

## **HOW THE EXPERTFIT DISTRIBUTION-FITTING PACKAGE CAN MAKE YOUR SIMULATION MODELS MORE VALID**

Averill M. Law  
Michael G. McComas

Averill M. Law and Associates, Inc.  
P.O. Box 40996  
Tucson, AZ 85717, U.S.A.

### **ABSTRACT**

In this paper, we discuss the critical role of simulation input modeling in a successful simulation study. Two pitfalls in simulation input modeling are then presented and we explain how any analyst, regardless of their knowledge of statistics, can easily avoid these pitfalls through the use of the ExpertFit distribution-fitting software. We use a set of real-world data to demonstrate how the software automatically specifies and ranks probability distributions, and then tells the analyst whether the “best” candidate distribution is actually a good representation of the data. If no distribution provides a good fit, then ExpertFit can define an empirical distribution. In either case, the selected distribution is put into the proper format for direct input to the analyst’s simulation software.

### **1 THE ROLE OF SIMULATION INPUT MODELING IN A SUCCESSFUL SIMULATION STUDY**

In this section we describe simulation input modeling and show the consequences of performing this critical activity improperly.

#### **1.1 The Nature of Simulation Input Modeling**

One of the most important activities in a successful simulation study is that of representing each source of system randomness by a probability distribution. For example in a manufacturing system, processing times, machine times to failure, and machine repair times should generally be modeled by probability distributions. If this critical activity is neglected, then one’s simulation results are quite likely to be erroneous and any conclusions drawn from the simulation study suspect – in other words, “garbage in, garbage out.”

In this paper, we use the phrase “simulation input modeling” to mean the process of choosing a probability

distribution for each source randomness for the system under study and of expressing this distribution in a form that can be used in the analyst’s choice of simulation software. In Sections 2 and 3 we discuss how an analyst can easily and accurately choose an appropriate probability distribution using the ExpertFit software. Section 4 discusses important features that have recently been added to ExpertFit.

#### **1.2 Two Pitfalls in Simulation Input Modeling**

We have identified a number of pitfalls that can undermine the success of a simulation study (Law and Kelton 2000). Two pitfalls that directly relate to simulation input modeling are discussed in the following two sections [see our Web Site <[www.averill-law.com](http://www.averill-law.com)> (“ExpertFit Distribution Fitting Software”) for further discussion of pitfalls, and for a more comprehensive discussion of ExpertFit, in general].

##### **1.2.1 Pitfall Number 1: Replacing a Distribution by its Mean**

Simulation analysts have sometimes replaced an input probability distribution by its perceived mean in their simulation models. This practice may be caused by a lack of understanding of this issue on the part of the analyst or by lack of information on the actual form of the distribution (e.g., only an estimate of the mean of the distribution is available). Such a practice may produce completely erroneous simulation results, as is shown by the following example.

Consider a single-server queueing system (e.g., a manufacturing system consisting of a single machine tool) at which jobs arrive to be processed. Suppose that the mean interarrival time of jobs is 1 minute and the mean service time is 0.99 minute. Suppose further that the interarrival times and service times each have an exponential distribution. Then it can be shown that the

long-run mean number of jobs waiting in the queue is *approximately 98*. On the other hand, suppose we were to follow the dangerous practice of replacing a source of randomness with a constant value. If we assume that each interarrival time is *exactly* 1 minute and each service time is *exactly* 0.99 minute, *then each job is finished before the next arrives and no job ever waits in the queue!* The variability of the probability distributions, rather than just their means, has a significant effect on the congestion level in most queueing-type (e.g., manufacturing) systems.

### 1.2.2 Pitfall Number 2: Using the Wrong Distribution

We have seen the importance of using a distribution to represent a source of randomness. However, as we will now see, the actual distribution used is also critical. It should be noted that many simulation practitioners and simulation books widely use normal input distributions, even though in our experience this distribution will *rarely* be appropriate to model a source of randomness such as service times.

Suppose for the queueing system in Section 1.2.1 that jobs have exponential interarrival times with a mean of 1 minute. We have 200 service times that have been collected from the system, but their underlying probability distribution is unknown. Using ExpertFit, we fit the best Weibull distribution and the best normal distribution (and others) to the observed service-time data. However, as shown by the analysis in Section 6.7 of Law and Kelton (2000), the *Weibull distribution* actually provides the best overall model for the data.

We then made a *very long* simulation run of the system using *each* of the fitted distributions. The average number of jobs in the queue for the Weibull distribution was 4.41, which should be close to the average number in queue for the actual system. On the other hand, the average number in queue for the normal distribution was 6.13, corresponding to a *model output error of 39 percent*. It is interesting to see how poorly the normal distribution works, given that it is the most well-known distribution.

We will see in Section 2 how the use of ExpertFit makes choosing an appropriate probability distribution a quick and easy process.

### 1.3 Advantages of Using ExpertFit

With the assistance of ExpertFit, an analyst, regardless of their prior knowledge of statistics, can avoid the two pitfalls introduced above. When system data are available, a complete analysis with the package takes just minutes. The package identifies the “best” of the candidate probability distributions, and also tells the analyst whether the fitted distribution is good enough to actually use in the simulation model. If none of the candidate distributions

provides an adequate fit, then ExpertFit can construct an empirical distribution. In either case, the selected distribution can be represented automatically in the analyst’s choice of simulation software. Appropriate probability distributions can also be selected when no system data are available. For the important case of machine breakdowns, ExpertFit will specify time-to-failure and time-to-repair distributions that match the system’s behavior, even if the machine is subject to blocking or starving.

## 2 USING EXPERTFIT WHEN SYSTEM DATA ARE AVAILABLE

We consider first the case where data are available for the source of randomness to be represented in the simulation model. Our goal is to give an overview of the capabilities of ExpertFit – a demo disk with a thorough discussion of program operation is available from the authors.

We have designed ExpertFit based on our 22 years of research and experience in selecting simulation input distributions. The user interface employs four tabs that are typically used sequentially to perform an analysis. Furthermore, the options in each tab have default settings to promote ease of use. All graphs are designed to provide definitive comparisons and to minimize possible analyst misinterpretation. For example, the following features are available:

- Multiple distributions can be plotted on the same graph
- Error graphs are automatically scaled so that the visual display of an error reflects the severity of the error
- Whenever possible, bounds for an acceptable error are displayed.

These software features make it easy for an analyst to perform an accurate and thorough analysis of a data set, regardless of their prior knowledge of statistics. On the other hand, the user interface is completely flexible so that an experienced analyst can easily access the full set of available tools for performing a comprehensive and complete analysis, in any order desired.

The first data-analysis tab has options for obtaining the data set and for displaying its characteristics. An analyst can read a data file, manually enter a data set, paste in a data set from the Clipboard, or import a data set from Excel. Once a data set is available, a number of graphical and tabular sample summaries can be created, including histograms, sample statistics, and plots designed to assess the independence of the observations.

The data set we have chosen for this example consists of 622 processing times for parts, which were provided to us by a major automobile manufacturer.

At the second tab distributions are fit to the data set. For the recommended automated-fitting option, the only information required by ExpertFit to begin the fitting and evaluation process is a specification of the range of the underlying random variable. Since all we know about the data is that the values are non-negative, we accepted the default limits of “zero” and “infinity.” ExpertFit responds by fitting distributions with a range starting at zero and also distributions whose lower endpoint was estimated from the data itself. These candidate models were then automatically evaluated and the results screen shown in Figure 1 was displayed.

ExpertFit fit and ranked 24 candidate models, with the three best-fitting models being displayed on the screen along with their scores. The displayed scores are calculated using a proprietary evaluation scheme that is based on our 22 years of experience and research in this area, including the analysis of 35,000 computer-generated data sets. Results from the heuristics that we have found to be the best indicators of a good model fit are combined and the resulting numerical evaluation is normalized so that 100 indicates the best possible model and 0 indicates the

worst possible model. These scores are *comparative* in nature and do not give an overall assessment of the quality of fit. ExpertFit provides a separate *absolute* evaluation of the quality of the representation provided by the best-ranked model. This absolute evaluation is absolutely critical because, perhaps, one third of all data sets are not well represented by a standard theoretical distribution. *Furthermore, ExpertFit is the only software package that provides such a definitive absolute evaluation.*

In Figure 1 we see that the Inverted Weibull distribution (with a range starting at zero) is the best model for the processing-time data. Furthermore, the Absolute Evaluation is “Good,” which indicates that this distribution is good enough to use in a simulation model.

However, it is generally desirable to confirm the quality of the representation using the third tab. Although the Inverted Weibull distribution may be unfamiliar to you, it can be used in almost all simulation packages since it is the inverse of a Weibull random variable. It should also be noted that ExpertFit completed the entire analysis without any further input from the analyst. After automated fitting, the analyst is automatically transferred to the third tab,

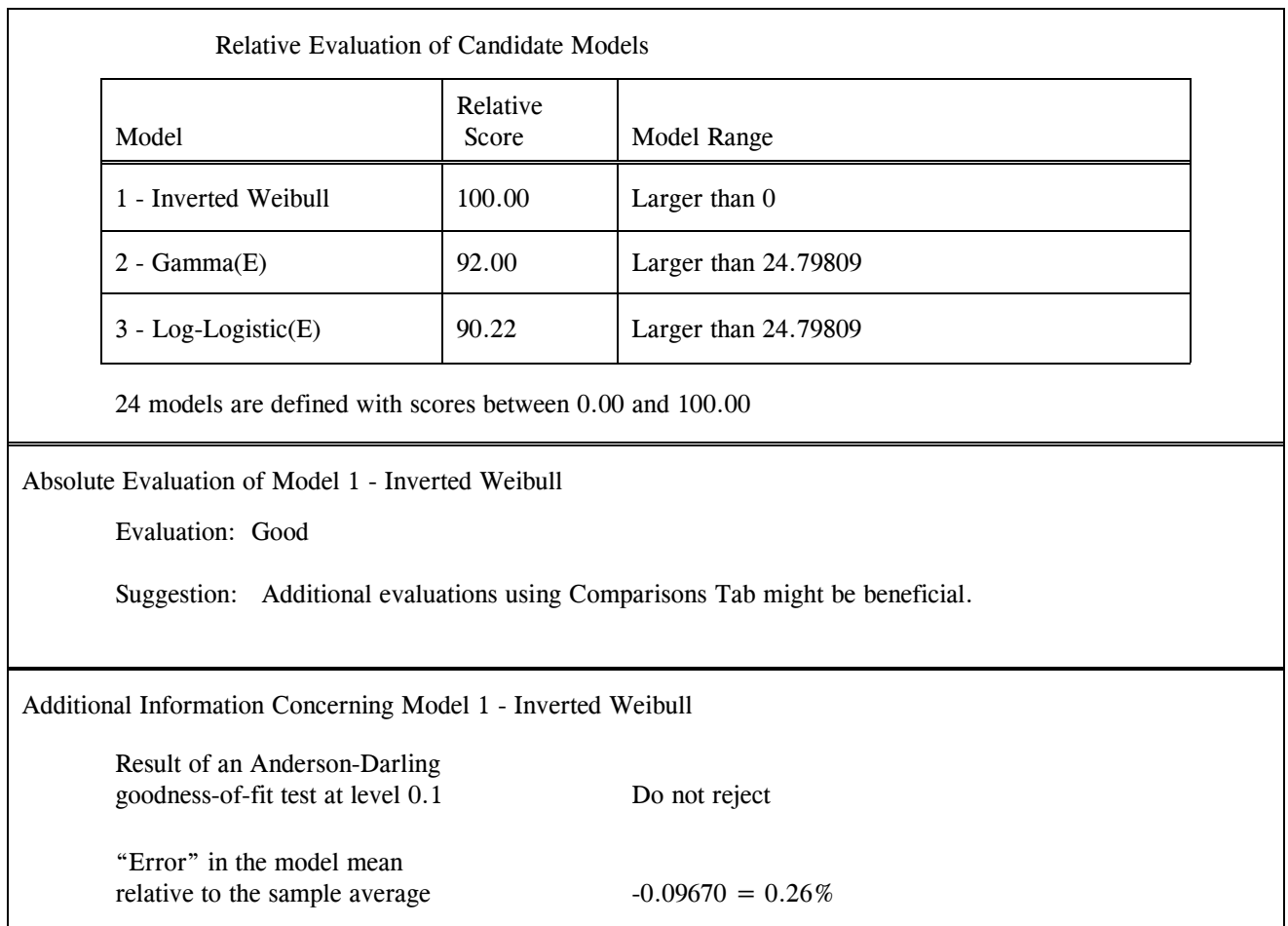


Figure 1: Evaluation of the Candidate Models for the Processing-Time Data

where the specified models can be compared to the sample to confirm the quality of fit (if additional confirmation is desired). Two of our favorite comparisons are the density/histogram overplot and the distribution function differences plot, which are shown in Figures 2 and 3, respectively. In the former case, the density function of the Inverted Weibull distribution has been plotted over a histogram of the data (a graphical estimate of the true density function). This plot indicates that the Inverted Weibull distribution is a good model for the observed data. The distribution function differences plot graphs the differences between a sample distribution function (a graphical estimate of the true distribution function) and the distribution function of the Inverted Weibull distribution. Since these vertical differences are small (i.e., within the horizontal error bounds), this also suggests that the

Inverted Weibull distribution is a good representation for the data. Note that the third tab also allows the analyst to perform several goodness-of-fit tests such as the chi-square and Kolmogorov-Smirnov tests. ExpertFit includes an option in the fourth tab for displaying the representation of the Inverted Weibull distribution using different simulation packages. We show in Figure 4 the representations for four of the simulation packages supported by ExpertFit.

For some data sets, no candidate model provides an adequate representation. In this case we recommend the use of an empirical distribution. Note that ExpertFit allows an empirical distribution to be based on all data values or on a histogram to reduce the information that is needed for specification. We show a histogram-based representation (with 20 intervals) for two simulation packages in Figure 5.

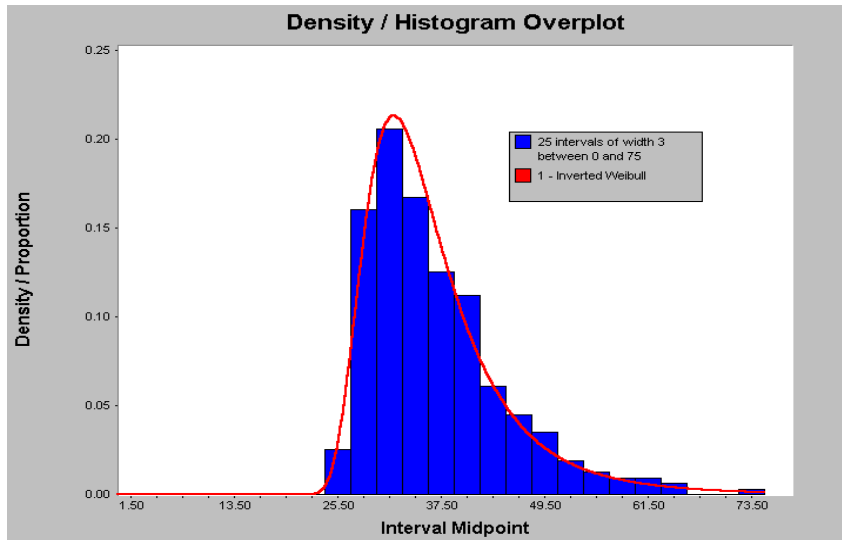


Figure 2: Density/Histogram Overplot for the Processing-Time Data

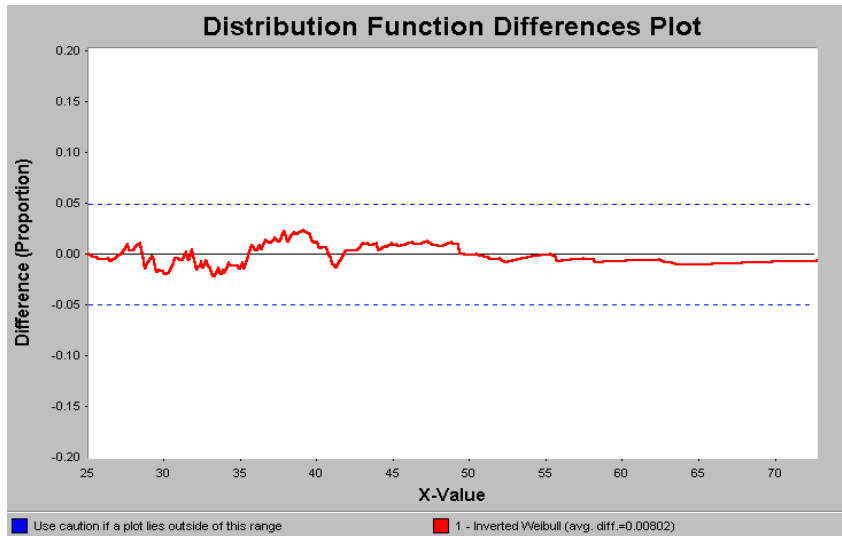


Figure 3: Distribution Function Differences Plot for the Processing-Time Data

Simulation Software	Representation
GPSS/H 3	RVIWEIB(<stream>,6.272056, 32.834140)
ProModel	InvWeibull(6.272056, 32.834140, <stream>, 0.000000)
Taylor ED	1./weibull(0.028324, 6.272056)
WITNESS	1./WEIBULL(6.272056, 0.030456, <stream>)

Figure 4: Simulation Software Representations of the Inverted Weibull Distribution

Simulation Software	Representation
Arena	CONT(0.0000,24.800000, 0.0322,27.185000, 0.1576,29.570000, 0.3183,31.955000, 0.4791,34.340000, 0.5981,36.725000, 0.6945,39.110000, 0.7942,41.495000, 0.8457,43.880000, 0.8778,46.265000, 0.9068,48.650000, 0.9421,51.035000, 0.9550,53.420000, 0.9711,55.805000, 0.9807,58.190000, 0.9839,60.575000, 0.9904,62.960000, 0.9968,65.345000, 0.9968,67.730000, 0.9968,70.115000, 1.0000,72.500000)
AutoMod	continuous(0.0000:24.800000,0.0322:27.185000,0.1576:29.570000, 0.3183:31.955000,0.4791:34.340000,0.5981:36.725000,0.6945:39.110000, 0.7942:41.495000,0.8457:43.880000,0.8778:46.265000,0.9068:48.650000, 0.9421:51.035000,0.9550:53.420000,0.9711:55.805000,0.9807:58.190000, 0.9839:60.575000,0.9904:62.960000,0.9968:65.345000,0.9968:67.730000, 0.9968:70.115000,1.0000:72.500000)

Figure 5: Simulation Software Representations of the Empirical Distribution Function

### 3 USING EXPERTFIT WHEN NO DATA ARE AVAILABLE

Sometimes a simulation analyst must model a source of randomness for which no system data are available. ExpertFit provides two types of analyses for this situation. A general activity time (e.g., a service time) can be modeled in ExpertFit by using a triangular or beta distribution. In the case of a triangular distribution, the analyst specifies the distribution by giving subjective estimates of the minimum, maximum, and most-likely activity times.

ExpertFit will also help the analyst specify time-to-failure and time-to-repair distributions for a machine that randomly breaks down. In this case, the analyst gives, for example, subjective estimates for the percentage of time that the machine is operational (e.g., 90 percent) and for the mean repair time.

### 4 NEW FEATURES IN EXPERTFIT

The following are new ExpertFit features:

- A comprehensive library of probability distributions has been added that represent sources of randomness for many different applications. This library was developed by analyzing data from a large number of real-world simulation projects. This feature will be extremely useful when modeling a system for which little or no data exist.

- A batch-mode capability has been added to the Professional Version (supports 31 different simulation packages) of ExpertFit that allows one to enter and analyze a large number of data sets in a matter of seconds with only a few keystrokes.
- A distribution viewer has been added that allows one to see characteristics of a distribution without entering any data. By using a slider bar for each parameter, you can interactively and quickly change the distribution being viewed.
- A scroll bar has been added for interactively changing the histogram interval widths – this makes finding the “optimal” histogram much faster.

### 5 CONCLUSION

ExpertFit can help you develop more valid simulation models than if you use a standard statistical package, an input processor built into a simulation package, or hand calculations to determine input probability distributions. ExpertFit uses a sophisticated algorithm to determine the best-fitting distribution and, furthermore, has 39 built-in standard theoretical distributions. On the other hand, a typical simulation package contains roughly 10 distributions.

ExpertFit can represent most of its 39 distributions in 31 different simulation packages such as Arena, AutoMod, Extend, GPSS/H, Micro Saint, OPNET Modeler, ProModel, SES/workbench, SIMPLE++ (eM-Plant), SIMPROCESS, SIMUL8, Taylor ED, and WITNESS, *even*

though the distribution may not be available in the simulation package itself.

## REFERENCE

Law, A. M. and W. D. Kelton. 2000. *Simulation Modeling and Analysis*, 3d ed., New York: McGraw-Hill.

## AUTHOR BIOGRAPHIES

**AVERILL M. LAW** is President of Averill M. Law & Associates, Inc. (Tucson, Arizona), a company specializing in simulation consulting, training, and software. He has been a simulation consultant to more than 100 organizations, including Andersen Consulting, Boeing, Cellular One, DMSO, Kimberly-Clark, M&M/Mars, 3M, Xerox, the U.S. Air Force, and the U.S. Army. He has presented more than 325 simulation short courses in 17 countries, and delivered more than 100 talks on simulation modeling at technical conferences. He is the author (or coauthor) of three books and numerous papers on simulation, manufacturing, operations research, statistics, and communications, including the textbook *Simulation Modeling and Analysis* that is used by more than 70,000 people worldwide. He is the developer of the ExpertFit software package for selecting simulation input probability distributions, and he has developed several simulation videotapes. Dr. Law wrote a regular column on simulation for *Industrial Engineering* magazine from 1990 through 1991. He has been a tenured faculty member at the University of Wisconsin and the University of Arizona. Dr. Law has a Ph.D. in industrial engineering and operations research from the University of California at Berkeley. His email and web addresses are <averill@ix.netcom.com> and <www.averill-law.com>.

**MICHAEL G. MCCOMAS** is Vice President for Consulting Services of Averill M. Law & Associates, Inc. He has considerable simulation modeling experience in a wide variety of application areas, and is the coauthor of seven published papers on simulation. His educational background includes an M.S. in systems and industrial engineering from the University of Arizona. His email and web addresses are <averill@ix.netcom.com> and <www.averill-law.com>.