Gupta, 1980) which bring with them the problem of the sampling uncertainty introduced by point measurements of the rainfall field (e.g. Wilson et al., 1979; Bras, 1979). Stochastic modelling of the energy inputs does not seem to have been attempted in the context of hydrologic models, probably because the amount of noise in their components (temperature, radiation, albedo) is thought to be much lower than that in the precipitation and can, in the first approximation, be neglected.

In the second category, attention has been directed almost exclusively to the effect of various hydrodynamic subsystems, such as storage reservoirs and their systems, channels, hillslopes, etc., on the stochastic properties of their inputs. A survey of these efforts is given in Klemeš (1978) and some newer results appear, and are referenced in, for example, Singh and Birsoy (1977), Sagar (1979), Anis et al. (1979), Smith and Freeze (1979a, b) Freeze (1980), Glynn (1981).

A typical objective of such investigations is to seek the relationship between some statistical parameters of the (precipitation) input and the corresponding parameters of the (runoff) output which is in concert with the central objective of stochastic RRR modelling: the determination of the unknown stochastic properties of runoff from the known properties of rainfall.

As an example, one may consider the effect of a storage reservoir on statistical parameters of an input process $X$ routed through it. In the simplest case the routing mechanism can be represented by the momentum equation reduced to the familiar form:

$$Y = aS^b \tag{1}$$

where $Y$ is the output rate, $S$ is the instantaneous water storage and a and b are positive constants. When the problem is cast in discrete time so that $X$ is a time series $X_1, X_2, \ldots, X_i$ (which is often the case in stochastic modelling), then $Y$ and $S$ in (1) have a subscript i and both are in volume dimensions.

Figure 1 shows the influence of the constants a and b on the coefficient of variation $C_v$ and the coefficient of skewness $C_s$ of output $Y$ for reservoir input $X$ represented by a random series with a two-parameter lognormal distribution with mean $E[X] \approx 1$ and variance $VAR[X] = 1$ so that $C_v[X] \approx 1$ and $C_s[X] = 4$ (this follows from the well known relation $C_s = 3C_v + C_v^3$ valid for the lognormal distribution). For the linear reservoir (b=1), the relationship between the two input and output parameters and the constant a is (Klemeš, 1978).

$$C_s[Y] = \frac{(a + 2)^2}{a^2 + 3a + 3} \frac{C_s[x]}{C_v[x]} C_v[Y] \tag{2}$$

In Fig. 1 this relationship is represented by the solid curve drawn between the origin $(a \to 0)$ and the input $C_v$ and $C_s$ coordinates $(a \to \infty)$. The four dashed curves represent similar relationships for the

nonlinear reservoir with b = 1/5, 1/2, 2, and 5, respectively, as obtained by Monte Carlo simulation (a closed—form mathematical relation analogical to eq. 2 has so far not been developed; the attendant mathematical difficulty will be explained later). Fig. 1 shows that, for a nonlinear reservoir, the relationship between the two output parameters is quite complex. It is interesting to note that while the output $C_V$ is always lower than its input value, the output $C_S$ can be both lower and higher than its input value depending on a and b, and can become negative even for a highly positively skewed input (Klemeš, 1970).
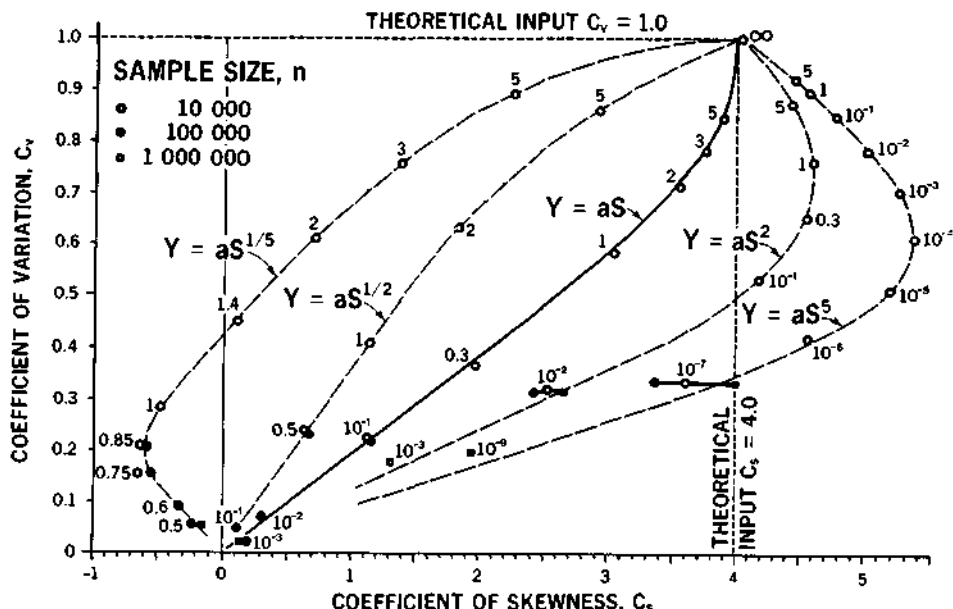


Fig. 1    Relationship between coefficient of variation $C_V$ and coefficient of skewness $C_S$ of output from nonlinear reservoir defined by $Y = aS^b$ (see eq. 1) for b = 1/5, 1/2, 1, 2 and 5, as obtained by Monte Carlo simulation for various values of a (shown by numbers) using the same input (random series with lognormal distribution with mean = 1, $C_V = 1$ and $C_S = 4$). Solid line shows theoretical solution for linear reservoir after Klemeš (1978).

The investigation of effects of various simple systems on stochastic inputs, while being a necessary stage in the research program, cannot be regarded as more than a "warm—up". The main event can start only after these simple systems have been uniquely identified with their prototypes in the basin. This will not be simple and almost certainly not possible via classical stochastic analysis as will be demonstrated in the next section. The identification will require both theoretical analysis of reducibility and separability of mathematical formulations of the physical mechanisms involved and empirical verification of results in the field (experimental research basins). For example, a specific river may be representable by a linear channel at low flows, at higher flows its behaviour may be well approximated by a cascade of linear reservoirs, while at extremely high flows one or two nonlinear reservoirs may provide the most adequate simple representation.

It is unfortunate that the most common present approach to model

conceptualization is to postulate the conceptual elements rather arbitrarily on the basis of a high-school image of the hydrologic cycle and then calibrate their parameters by minimizing the difference between the modelled and recorded outputs. Needless to say, such a conceptual model is essentially a convoluted operational model – an inflated interpolation formula – disguised in a physically sounding jargon.

An indispensable component of any hydrologic model worthy of that name must be a submodel of the processes operating on energy inputs and its coupling with the hydrodynamic submodel. This aspect has so far been almost totally neglected despite the fact that energy inputs govern the long-term water balance of the basin, as well as the initial moisture conditions and thus basin responses to individual precipitation events.

## RESULTS OF SOME SIMPLE EXPLORATIONS OF NONLINEAR STOCHASTIC SYSTEMS

The prevailing concept of stochastic hydrologic modelling has been its identification with the fitting of operational stochastic models to historic hydrologic series. This concept has been considerably reinforced during the past five years or so by a flood of papers on ARIMA modelling inspired by the undoubtedly important book by Box and Jenkins (1970). However, the main emphasis of this book is on short-term forecasting and control, i.e. on stochastic interpolation. This aspect has not been much emphasized in the publicity campaign in hydrologic literature. On the contrary, more emphasis has so far been given to applications of Box-Jenkins techniques to hydrologic extrapolation, (e.g. to the modelling of long hydrologic series and their long-term properties such as the Hurst effect) than to forecasting. Hydrologic forecasting, of course, is extrapolation in time; however, in the stochastic framework it is cast as an interpolation problem in the sense that the forecast period is treated simply as a time lag between two observations and an operational statistical relationship between them is established on the basis of correlation between the historic pairs of observations with the same lag. It is obvious that a forecasting model has the same limitations as those listed in the second section of this paper. It cannot correctly portray long-term properties or extreme values of the process being modelled if the historic series does not reflect them. Its application for short-term forecasting can, however, be very useful because the bulk of routine forecasts, i.e. paired observations, will be within the range of values of similar pairs on record. The danger is that forecasts of the occasional extreme events may be grossly inaccurate. On the other hand, even here there is a distinct possibility that the inaccuracy could be significantly reduced by a physically based model.

From the recent literature, one gets the impression that the most important thing in designing a reliable stochastic model of a hydrologic process is not so much the amount and quality of the available empirical data, (saying nothing about an understanding of the hydrological mechanisms involved) but simply a strict observance of detailed mathematical rules in the "identification, calibration and diagnostic checking (verification)" stages of model design. The purpose of this section is to demonstrate that it would be naive to expect that this approach must necessarily lead to stochastic RRR models with the correct structure and thus contribute to hydrologic knowledge. ("The words identification, calibration, and verification are misleading because of their connotation of greater understanding of and control over the physical processes than actually exist. Perhaps the words selection, parameter estimation and acceptance are more in line with the true capabilities of modelling" – Matalas and Maddock, 1976.)

It is known from theory as well as from experience that many processes of the hydrologic cycle are nonlinear. The well known simple example is the nonlinearity of the relation between the volume of water in storage and the corresponding rate of spontaneous discharge (eq. 1). This example will be used to demonstrate that stochastic properties of input and output as estimated from observations are a poor indicator of the dynamic structure of the system. The relationships between the input and output $C_V$ and $C_S$ shown in Fig. 1 will be used as a basis for the demonstration.

## The Problem of Inseparability of Parameter and Model Uncertainty

In stochastic modelling it is customary to separate parameter uncertainty from model uncertainty for the simple reason that the structure of the model is usually postulated a priori which makes it possible to separate out parameter uncertainty and study it as a sampling problem of the given model. A satisfactory method for separating model uncertainty does not seem to be feasible because the problem of model and parameter uncertainty is asymmetrical: while a given model is fully specified by a small number of parameters, a given set of parameters can be used for the specification of an unlimited number of different models. For example, the first three statistical moments can be common to many three-parameter probability distribution models. Thus even if we knew the population values of these three moments we would not be able to make an exhaustive analysis of model uncertainty.

In practice, even the theoretically feasible first part of the problem is out of reach of purely stochastic analysis because the latter cannot unambiguously specify the correct model: its structure must either be deduced from a physical theory or else, and this is the course normally taken in stochastic analysis, rather arbitrarily selected from the repertoire of models that are currently in use. This selection, euphemistically called stochastic model identification, is based on a tortuous cross-examination of the various inherently uncertain sample parameters so that the uncertainty of the chosen model is directly proportional to the parameter uncertainty.

## Parameter Uncertainty of Nonlinear Model

A mathematical analysis of parameter uncertainty of a nonlinear model given by eq. 1 is a difficult problem because it does not seem possible to obtain the parameters in an explicit form. For example, output variance involves the second, the $(1/b)$th and $(1/b + 1)$st moments of output, and output skewness involves the following output moments: 3rd, $(1/b)$th, $(1/b + 2)$nd, $(2/b)$th and $(2/b + 1)$st. Qualitatively it is immediately obvious that the uncertainty will be high because of the order of moments involved. Thus for $b = 0.5$ the coefficient of output skewness would involve the 3rd, 4th and 5th moment of output and for $b = 0.2$ the 3rd, 5th, 7th, 10th and 11th moment.

An indication of parameter uncertainty for the nonlinear case can be obtained from the linear case where the exact population parameters of output can be computed from specified input parameters (for independent input see Klemeš, 1978). Table 1 shows an example of the rate of convergence of the sample coefficient of skewness $C_S[Y]$ of output from a linear reservoir $Y = aS$ for the case shown in Fig. 1. The table gives the differences between the sample estimate of $C_S[Y]$ and its true value, expressed in percent of the true value. The dependence on a (a determines the output serial correlation) is quite obvious but it is clear that the sample values fluctuate widely and that the convergence

rate is slow. For the nonlinear case such a comparison cannot be made but one can get some idea from a comparison of fluctuations of sample estimates of $C_S$ with similar fluctuations for the linear case. Such a comparison is shown in Fig. 2 for two reservoirs which reduce the input variability to approximately the same value (see Fig. 1). The result indicates that, for the nonlinear case, even a sample size n = $10^6$ does not give a clear indication of the population value of the coefficient of skewness.

Table 1   Differences between Samples Values of Skewness Coefficient of Output from a Linear Reservoir and Population Value (in percent of Population Value)

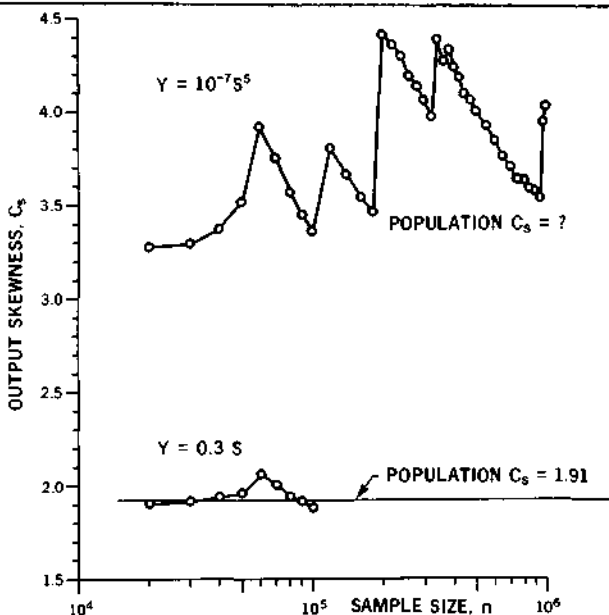| Sample | Coefficient a in eq. Y = aS | | | | |
|---|---|---|---|---|---|
| Size n | 0.03 | 0.1 | 0.3 | 1.0 | 3.0 |
| 20 | -2200 | -394 | 170 | 170 | 197 |
| 40 | - 246 | -118 | - 22 | 31 | 116 |
| 60 | - 35 | -123 | - 87 | -44 | 52 |
| 80 | 45 | - 78 | - 61 | - 30 | 65 |
| 100 | 121 | - 99 | - 84 | - 38 | 59 |
| 200 | 130 | 15 | - 14 | - 5.8 | 83 |
| 400 | 22 | - 3.5 | 3.8 | 5.3 | 87 |
| 2 000 | - 11 | 24 | 24 | 13 | 5.5 |
| 4 000 | - 8.4 | 10 | 9.5 | 5.4 | 2.2 |
| 6 000 | - 30 | 3.9 | 7.3 | 4.4 | 1.8 |
| 8 000 | - 49 | - 3.9 | 3.2 | 2.6 | 1.2 |
| 10 000 | - 54 | - 5.2 | 1.9 | 1.5 | 0.69 |
| 20 000 | - 48 | - 8.2 | - 1.1 | 0.23 | 0.17 |
| 40 000 | - 15 | 1.4 | 2.7 | 1.7 | 0.72 |
| 60 000 | - 16 | 2.0 | 2.7 | 0.82 | 0.13 |
| 80 000 | - 7.5 | 3.2 | 2.4 | 0.69 | 0.20 |
| 100 000 | - 4.8 | 2.5 | 1.6 | 0.52 | 0.19 |



Fig. 2   Convergence of sample skewness of output from a nonlinear reservoir Y = $10^{-7}$ $S^5$ and from a linear reservoir Y = 0.3 S. (Outputs from both reservoirs have approximately equal coefficients of variation, 0.34 and 0.36, respectively - compare Fig. 1).

147

A correct identification of model structure by stochastic analysis is virtually impossible, especially if nonlinear elements cannot be ruled out. As Moran (1975) observed,

"... entirely different physical models may lead to the same stochastic behaviour and even if they do not, the difference may be too small for discrimination in a reasonable sample... Suppose into some system there is an observed input $X(t)$ $(-\infty < t < \infty)$, and an observed output $Y(t)$. We take $Y(\tau)$ to be a functional, $F(X(t))$... of $X(t)$ for $t < \tau$. We seek to determine as much as possible of the structure of $F(X(t))$ from an observed record. This is the situation in studying rainfall runoff, and $F(X(t))$ may be highly nonlinear. Identifiability is then a key question... it is clear that nonlinearity is an all pervading problem and here we are confronted, if not with a brick wall, at any rate with a hill of rapidly increasing slope."

As an illustration consider the following problem. In Fig. 1, an output with a given pair of $C_v$ and $C_s$ can be obtained from the input shown (random lognormal series with $C_v = 1$ and $C_s = 4$) in a unique way by one reservoir and in a countless number of different ways by two or more reservoirs. For example, one linear reservoir defined by $Y = 0.26\ S$ yields an output, $Y_a$ say, with $C_v = 0.34$ and $C_s = 1.85$; an output $Y_b$ with the same parameters can be obtained when the input is routed through a cascade of 2 nonlinear reservoirs, $Y_1 = 3.5 \times 10^{-7}\ S^5$, and $Y_2 = 2.43\ S_2^{1/5}$. Nevertheless, the two outputs, $Y_a$ and $Y_b$, are significantly different. This is apparent from the distribution of the differences $Y_a-Y_b$ and from their autocorrelation function (Fig. 3), as well as from the distributions of $Y_a$ and $Y_b$ (Fig. 4); however, the differences all but disappear in the autocorrelation functions of $Y_a$ and $Y_b$ (Fig. 5). It must be pointed out that autocorrelation functions (together with their various modifications and transformations) are the main tool in stochastic model identification. The reason why these tools are quite useless for real identification of a hydrological model is that they are based on linear theory. To try to identify nonlinearity of a system from the autocorrelation function of its output is about as effective as trying to determine the degree of a polynomial from a straight line fitted to its segment. While this is well known to mathematicians, stochastic hydrologists keep pushing their heads into the sand and devising new "efficient" identification procedures based on one and the same linear theory. Let us return once more to Moran (1975): "Nearly all work done on time series has been from the point of view of spectral analysis which is a linear theory and is not invariant under nonlinear transformations of observed values. Thus many questions remain unanswered".

## Diagnostic Checking

It could be argued that the final step in the recommended procedure for stochastic model design – the diagnostic checking – would reveal the inadequacy in model structure that may have slipped through the identification stage. Thus in the case discussed in the preceding section, the diagnostic check would show that the residuals $Y_a-Y_b$ are neither normally distributed nor independent (Fig. 3). Hence the cascade of two nonlinear reservoirs yielding the output $Y_b$ could not pass for an adequate model for a system which in reality consists of one linear reservoir and yields an output $Y_a$, or vice versa. While this is true, there still are Moran's "many questions (that) remain unanswered".
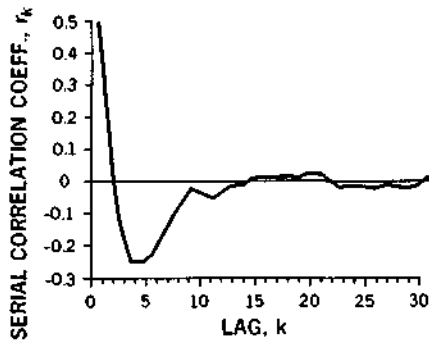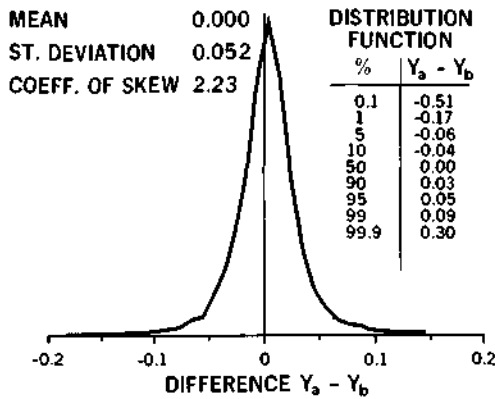
Fig. 3    Histogram and autocorrelation functions of differences between outputs $Y_a$ and $Y_b$ where $Y_a$ is output from one linear reservoir $Y = 0.26\ S$ and $Y_b$ is output from a cascade of two nonlinear reservoirs, $Y_1 = 3.5 \times 10^{-7}\ S^5$, $Y_2 = 2.43\ S^{1/5}$. In both cases input is the same as in Fig. 1 and sample size $n = 10\ 000$.
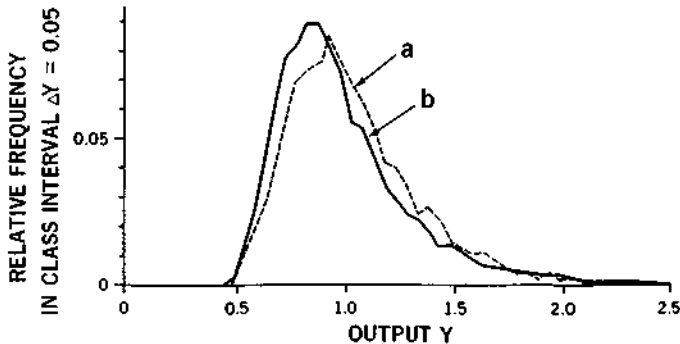


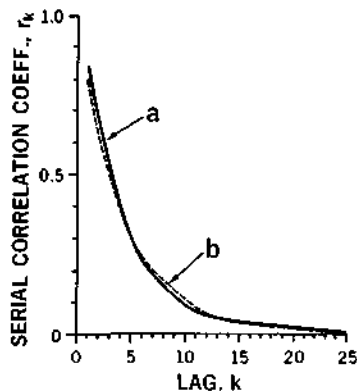Fig. 4    Histograms of outputs $Y_a$ and $Y_b$ (see legend to Fig. 3).

149

Fig. 5    Autocorrelation functions of outputs $Y_a$ and $Y_b$ (see legend to Fig. 3).

One problem is that it is a long way from showing that one particular model structure is inadequate to an identification of the correct structure.

Another problem is that of sample size which in practice will often not be large enough to render the departures from normality and independence statistically significant.

The most important practical problem, however, is that a model may well have a correct structure even if the common diagnostic checks fail. The reason for this is the unknown noises in the inputs and in the system. Consider the following example. Rainfall X is measured by point measurements and rainfall input X' into the model is computed from them in some "standard" manner (e.g. by the Thiessen polygon method). Let us assume that the correct structure of basin model is one nonlinear reservoir specified by the relation $Y = 0.02\ S^2$. The recorded streamflow, $Y_R$, which we assume to be measured with a high accuracy, then represents the true output from this nonlinear reservoir fed with the true rainfall X which we do not know. In our rainfall-runoff model, we route the input X' through this nonlinear reservoir and obtain a model output $Y_M$. To check the adequacy of our model, we subject the residuals $\Delta Y = Y_M - Y_R$ to diagnostic checking. Let us assume that the true rainfall X is a random series with a lognormal distribution with parameters $E[X] = 1$, $\sigma[X] = 1$, as in our previous examples, and that our computed rainfall input X' contains a random noise $z = z'- 0.1$, where $z'$ is lognormal with $E[z'] = 0.1$ and $\sigma[z'] = 0.1$. In our model, $X' = X + z$ is routed through the nonlinear reservoir to yield a noisy output $Y_M$ which must be different from the recorded $Y_R$ despite the fact that our basin model is perfect. The residuals $\Delta Y$ have parameters $E[\Delta Y] = 0$, $\sigma[\Delta Y] \approx 0.034$ and $C_S = 1.45$, and their distribution and autocorrelation functions are shown in Fig. 6. Their properties clearly indicate that the model structure is wrong which, of course, is not true. By fudging the model to give residuals with "desirable" properties we would simply distort the basin model by making its parameters and structure compensate for the unknown noise in our rainfall input. Such a model might be adequate for interpolation of missing streamflow data and for runoff simulation for rainfall conditions generally the same as in the past, including the same method of rainfall measurement and rainfall input computation. The use of such a model, for instance, for an evaluation of the impact of different rainfall patterns (say due to climatic change), or for runoff simulation with rainfall data obtained by different
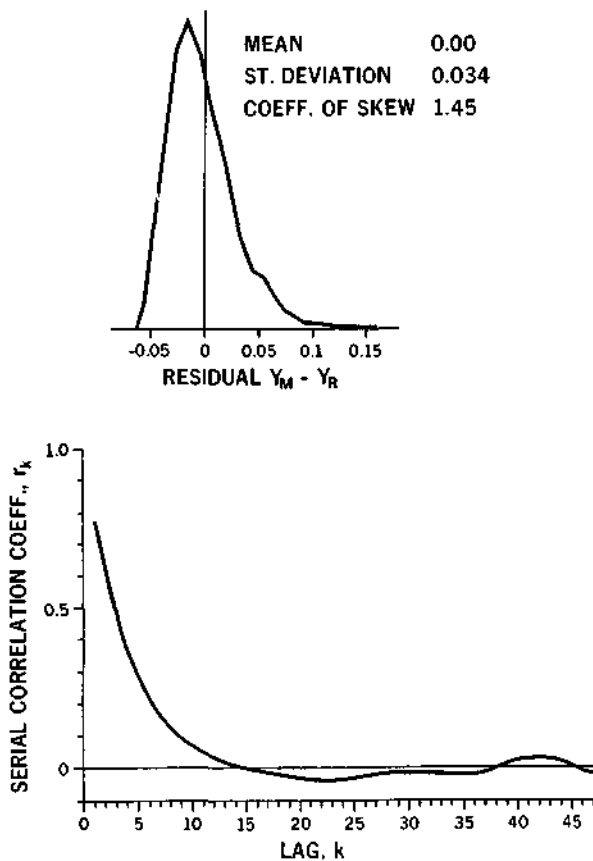
Fig. 6    Histogram and autocorrelation functions of residuals of modelled
output $Y_M$ from true output $Y_R$ from a nonlinear reservoir
$Y = 0.2\ S^2$.  Modelled output $Y_M$ was obtained from input
contaminated by random lognormally distributed noise with zero
mean and standard deviation 0.1.

instruments or from a different network, or with basin rainfall computed
in a different manner, etc., may lead to large errors, especially as far
as extreme values are concerned.

To summarize, the methods for stochastic diagnostic checking
that are currently being advocated in the literature cannot be relied
upon because they imply assumptions which are rarely satisfied in
practice.

## CONCLUSIONS

Stochastic modelling of RRR is in its early infancy.  The models

in current use fall by and large into the category of operational models and can be considered merely as interpolation formulae of hydrologic variables within the period of record. Within this scope they can be useful for various water resource management purposes, such as short-term hydrologic forecasting, reconstruction of missing data, etc. Their use for extrapolation and for any estimates which by their nature are not verifiable by measurement and observation (e.g. 100-year flood, estimation of runoff properties for changed climatic conditions, etc.) may lead to large errors. Although these models have no intrinsic hydrologic value in the sense of their contribution to hydrologic knowledge, they can, in the absence of better tools, be of use in water management, not as recipes for truth, but as a basis for establishing conventions and standardization which are useful for planning and design even if their correctness cannot be guaranteed (if our 100-year floods are in fact 300-year floods it does not really matter; what matters is to have a reasonable convention for the computation of a reasonably rare flood).

In order to be transformed into hydrologic models, stochastic models of RRR must be conceptualized, i.e. put on a physical basis. This however is also true for deterministic RRR models which, even when called conceptual, are in reality to a great extent only disguised operational models and cannot be relied on in situations for which they have not been "calibrated".

The first prototype of a physically based stochastic RRR model is Eagleson's (1978) model for annual water balance. It shows both the potentials and the difficulties of this approach. The way to progress leads through a thorough understanding of the impacts of various physical (hydrodynamical, thermodynamical, etc.) systems on stochastic properties of their inputs, and of the consequences of the various uncertainties inherent in the systems for the validity of the empirical laws describing their behaviour. These laws (e.g. Darcy's law) are often "macroscopic", i.e. involve relations between averages of stochastic properties and may not be applicable for model formulations at deeper levels because of nonlinear transformations that have been involved in the formation of the averages to which our physical laws often relate. Hence the proper structure of stochastic RRR models cannot be identified by classical methods of stochastic time series analysis but must be derived from physical theory. The theory of stochastic analysis has been worked out chiefly for the simplest mathematical assumptions of normality and linearity. Whenever it is stretched beyond this framework, the results are uncertain and often misleading. It is a rather humorous paradox that, while mathematicians are seeking refuge from mathematical difficulties of stochastic analysis in the physical mechanisms underlying complex stochastic processes, physical scientists and engineers are looking up to simplistic stochastic analysis for an enlightenment on their complex physical systems.


# REFERENCES

Anis, A.A., Lloyd, E.H. and Saleem, S.D. 1979. The Linear Reservoir with Markovian Inflows. Water Resources Research, 16(6), 1623-1627.

Berlinski, D. 1976. On Systems Analysis. MIT Press, Cambridge, Massachussetts.

Box, G.E.P. and Jenkins, G.M. 1970. Time Series Analysis: Forecasting and Control. Holden-Day, San Francisco.

Bras, R.L. 1979. Sampling of Interrelated Random Fields: The Rainfall-Runoff Case. Water Resources Research, 15(6), 1767-1780.

Clarke, R.T. 1980. Bivariate Gamma Distributions for Extending Annual Streamflow Records from Precipitation: Some Large-Sample Results. Water Resources Research, 16(5), 863-870.

Duckstein, L., Fogel, M. and Bogardi, I. 1979. Event-Based Models of Precipitation for Semiarid Lands. In: Proc. Symp. The Hydrology of Areas of Low Precipitation, Canberra, IAHS Publ. No. 128, pp. 51-64.

Eagleson, P.S. 1978. Climate, Soil, and Vegetation, 1 to 5, Water Resources Research, 14(5), 705-776.

Freeze, R.A. 1979. A Stochastic-Conceptual Analysis of Rainfall-Runoff Processes on a Hillslope. Water Resources Resarch, 16(2), 391-408.

Glynn, J. 1981. Negatively Skewed Runoff from Nonlinear Reservoir. National Hydrology Research Institute, Environment Canada, Ottawa (Unpublished report).

Gupta, V.K. and Wagmire, E.C. 1979. A Stochastic Kinematic Study of Subsynoptic Space-Time Rainfall. Water Resources Research, 15(3), 637-644.

Klemeš, V. 1970. Negatively Skewed Distribution of Runoff. In: Symposium of Wellington (N.Z.), pp. 219-236, IASH Publ. No. 96.

Klemeš, V. 1978. Physically Based Stochastic Hydrologic Analysis. In: V.T. Chow (Editor), Advances in Hydroscience, Vol. 11, pp. 285-355, Academic Press, New York.

Klemeš, V. 1981. Empirical and Causal Models in Hydrology. In: Panel on Scientific Basis of Water Resource Management, Scientific Basis of Water Management (tentative title), National Academy of Sciences, Washington, D.C. (in press).

Matalas, N.C. and Maddock III, T. 1976. Hydrology Semantics. Water Resources Research, 12(1), 123.

Mead, D.W. 1950. Hydrology. Second edition. McGraw-Hill, New York (first published in 1917).

Moran, P.A.P. 1957. The Statistical Treatment of Flood Flows. Trans. Amer. Geophys. Union, Vol. 38, 519-523.

Moran, P.A.P. 1975. The Future of Stochastic Modelling. In: Proc. Internat. Congress of Mathematicians, Vancouver, 1974. Canad. Math. Congress, pp. 517-521.

Sagar, B. 1979. Solution of Linearized Boussinesq Equation with Stochastic Boundaries and Recharge. Water Resources Research, 16(3), 618-624.

Singh, V.P. and Birsoy, Y.K. 1977. Some Statistical Relationships Between Rainfall and Runoff. Journ. Hydrology, 34, 251-269.

Smith, L. and Freeze, R.A. 1979a. Stochastic Analysis of Steady State Groundwater Flow in a Bounded Domain, 1, One-Dimensional Simulations. Water Resources Research, 15(3), 521-528.

Smith, L. and Freeze, R.A. 1979b. Stochastic Analysis of Steady State Groundwater Flow in a Bounded Domain, 2, Two-Dimensional Simulations. Water Resources Research, 15(6), 1543-1559.

Sokolovskii, D.L. 1971. River Runoff. Israeli Program for Scientific Translations, Jerusalem. Available from U.S. Department of Commerce, Springfield, Virginia (Russian original published in 1968).

Sorman, A.U. 1975. Characteristics of Rainfall Cell Patterns in the Southeast Coastal Plain Areas of the USA, and a Computer Simulation Model of Thunderstorm Rainfall. Proc. Symp. Application of Mathematical Models in Hydrology and Water Resources Systems, Bratislava. IAHS Publ. No. 115, pp. 214-221.

Waymire, E.E.C. and Gupta, V.K. 1980. The Mathematical Structure of Rainfall. Submitted to Water Resources Research.

Wilson, C.B., Valdes, J.B. and Rodriguez-Iturbe, I. 1979. On the Influence of the Spatial Distribution of Rainfall on Storm Runoff. Water Resources Research, 15(2), 321-328.