
ENE 3031

Computer Simulation

Week 9: Input Modeling

(attributed by Seong-Hee Kim and Barry L. Nelson)

Chuljin Park
Assistant Professor
Industrial Engineering
Hanyang University

Input Modeling

- Input models represent the **uncertainty** in a stochastic simulation.
- The fundamental requirements for an input model are:
 - It must be capable of representing the physical realities of the process.
 - It must be easily tuned to the situation at hand.
 - It must be amenable to random variate generation.



- There is no “true” model for any stochastic input. The best that we can hope is to obtain an approximation that yields useful results.
- A key distinction in input modeling problems is the presence or absence of data:
 - When we have data, then we *fit* a model to the data. Good software is available for this.
 - When no data are available then we have to creatively use what we can get to construct an input model.

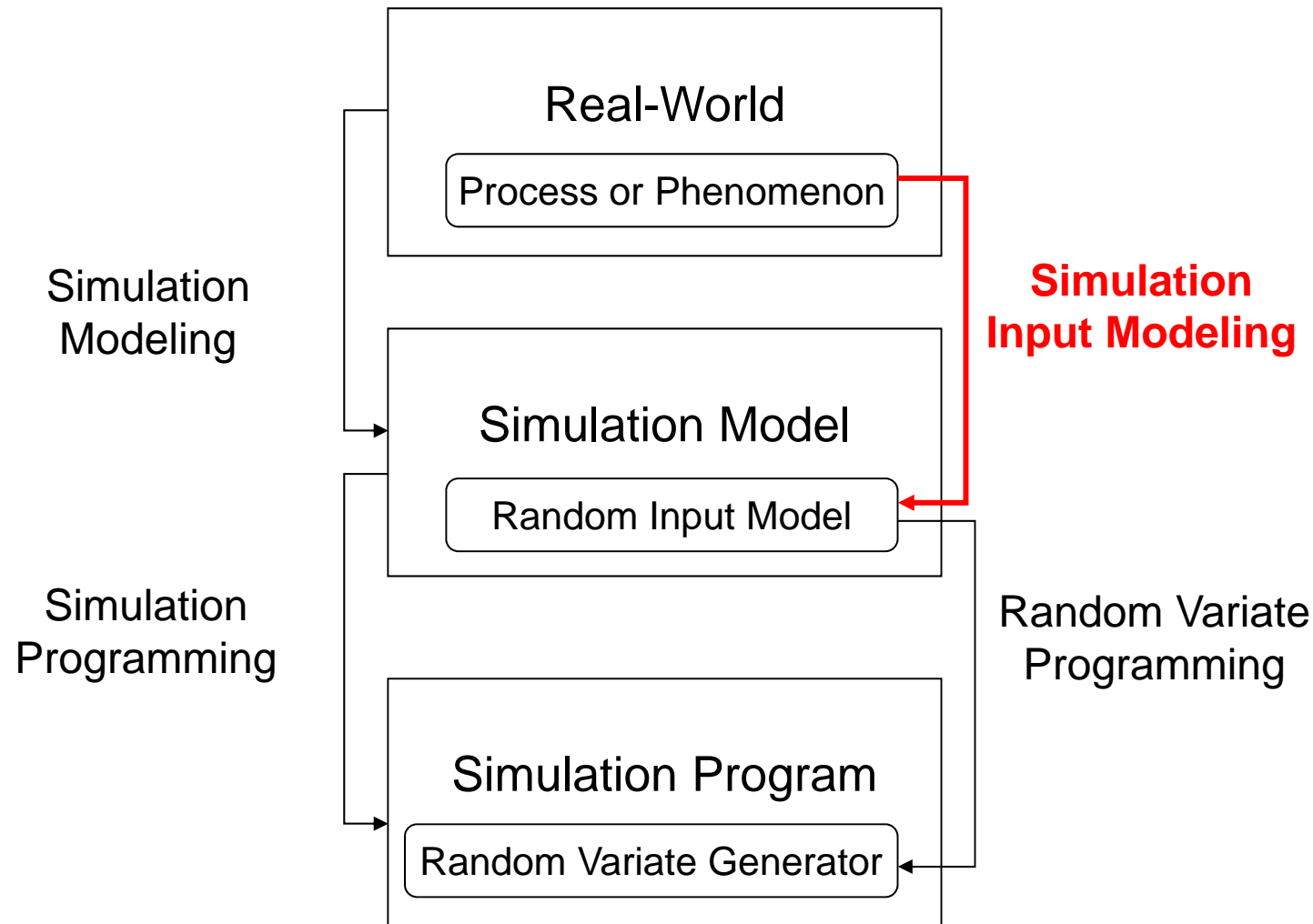


Outline

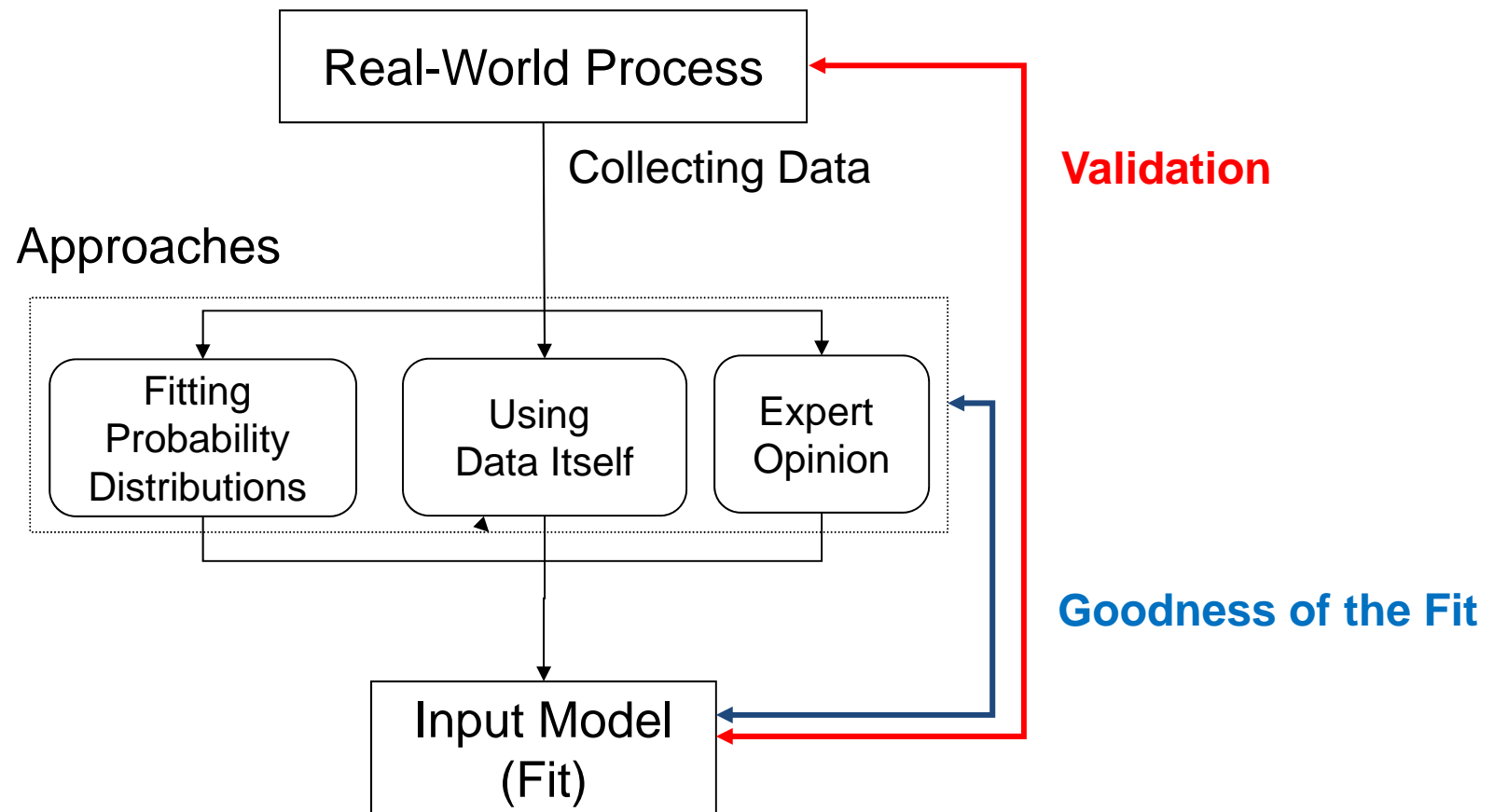
- Input modeling with data
 - physical basis for distributions
 - fitting and checking
 - ExpertFit
- Input modeling without data
 - sources of information
 - incorporating expert opinion



Simulation Model Development



Input Model Development



Why do we use input models?

Reliability Example

- Suppose you are a supplier of a component that is supposed to last for one year, a component that you know has a mean time to failure of 2 years.
- A client is willing to pay \$1000 for your component, but wants you to pay a penalty of \$5000 if failure occurs in less than one year
- Ignore the uncertainty: \$1000 profit per component
- Life time \sim Expo with mean 2 years : \$967 loss



Does the particular input model matter?

YES, IT DOES!

Reliability Example (Continued)

- Life Time \sim Expo with mean 2 years: \$967 loss
- Life Time \sim Unif(0,4years) : \$250 loss
- Uniform has the right mean (2 years) but it results in much smaller loss and causes you to underprice the component.



Input Modeling with Data

1. Select one or more candidate distributions, based on **physical characteristics** of the process and **graphical examination** of the data.
2. Fit the distribution to the data (determine values for its **unknown parameters**).
3. Check the fit to the data via tests and graphical analysis.
4. If the distribution does not fit, select another candidate and go to 2, or use an **empirical distribution**.



Physical Basis for Distributions

- Most probability distributions were invented to represent a particular physical situation.
- If we know the physical basis for a distribution, then we can match it to the situation we have to model
- A number of examples follow...



- **binomial**: Models the number of successes in n trials, when the trials are independent with common success probability, p . *Example: the number of defective components found in a lot of n components.*
- **negative binomial** Models the number of trials required to achieve k “successes.” *Example: the number of components that we must inspect to find 4 defective components.*



- **Poisson** Models the number of independent events that occur in a fixed amount of time or space. *Ex: number of customers that arrive to a store during 1 hour, or number of defects found in 30 cubic meters of sheet metal.*
- **Normal**: Models the distribution of a process that can be thought of as the *sum* of a number of component processes. *Ex: the time to assemble a product which is the sum of the times required for each assembly operation.*

- **lognormal** Models the distribution of a process that can be thought of as the *product* of a number of component processes. *Example: the rate of return on an investment, when interest is compounded, is the product of the returns for a number of periods. Also widely used to model stock prices.*
- **exponential**: Models the time between independent events, or a process time which is *memoryless*. *Example: the time to failure for a system that has constant failure rate over time. Note: if the time between events is exponential, then the number of events is Poisson.*



- **Erlang** The sum of k identical exponential random variables. A special case of the gamma...
- **gamma** An extremely flexible distribution used to model nonnegative random variables.
- **beta** : An extremely flexible distribution used to model bounded (fixed upper and lower limits) random variables.
- **Weibull** : Models the time to failure for components; can model increasing or decreasing failure rate hazard. *Ex: the time to failure for a disk drive.*

- **discrete or continuous uniform** Models complete uncertainty, since all outcomes are equally likely.
- **triangular** Models a process when only the minimum, most likely and maximum values of the distribution are known. *Ex: the minimum, most likely and maximum inflation rate we will have this year.*
- **empirical** Reuses the data themselves by making each observed value equally likely. Can be interpolated to obtain a continuous distribution.

Fitting

- Common methods for fitting distributions are *maximum likelihood*, *method of moments*, and *least squares*.
 - While the method matters, the variability in the data often overwhelms the differences in the estimators.
 - Remember: There is no “true distribution” just waiting to be found!



- Ways to check fit include the χ^2 , K-S and Anderson-Darling tests, and density-histogram and q - q plots.
 - Beware of goodness-of-fit tests because they are unlikely to reject *any* distribution when you have little data, and are likely to reject *every* distribution when you have lots of data.
 - Tests represent lack of fit by a *summary statistic*, while plots show where the lack of fit occurs and whether it is important.
 - χ^2 tests and density-histogram plots are sensitive to how we group the data.



Tests and p -values

- In the typical test...
 H_0 : the chosen distribution fits
 H_1 : the chosen distribution does not fit
- Test statistics calculate difference between data and a chosen distribution.
- Thus, a small test statistics supports a good fit.
- The p -value of a test is the Type I error level (significance) at which we would *just reject H_0 for the given data*.
- Thus, a large (> 0.10) p -value supports H_0 that the distribution fits.



ExpertFit

- Stand-alone software package for fitting distributions to data.
- Data > Data Summary: good for getting sample mean and sample variance
- Models > Automated Fitting: sort candidates from the best
- Models > View/Delete Models > Show Model Parameters: give fitted distributions



- Reports *test-statistic* for
 - Anderson-Darling test
 - Kolmogorov-Smirnov test
 - χ^2 test
- Also provide density-histogram, p-p and q-q plots.



Data Summary/Fitted Model

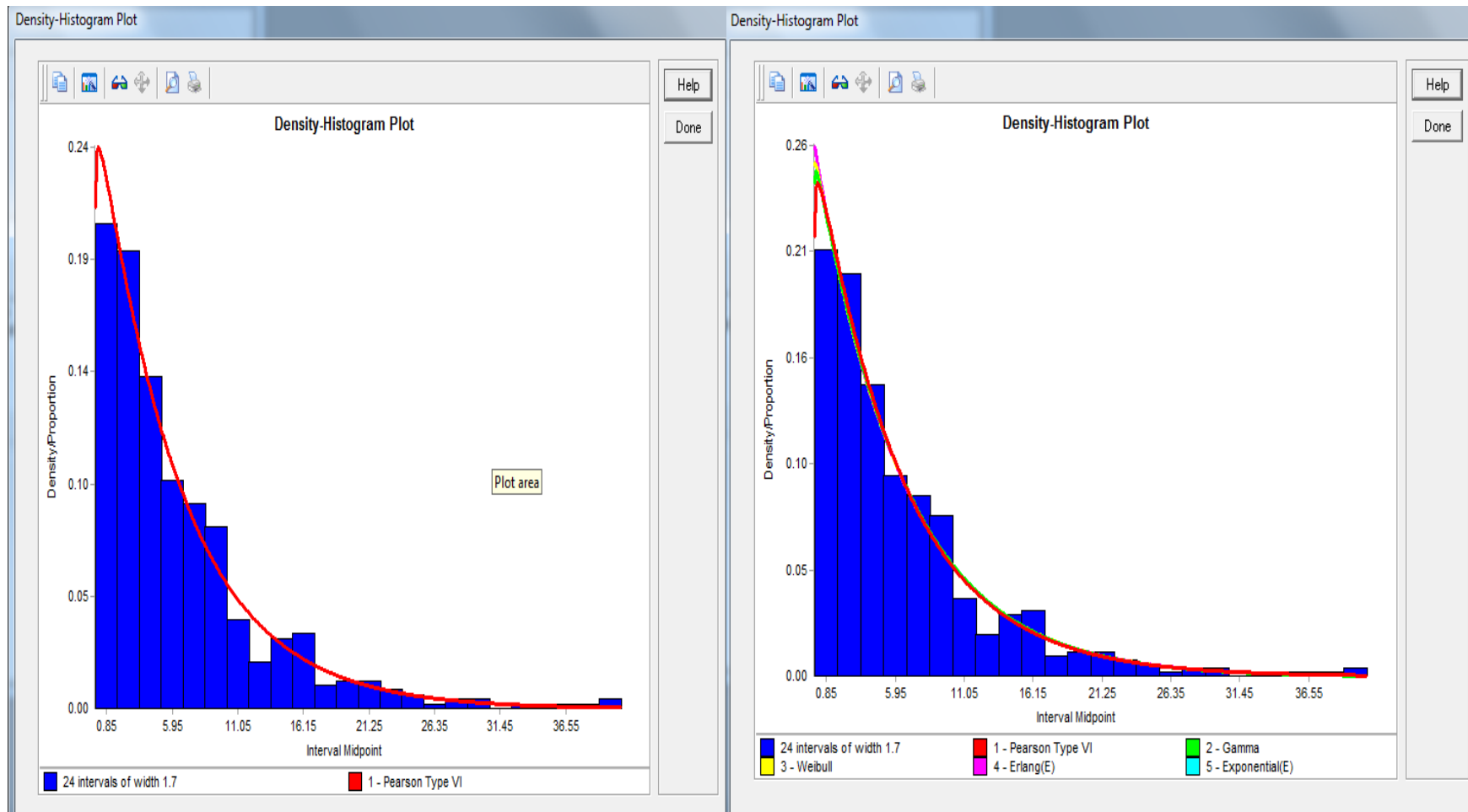
Data-Summary Table

Data Characteristic	Value
Source file	3044a
Observation type	Real valued
Number of observations	500
Minimum observation	0.00417
Maximum observation	40.21170
Mean	6.56469
Median	4.46179
Variance	44.04246
Coefficient of variation	1.01093
Skewness	2.03511

Fitted Models

Model	Parameters		
6 - Erlang	Location	Default	0.00000
	Scale	ML estimate	6.56469
	Shape	ML estimate	1
7 - Exponential	Location	Default	0.00000
	Scale	ML estimate	6.56469
8 - Beta	Lower endpoint	OPT estimate	1.43349 e -4
	Upper endpoint	OPT estimate	80.42180
	Shape #1	ML estimate	0.93629
	Shape #2	ML estimate	10.45035
9 - Johnson SB	Lower endpoint	OPT estimate	0.00000
	Upper endpoint	OPT estimate	43.87742
	Shape #1	ML estimate	1.54858
	Shape #2	ML estimate	0.67884
10 - Lognormal	Location	Default	0.00000
	Scale	ML estimate	3.71528
	Shape	ML estimate	1.29332
11 - Random Walk	Location	Default	0.00000
	Scale	ML estimate	1.79107

Density-Histogram Plot



Usage Notes

- The “Automated Fitting” option tries all relevant distributions and ranks distributions from the best.
- Be sure to try different numbers of histogram cells; it affects the *test-statistic* of the χ^2 test, and your perception of the fit.



- Since exponential is a special case of Erlang which is a special case of gamma, “Automated Fitting” rarely selects exponential or Erlang. Similarly, exponential is a special case of Weibull.
- Raw data can be read in from text files (looks for .dat), one value per line.



Distributions in Simio

Random.Bernoulli	Random.Lognormal
Random.Beta	Random.NegativeBinomial
Random.Binomial	Random.Normal
Random.Erlang	Random.PearsonVI
Random.Exponential	Random.Pert
Random.Gamma	Random.Poisson
Random.Geometric	Random.Triangular
Random.JohnsonSB	Random.Uniform
Random.JohnsonUB	Random.Weibull
Random.LogLogistic	

Usage Notes

- Parameter definitions often do not match with Simio definitions.
- Check [Help > User's Guide](#) before setting parameters.
- Be sure to check if your entered input distribution generates random variates from the right distribution.

q-q Plot

- One way to generate data from cdf F is via

$$Y = F^{-1}(R)$$

- The *q-q* plot displays the *sorted* data

$$Y_1 \leq Y_2 \leq \dots \leq Y_n$$

vs.

$$F^{-1}\left(\frac{j-1/2}{n}\right), j = 1, 2, \dots, n$$



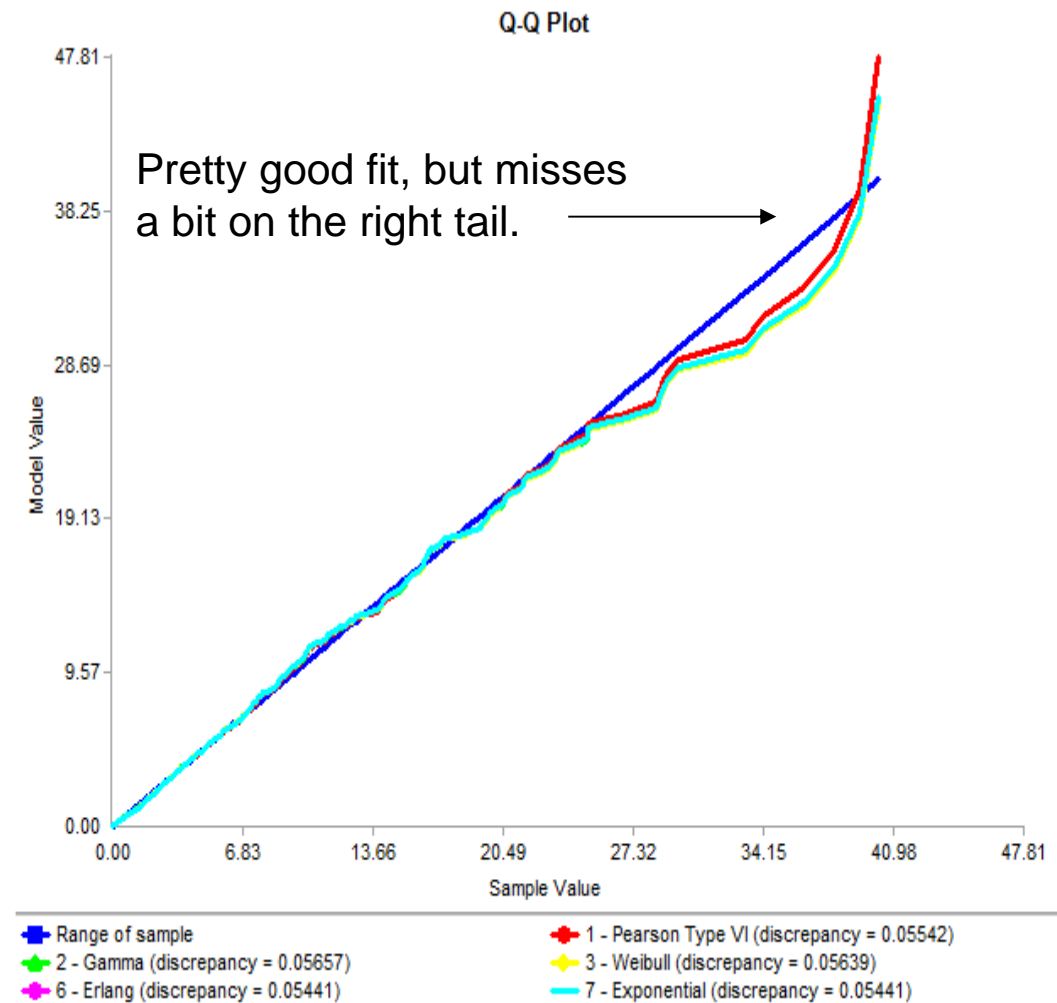
Features of the q - q Plot

- It does not depend on how the data are grouped.
- It is much better than a density-histogram when the number of data points is small.
- Deviations from a **straight line** show where the distribution does not match.



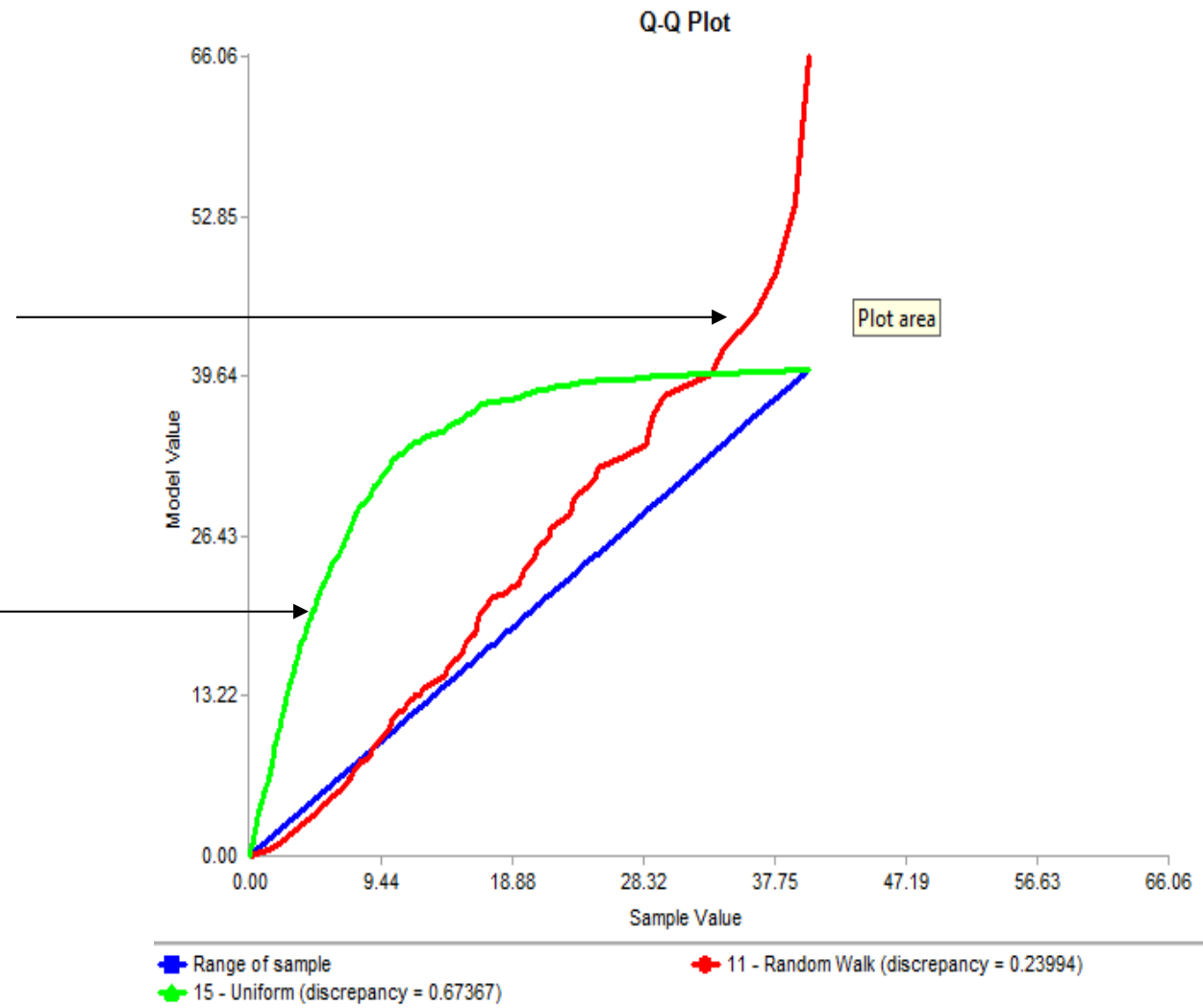
Mode (from the top menu) > Advance > Comparisons
> Probability Plots > Q-Q Plot

A straight line implies the *family* of distributions is correct; a 45° line implies correct *parameters*.



misses badly
in the right tail

a curved line
implies a
wrong dist'n
family



Fitting with GoF tests & Q-Q plots

Chi-square test

Features:

- A formal comparison of a histogram or line graph with the fitted density or mass function
- Sensitive to how we group the data.

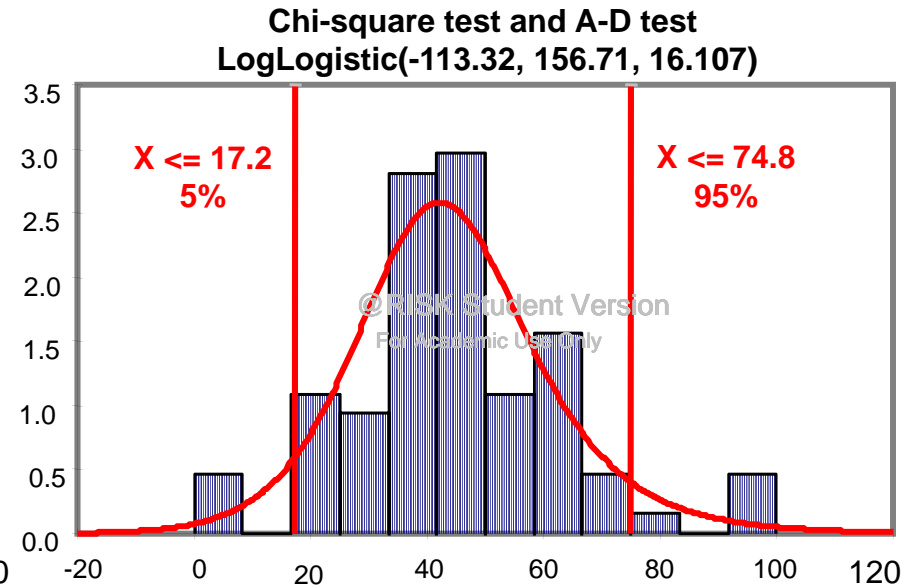
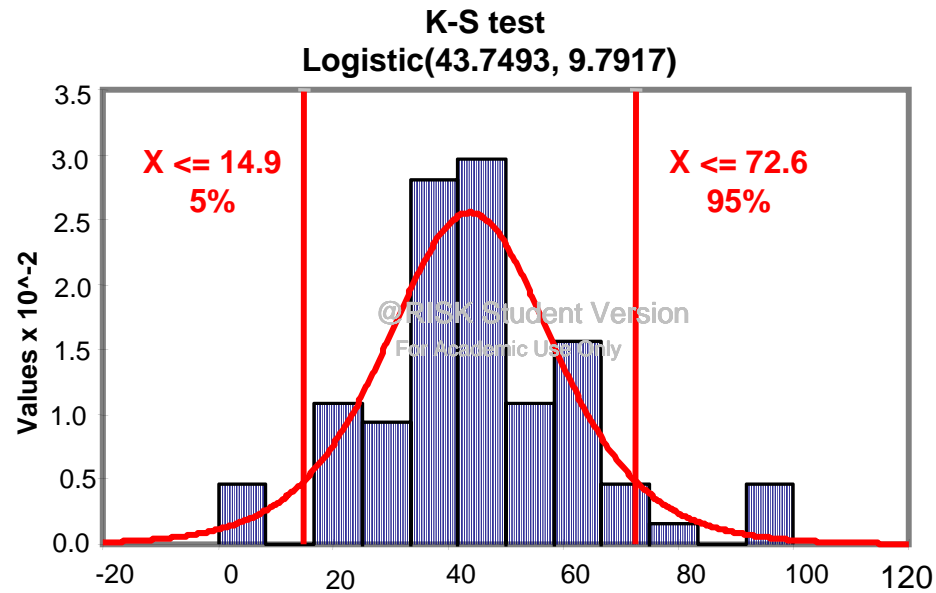
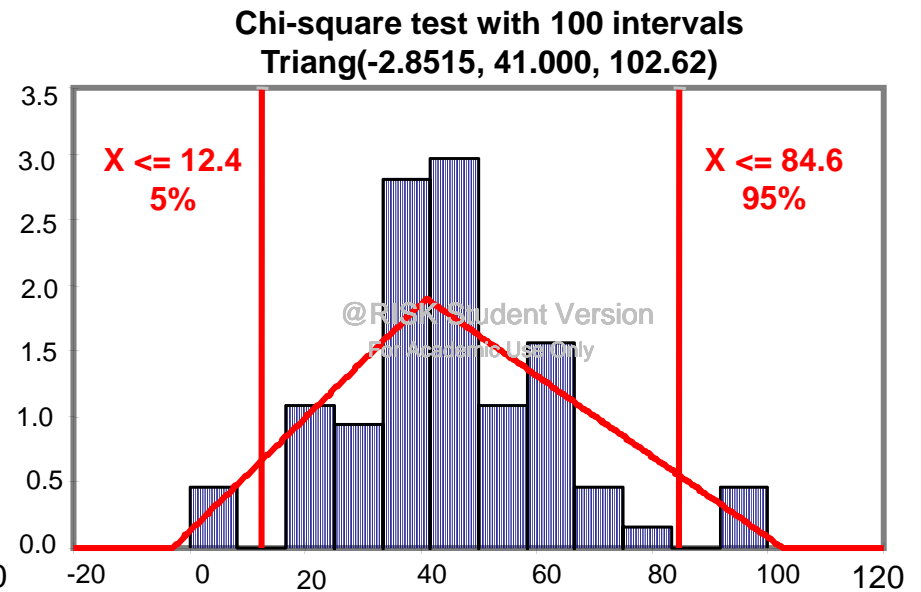
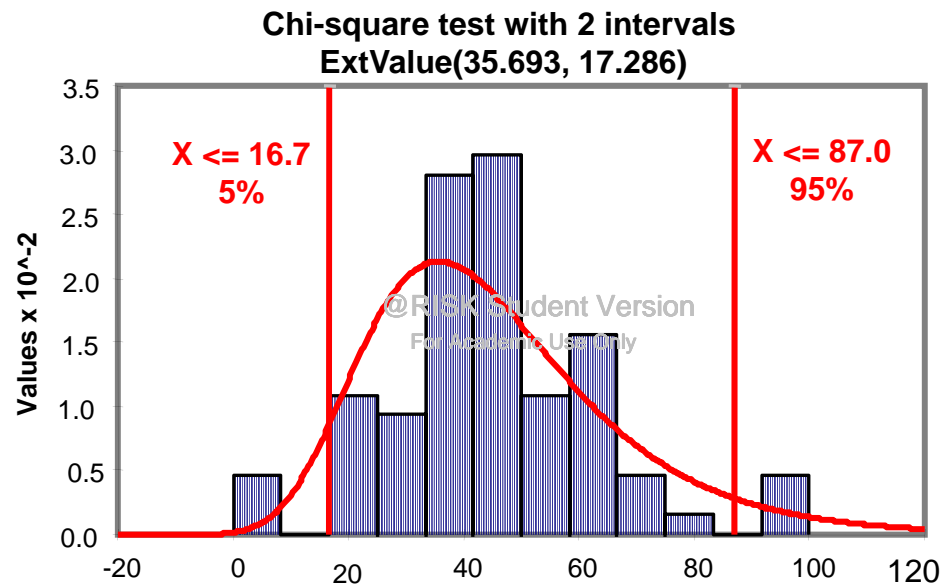
K-S and A-D tests

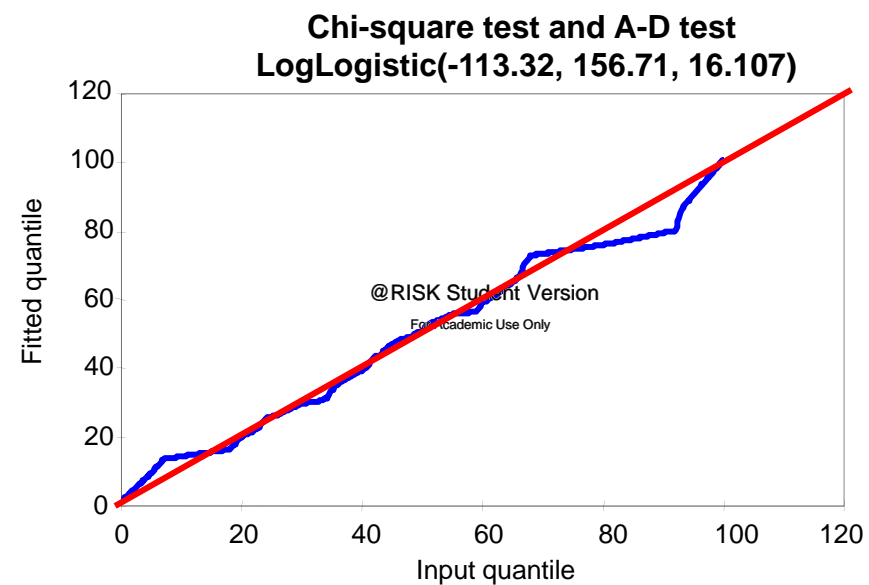
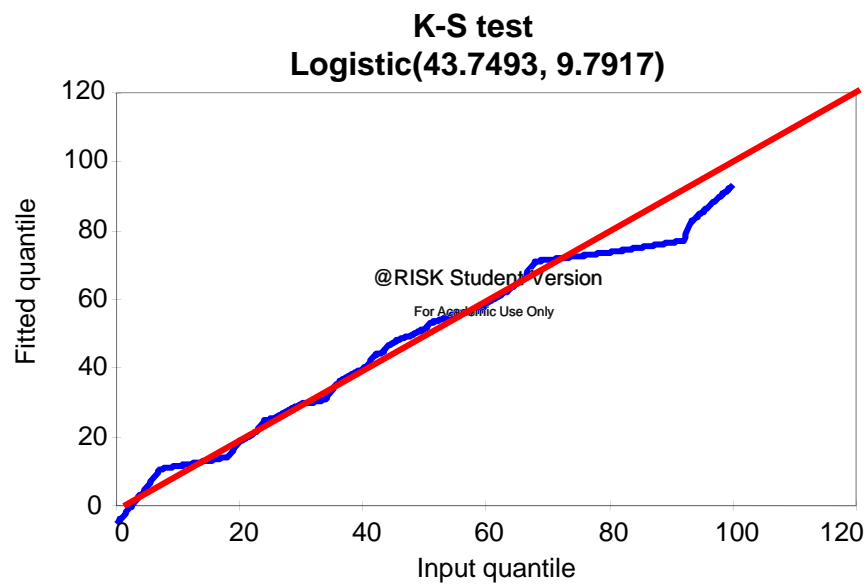
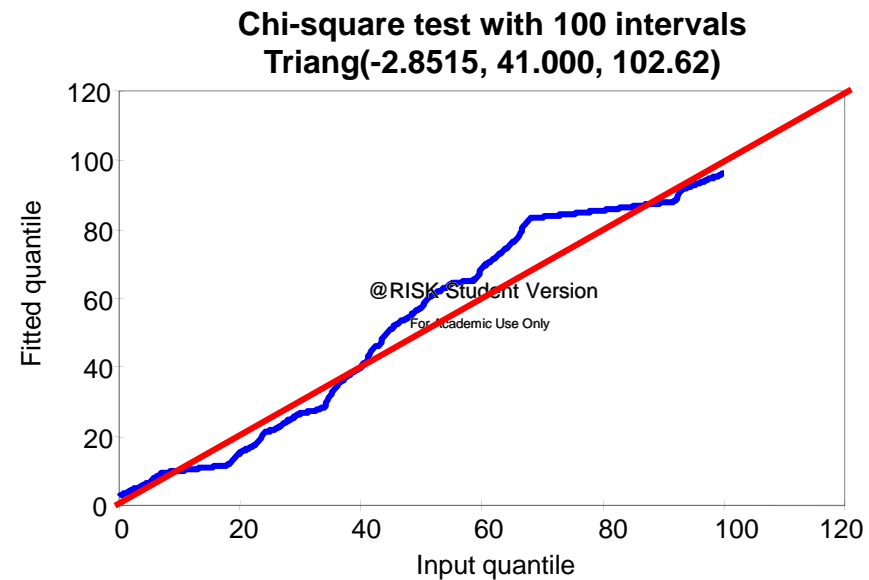
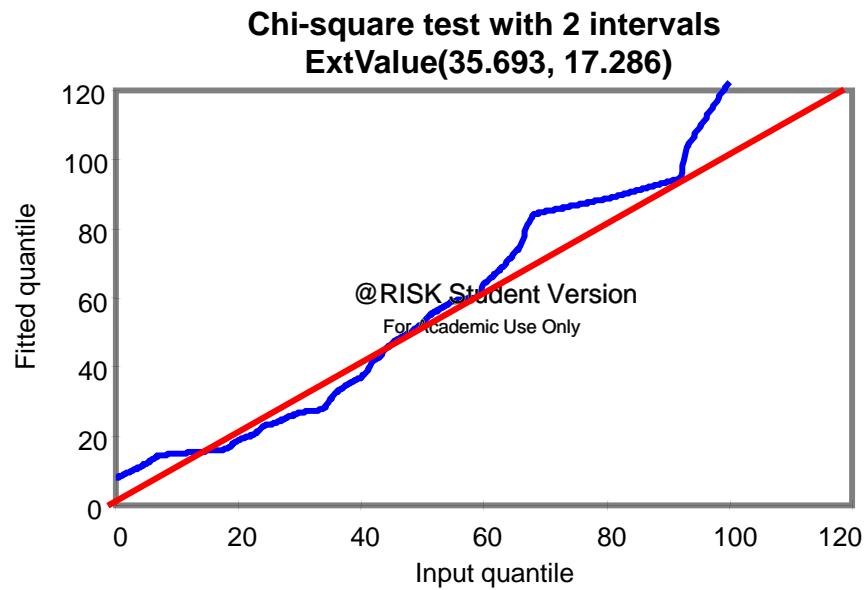
Features:

- Comparison of an empirical distribution function with the distribution function of the hypothesized distribution.
- Does not depend on the grouping of data.
- A-D detects discrepancies in the tails and has higher power than K-S test

- Beware of goodness-of-fit tests because they are unlikely to reject *any* distribution when you have little data, and are likely to reject *every* distribution when you have lots of data.
- **Avoid histogram-based summary measures, if possible, when asking the software for its recommendation!**







Using the Data Itself

- When might we want to use the data itself?
 - When no standard distribution fits well
 - When we have no justification for a standard distribution
 - When there is too little data to distinguish between standard distributions
- Example: Fit data 2.1, 5.7, 3.4, 8.1 (min)



Empirical Distribution

- Each data point is equally likely to be *resampled*.
- In Simio:

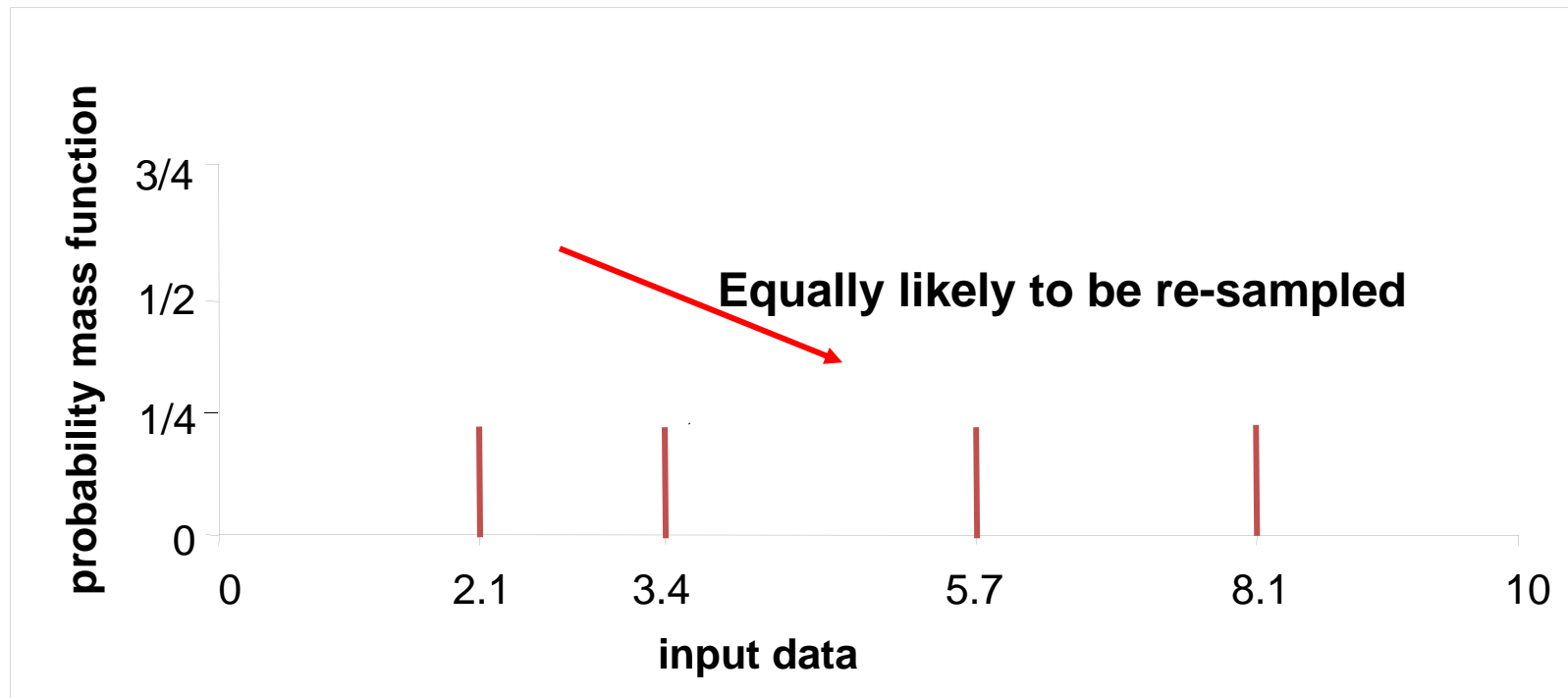
Random.Discrete(2.1, .25, 3.4, .5, 5.7, .75, 8.1, 1)

Random.Discrete ($X_1, 1/n, X_2, 2/n, \dots, X_n, 1$)



Empirical Distribution

Objective Fit an input model to data 2.1, 5.7, 3.4, 8.1 via **empirical distribution function**.



Pros and Cons

- As the sample size n goes to infinity, the empirical distribution converges to “the truth.”
- No assumed distribution need be selected.
- Only the values we saw can appear again; no tails and nothing in the gaps.



Interpolated Empirical

- To fill in gaps, we linearly interpolate between the **sorted** data points.

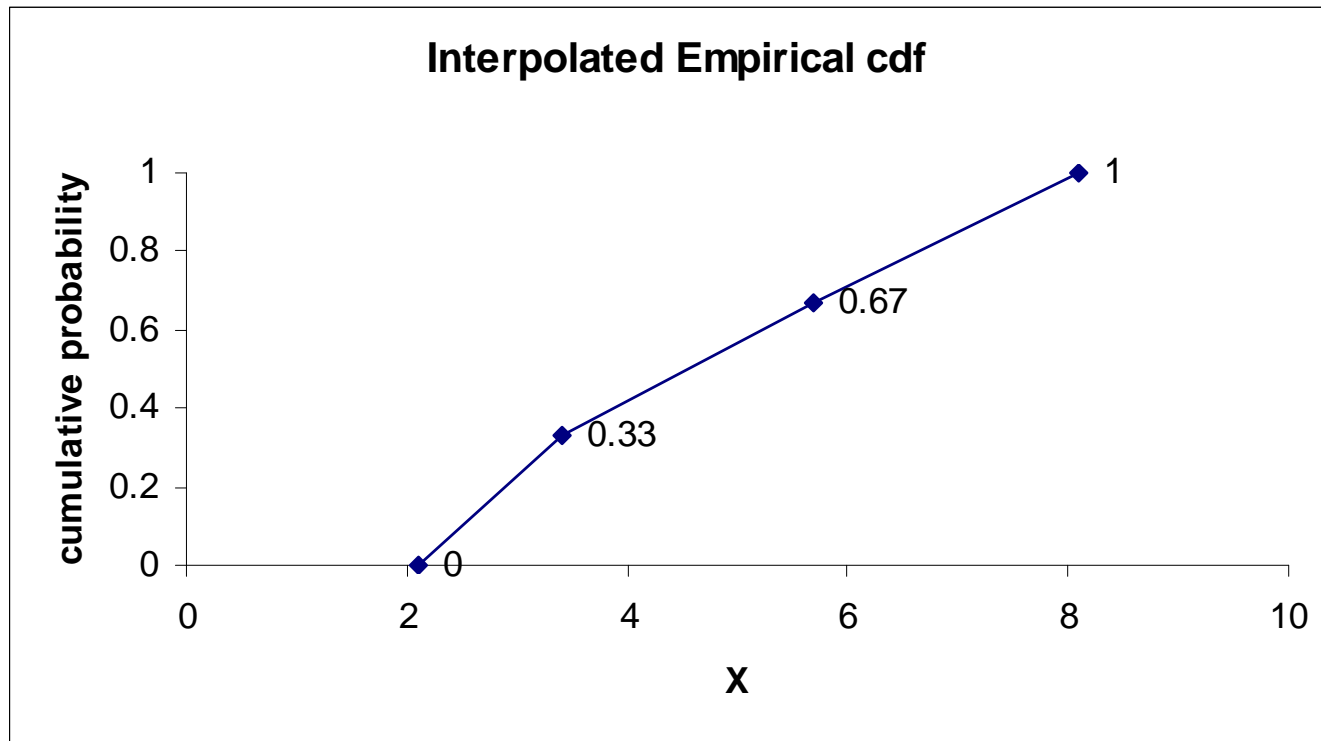
- In Simio:

Random.Continuous(2.1, 0, 3.4, .33, 5.7, .67, 8.1, 1)

CONT(X_1 , 0, X_2 , $1/(n-1)$, X_3 , $2/(n-1)$, ..., X_n , 1)



Example



Input Modeling without Data

- We have to use anything we can find...
 - Engineering data, standards and ratings can provide central values
 - **Expert opinion**
 - Physical or conventional limitations can provide bounds
 - Physical basis of the process can suggest appropriate distribution families



Breakpoints Method

- Useful for modeling quantities with a large number of possible outcomes, like quarterly sales volume or aggregate number of overtime hours.
- Minimum information needed: **smallest and largest possible values.**
 - Ex: sales of XYZ-123 will be no less than 1000 units, but no more than 5000 units.



– Random.Uniform(1000,5000)

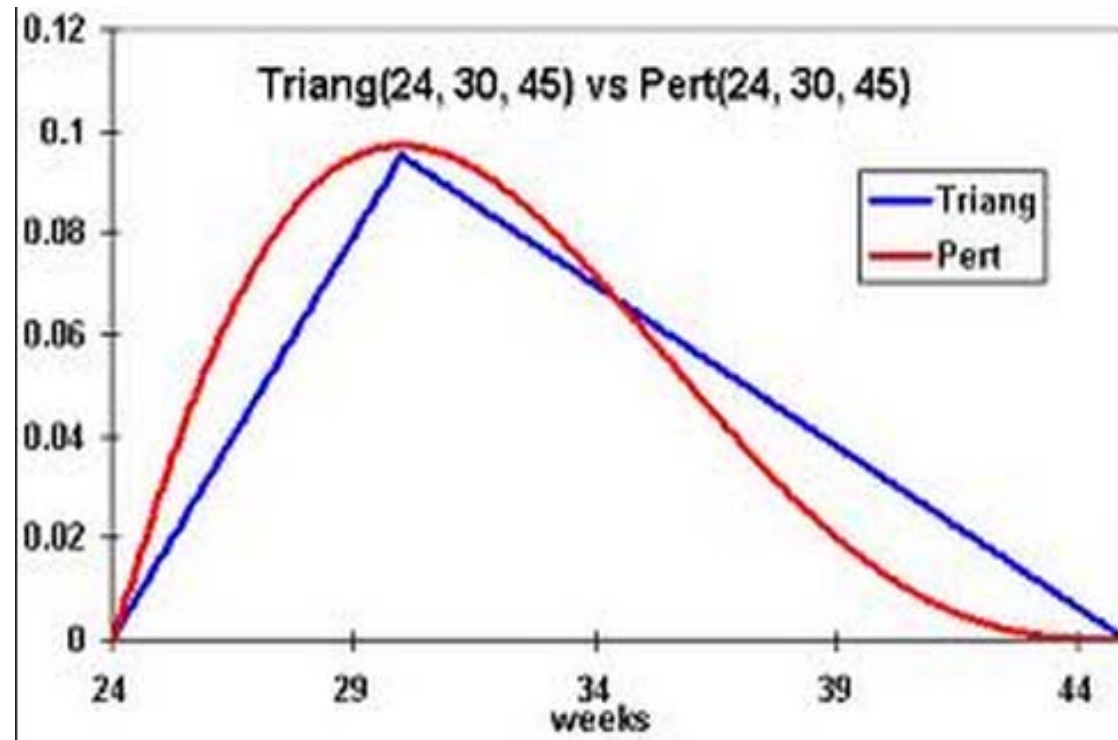
– Comments:

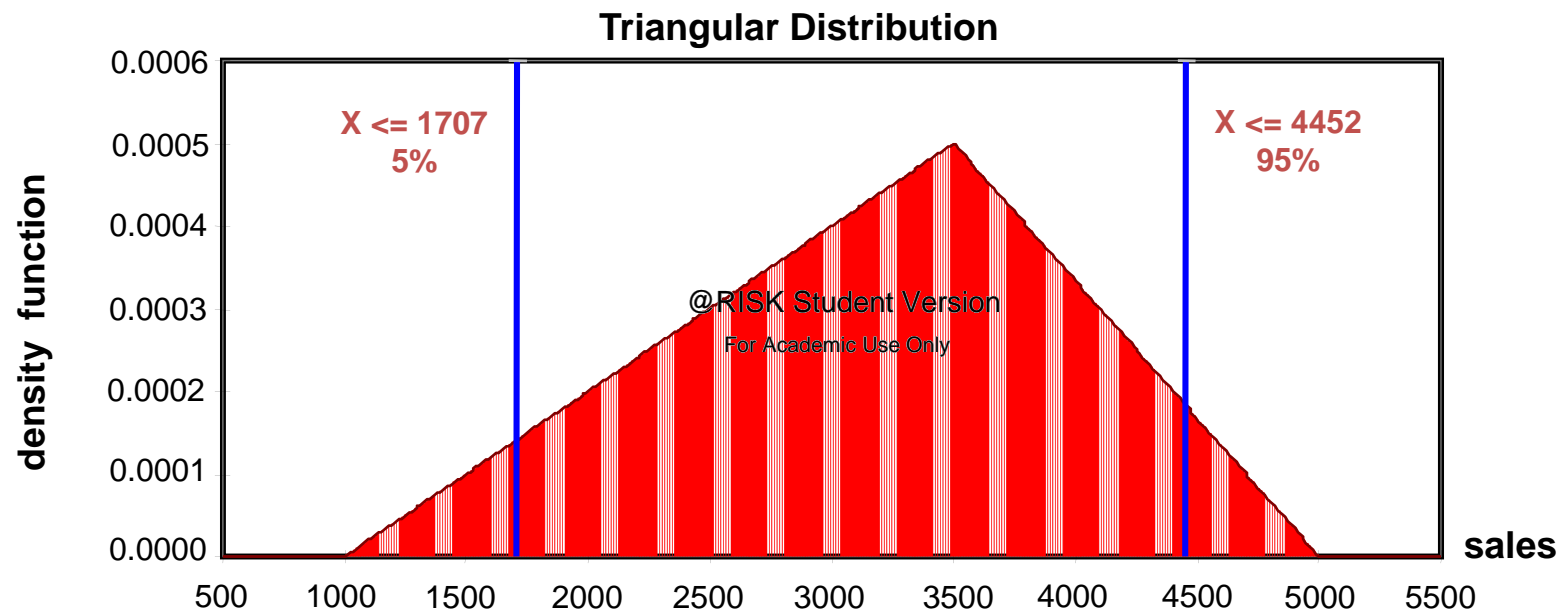
- This is typically a poor model since the probability is spread out evenly from low to high. Thus, the extremes are just as likely as the middle values.
- However, if you want to be conservative (maximum uncertainty) or you cannot justify any additional information, then this model may be reasonable.



- Better information: **smallest and largest possible values, and most likely value.**
 - Ex: sales of XYZ-123 will be no less than 1000 units, no more than 5000 units, and is most likely to be 3500 units.
 - **Random.Triangular(1000,3500,5000)**
 - **Random.Pert(1000,3500,5000)**
 - Pert is a beta distribution. Be careful with betas; plot the resulting distribution.







- Best information: **smallest and largest possible values plus 1-3 breakpoints** (values and a percentage chance of being less than that value).

– Ex: sales of XYZ-123 will be between 1000 and 5000, and...

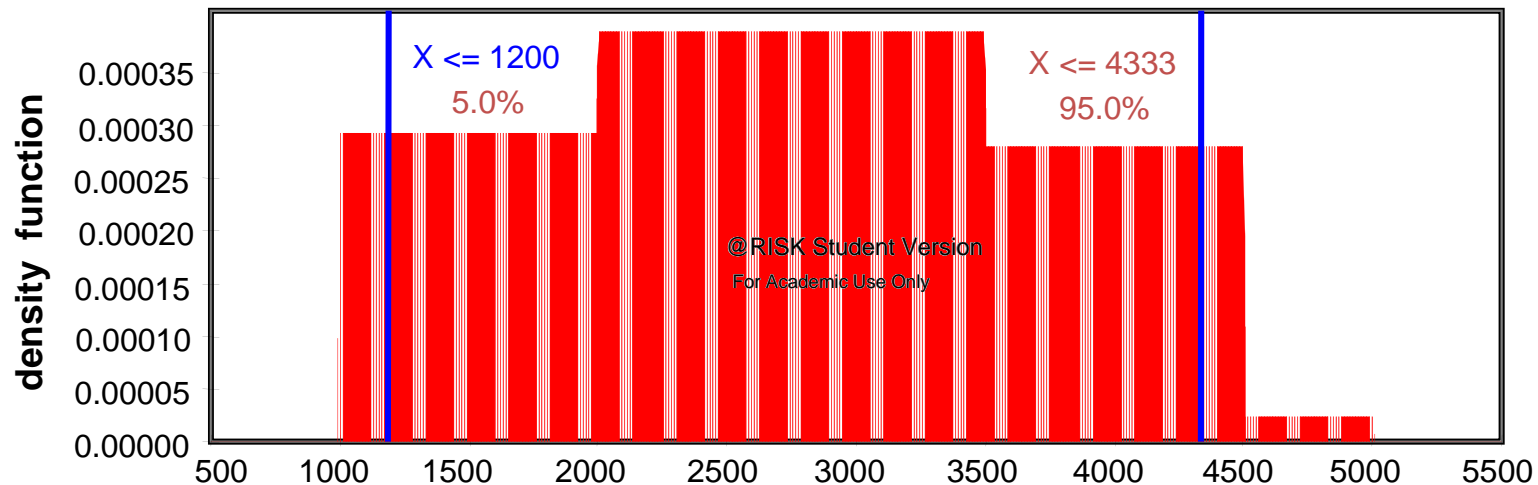
<u>sales</u>	<u>chance of being \leq sales</u>
2000	25%
3500	75%
4500	99%



– **Random.Continuous(1000, 0, 2000, .25, 3500, .75, 4500, .99, 5000, 1)**

– **Comments:**

- Use only as many breakpoints as you can confidently get. Three is usually the maximum if no data are available.
- Try to get breakpoints near the extremes if possible, since the extremes are often not realistic.
- Sometimes it is easier to get people to give the chance of exceeding a value, rather than being less than a value.



Mean & Variability Method

- Useful for modeling quantities with a large number of possible outcomes, like quarterly sales volume or aggregate number of overtime hours.
- Also useful for modeling the variability in percentage changes.



- Minimum data: mean value and an average percentage variation around that mean.
 - Example: Last year we sold 10,000 units of ABC-000. This year we expect a 15% increase, with a typical swing of 5% above or below that value.
 - Comments: **Random.Normal (11500, 0.05*11500)**
 - May need to round this value.
 - If the % swing > 33%, then negative values are possible.



- Better data: mean value, an average percentage variation around that mean and upper and lower limits.
 - Ex: Last year we sold 10,000 units of ABC-000. This year we expect a 15% increase, with a typical swing of 5% above or below that value. But we won't sell less than 8000 units, or more than 15,000 under any conditions.

$\text{Math.Min}(15000, \text{Math.Max}(\text{NORM}(11500, 0.05*11500), 8000))$



Don't Forget Correlations

- If data exists, calculate sample correlation.
- If not,
 - the direction (positive or negative) is typically easy to specify.
 - One way to scale the value is to ask what percentage (say P%) of the time the two inputs move together.
 - Then estimate the correlation by $\rho = P/100$ with the appropriate sign.



Sensitivity Analysis

- *Sensitivity analysis* (varying the parameters of the input model) is especially important when the model is not based on data. Pay special attention to the standard deviation and bounds or limits.



- Sensitivity analysis examples:
 - **Random.Continuous(1000, 0, 2000, .25, 3500, .75, 4500, .99, 5000, 1)**
Might want to vary the 1000 and/or 5000 in and out to see if there is a marked change in the output results.
This is also very important for **Triang**.
 - **Random.Normal(11500, 0.05*11500)**
Might want to vary the standard deviation up and down to see if there is a marked change in the output results.
- Concentrate sensitivity analysis on those inputs to which the outputs are most sensitive.



Concluding Remarks

- Use input models to represent uncertainty in simulation
- The particular input model chosen matters!
- Selection of the an input model is **not an exact science - no right answer**, but the issues to consider are
 - theoretical vs. empirical data
 - physical basis of the distribution
 - assessment of the goodness of a fit
 - independence vs. dependence
 - stationarity vs. non-stationarity (to be covered later)
- Assess the sensitivity of simulation output results to input models chosen
- Use expert opinion whenever you can

You are smarter than the software!

