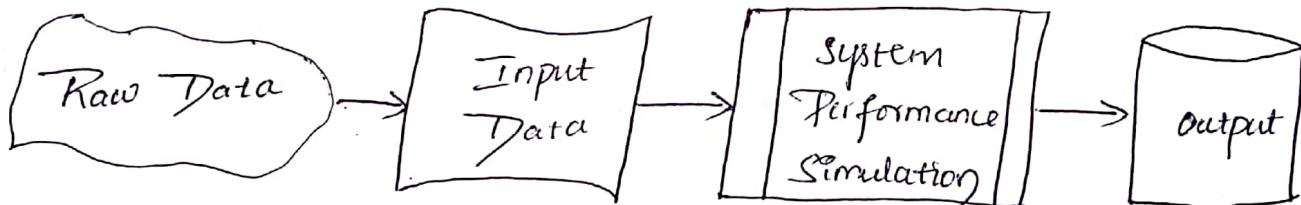


Write a short note on the following topics

1. Data Collection

- ⇒① Data Collection is one of the biggest task in solving real problem.
- ⇒② It is one of the most important and difficult problems in simulation.
- ⇒③ Even if when data are available, they have rarely been recorded in form that is directly useful for simulation input modelling.
- ⇒④ "Garbage-in, Garbage-out" is a basic concept in Computer Science and it applies equally in area of discrete system simulation.
- ⇒⑤ Many are fooled by a pile of computer output or a sophisticated animation, as if there were the absolute truth.



- ⇒⑥ Even when Model Structure is valid simulation results can be misleading, if the input data is
 - Inaccurately collected
 - Inappropriately analyzed
 - Not representative of the environment.

* Suggestion that enhance data Collection:

- ① Plan ahead : Pre observing session.
- ② Analyze data being collected.
- ③ Check for variable relationship.
- ④ Check for auto correction

② Identifying the Distribution with Data?

⇒ A frequency distribution or histogram is useful in identifying the shape of a distribution.

① A histogram is constructed as follows:

1. Divide the range of the data into intervals [usually of equal width; however, unequal width may be used if tight of frequencies are adjusted]

2. Label the horizontal axis to conform to the intervals stated.

3. Determine the frequencies of occurrence within each interval.

4. Label the vertical axis so that total occurrence can be plotted for each interval.

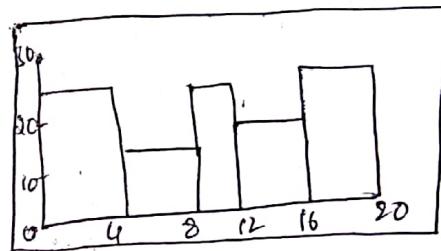
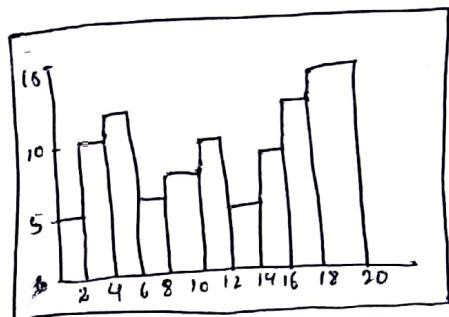
5. Plot the frequencies so that total occurrence on the vertical axis.

② The number of class intervals depends on

→ The use of observation.

→ The dispersion of data.

→ Suggested by data no of intervals.



③ for continuous data:-

→ Corresponds to probability density function of theoretical distribution.

④ for discrete data:-

→ Corresponds to probability mass function.

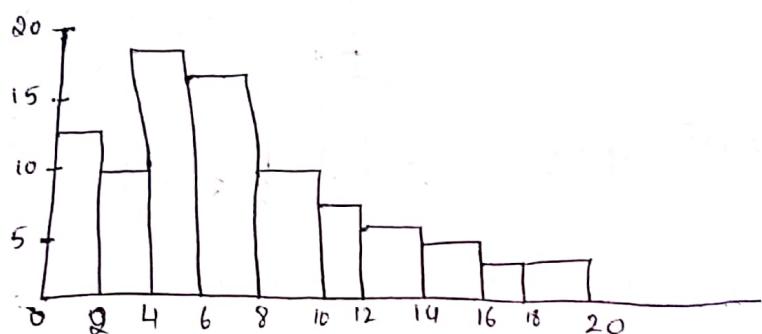
⑤ if few data points are available.

→ Combines adjacent cell to eliminate the ragged appearance of histogram.

⑥ figure above shows same data with different Interval size.

Eg :- The No. vehicle arriving at the northeast corner of intersection in 5 min b/w 7 AM & 7:05 AM was monitored for 5 days over 20 intact period. Table shows resulting data. The 1st entry in table indicate 12 to 15 min period during which 0 vehicle arrived & so on. - Red, 10 period during which 1 vehicle arrived & so on. The No. of automobiles in a discrete fashion. The No. of sample data, & the histogram may have a cell of each possible value in range data.

Arrival Per period	Frequency
0	12
1	10
2	19
3	17
4	10
5	8
6	7
7	5
8	5
9	3
10	3
11	1



③ Parameter Estimation :-

⇒ ① After a family of distribution has been selected, the next step is to estimate Parameter of distribution.

② If Observation in a Sample of size n are x_1, x_2, \dots, x_n . The sample Mean & sample Variance are :

• The Sample Mean is $\bar{x} = \frac{\sum_{p=1}^n x_p}{n}$

• The Sample Variance is $s^2 = \frac{\sum_{p=1}^n x_p^2 - \bar{x}_n^2}{n-1}$

• If data are discrete in a frequency distribution, Then we can re-write the equation as :

$$\bar{x} = \frac{\sum_{j=1}^k f_j x_j}{n}$$

and

$$s^2 = \frac{\sum_{j=1}^k f_j x_j^2 - n \bar{x}^2}{n-1}$$

where k is No. of distinct value of x
 f_j is observed frequency of value x_j of x .

③ If The Data are Continuous we "discretize" Then estimate The Mean $\bar{x} \approx \sum_j f_j m_j$ and Variance $s^2 \approx \frac{\sum_j f_j m_j^2 - n\bar{x}^2}{n-1}$

④ A parameter is an unknown Constant, but an estimator in Statistic.

Eg: Vehicle Arrival Example.

Table in Histogram of Vehicle Example can be Analysis to obtain

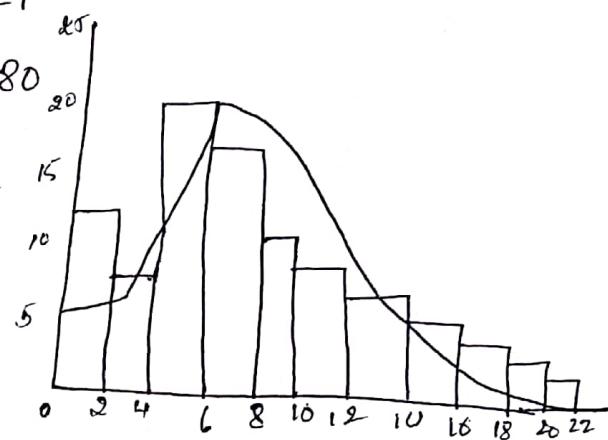
$$n_k = 100, f_p = 12, x_1 = 0, x_2 = 1$$

$$\sum_{j=1}^k f_j x_j = 364 \text{ and } \sum_{j=1}^k f_j x_j^2 = 2080$$

* The sample Mean and Variance are

$$\bar{x} = \frac{364}{100} = 3.64$$

$$s^2 = \frac{2080 - 100 \times (3.64)^2}{99} = 7.63$$



⑤ The Histogram Suggests x to have poisson Distribution.

- However, Most Sample is Not equal to Sample Variance.

Theoretically: Poisson With Parameter $\lambda \Rightarrow \mu = \sigma^2 = \lambda$

Reason: Each estimator is random Variable, its Not Perfect.

⑥ Goodness of fit Tests:-

- Provides helpful guidance for evaluating Suitability of a Input model
- There is no single Correct distribution in real application exists.
- If Very little data are available, it is unlikely to reject any candidate distribution.
- If a lot of data are available, it is likely to reject all candidate distribution.
- Conduct hypothesis testing on Input data distribution using:
 - * Chi-Square test
 - . Kolmogorov-Smirnov test.

* Chi-square Test:

One procedure for testing hypotheses of random size n of random variable x follows specific distribution form in chi-square. The test Procedure begins by arranging n observation into set of k class interval or chi-statistic is given by $\chi^2_o = \sum_{i=1}^k \frac{(O_p - E_p)^2}{E_p}$
 Where $O_p \rightarrow$ Observed frequency
 $E_p \rightarrow$ Expected frequency.

The Expected frequency for each class interval is computed as $E_p = n p_p$

The Hypotheses are following:-

H_0 : The random variable x , Conforms The distributional assumption with parameter given by Parameter Estimate.

H_1 : The random variable x does not Conform.

① Each value of random variable Should be class interval, unless Combining is necessary and $p = p(x_p) = p(x=x_p)$

② for Continuous case with assumed $f(x)$ or assumed pdf $f(x)$, p_p can be Computed by,

$$p_p = \int_{a_{p-1}}^{a_p} f(x) dx = f(a_p) - f(a_{p-1})$$

Table Below are Made to aid in determining The No of class Interval for Continuous Data.

Sample size n	Number of class Intervals *
20	Do Not use chi-square
50	5 to 10
100	10 to 20
> 100	$\sqrt{n} + 0.5$

② Chi-Square Test with equal Probabilities.

- If Continous assumption is being tested, class Interval are equal in probability rather than equal in width of Interval Should be used.
- unfortunately, There is yet no Method for determining The probability associated with each Interval that Maximize the Power of that of given size

$$E_p = n \quad P_p \geq 5$$

③ Kolmogorov-Smirnov Goodness of fit Tests

This test is particularly useful when sample size are small & when no parameters have been estimated from data.

Eg Suppose 50 Interval interval Times are Collected Over the following 100 min Interval.

0.44, 0.53, 0.04, 2.14, 2.00, 0.30, 2.54, 0.52, 2.02, 1.89, 1.53, 0.21, 2.80, 0.04
 1.35, 8.32, 2.34, 1.95, 0.1, 1.42, 0.46, 0.09, 1.09, 0.75, 5.55, 3.93, 1.01, 2.26
 2.88, 0.67, 1.12, 0.26, 4.57, 5.37, 0.12, 3.19, 1.63, 1.46, 1.09, 2.06, 0.85, 0.83,
 2.44, 2.11, 3.15, 2.9, 6.58, 0.64

H_0 : The Interarrival Times are Exponentially distributed.

H_1 : The Interarrival Times are not exponentially distributed.

The Data were Collected Over The Interval 0 to $T = 100$ min can be shown $\{T_1, T_2, \dots\}$ if time exponential, The arrival times distributed on Interval $(0, T)$. The arrival times $T_1, T_1+T_2, T_1+T_2+T_3, \dots + T_{50}$ are obtained by adding times on a $(0, 1)$ Interval, The points will be $[T_1/T, (T_1+T_2)/T, \dots, (T_1+T_2+\dots+T_{50})/T]$.

0.0044	0.0097	0.301	0.0575	0.0775	0.0805	0.1059	0.1111	0.1313	0.1502
0.1655	0.1676	0.1956	0.1960	0.2095	0.2927	0.3151	0.3576	0.3866	0.3828
0.3553	0.3561	0.3670	0.3746	0.4300	0.4694	0.4796	0.5027	0.5315	0.5382
0.7553	0.7636	0.7880	0.7982	0.8206	0.8417	0.873	0.9022	0.968	0.944

following the procedure in Kolmogorov-Smirnov test studied earlier,
Procedure & $O^+ = 0.1054$ & $D = 0.0080$. Therefore This statistic is
 $D = \max(O^+, O^-) = \max(0.1054, 0.0080) = 0.1054$.

The Critical Value of D for level of significance of $\alpha = 0.05$ & $n = 50$
is $D_{0.005} = 1.36 / \sqrt{n} = 0.1923$.

But $D = 0.1054$, so The hypothesis that the Interarrival Times are
Exponentially distributed Cannot be rejected.

④ P-values & Best-fit :-

P value for testing statistics. The significance level at which one
would reject H_0 for the given statistics.

A Measure of fit, the larger the better,

- Large - p-value : good fit.
- Small p-value : poor fit.

⑤ Fitting non stationary Poisson process.

⇒ fitting a NSPP data is difficult, possible approaches:

① fit a very flexible model with lots of parameters.

② Approximate constant arrival rate over some basic Interval time,
but vary it from time Interval.

③ Suppose we need to model arrival rate time $[0, T]$, Our approach
is most appropriate when we can:

④ Observe the time Period repeatedly.

⑤ Count arrival/round arrival times.

⑥ Divide the Time Period into k equal Interval of length $\Delta t = T/k$

⑦ Over n Period of observation let c_{ij} be no of arrivals
during the Interval on j^{th} Period.

The estimated arrival rate during i^{th} time period ($i-1 \Delta t \leq s \leq i \Delta t$) is:

$$\lambda(t) = \frac{1}{n \Delta t} \sum_{j=1}^n C_{ij}$$

Eg: Divide 10 hours business day [8am, 6PM] into Interval $k=20$
whose length $\Delta t = 1/2$, & observe over $n=3$ days.

Time Period	No. of Arrival			Estimate Arrival Rate (arrival/hr)
	Days	Days	Days	
8:00 - 8:30	12	14	10	24
8:30 - 9:00	23	26	32	54
9:00 - 9:30	27	18	32	52
9:30 - 10:00	20	13	12	30

for instances,
 $1/3 (0.5) * (23 + 26 + 32) = 54$ arrival/hour

⑥ Selecting Input Models Without Data.

⇒ If data is not available, some possible source to obtain information about process are:

- ① Engineering data: often product or process has performance rating provided by manufacturer or company rules specific time or production standard.
- ② Expert Option: Pessimistic and Most-likely times, and they may know the variability as well.
- ③ Physical or Conventional limitations: physical limit on performance, limit or bounds that narrow the range of input process.

Eg: Production Planning Simulation

Input of sales volume of Varibes product is required, salesperson of Product xyz says that:

Interval [Sales]	PDF	cumulative frequency
1 $1000 \leq x \leq 2000$	0.1	0.10
2 $2000 \leq x \leq 2500$	0.65	0.75
3 $2500 \leq x \leq 4500$	0.24	0.99
4 $4500 \leq x \leq 5000$	0.01	1.00

- No fewer than 1000 units and no more than 5000 units will be sold.
- Given Experiences, says There is 90% chance of selling More than 2000 units, a 25% chance of selling more than 2500 unit & 1% chance of selling More than 4500 units.
- Translating these information into a Cumulative probability or being less than or Equal to Those goals Simulation Input.

⑦ Multivariate and Time Series Input Models.

→ The random Variable discussed until now were Considered to be Independent of any other Variables Within The Context of The problem.

- However , variables may be related.
- If they appear as I/p, The relationship Should Be investigated and taken into consideration.

⑧ Multivariate Input Models :

- fixed, finite No of random Variables $x_1, x_2 \dots x_n$
- for Eg. lead time and Annual demand for an Inventory Model.
- An Increase in demand results in Lead time Increase, Hence Variable are dependent.

⑨ Time Series Input Models:

- Infinit Sequence of random Variables Eg x_1, x_2, x_3, \dots
- for Eg:- Time b/w arrival of orders to buy & sell stocks.
- Buy and Sell Orders tends to occur in bursts hence, Times b/w arrivals are dependent.

⑧ Types of Simulation with respect to Output Analyses.

- ⇒ • Terminating Versus Non-Terminating Simulation.
- Terminating Simulation:
 - Runs for some duration of Time T_E , where E is a specified event that stops the simulation.
 - Starts at time 0 under well-specified initial Condition.
 - Bank eg: Opens at 8:30 AM (time 0) with no customer present at 8 of the 11 teller working & closes at 4:30 PM (Time $T_E = 480$)
 - The Simulation Analyst chooses to consider it a terminating system because the object of interest is one day's operation.

⑨ Stochastic Nature of Output Data.

⇒ Model output consists of one or more random variable (r, v)
Because the model is an I/O transformation and input variables are r.v's.

M/G/1 queuing eg:

→ Poisson arrival rate = 0.1 per min

Service time $\sim N(m=9.5, s=1.75)$

→ System Performance: long-run Mean queue length $L_Q(t)$.

→ Suppose we run a single simulation for a total of 5000 min.

• Divide the Time Interval $[0, 5000]$ into 5 equal subinterval

of 1000.

Average No of Customer in queue from time $(j-1)1000$ to $j(1000)$

• M/G/1 queuing

→ Batched Average queue length for 3 Independent replication

Batching Interval [Min]	Batch, j	Replication		
		1, y_{1j}	2, y_{2j}	3, y_{3j}
0, 1000	1	3.61	2.91	7.67
1000, 2000	2	3.21	9.00	19.53
2000, 3000	3	9.18	16.15	20.36
3000, 4000	4	6.92	24.53	8.11
4000, 5000	5	9.89	25.19	19.62
0, 6000		3.75	15.56	13.66

- Inherent Variability in Stochastic Simulation both within a single replication and across different Application
- The Average across 3 replications, can be regarded as Independent Observations, but Average within a replication $y_1 \dots y_{15}$ are Not.

⑩ Absolute Measure of Performance and their estimation.

- ⇒ • Consider The estimation of a performance Parameter, $q(\alpha, f)$ of Simulated system.
- Discrete Time Data: $[y_1, y_2 \dots y_n]$ with Ordinary Mean: q
- Continuous-time data: $\{y(t), 0 \leq t \leq T_E\}$ with time-weighted mean: f
- Point estimation of discrete time data
 - The Point Estimator: $\hat{\theta} = \frac{1}{n} \sum_{i=1}^n y_i$
 - It is unbiased if its Expected value is θ , p.e if $E(\hat{\theta}) = \theta$.

Is biased if:

* Point Estimator

- Point estimator for continuous-time data.
→ The point estimator: $\hat{y} = \frac{1}{T_E} \int_0^{T_E} y(t) dt$
→ Is biased in general where:
 - An unbiased or low-biass estimator is desired.
 - usually, system performance measure can be put into common framework by g or f :

Eg: The performance of days on which sales are lost through an out of stock situation let:

$$y(t) = \begin{cases} 1, & \text{if out of stock on day } t \\ 0, & \text{otherwise.} \end{cases}$$

- Performance Measure that does not fit: quantile or percentile.
- Estimating quantiles: The inverse of problem of estimating a proportion or probability. $\Pr\{Y \leq Q\} = P$.
- Consider a histogram of observed value y :
find Q such that 100% of the histogram is to the left of Q .

④ Confidence-Interval Estimation:

- To understand this fully, it is important to distinguish b/w Measures of error, risk, e.g. Confidence Interval versus Prediction Interval.
- Suppose the model is normal distribution with Mean μ , Variance s^2 (both unknown)
- Let y_{ij} be avg cycle time for parts produced on i^{th} replication in the simulation.

→ Avg cycle time will vary from day to day, but over the long-run
the average of the averages will be close to \bar{y} .

→ Sample Variance across R replication : $S^2 = \frac{1}{R-1} \sum_{r=1}^R (y_r - \bar{y})^2$

* Confidence - Interval Estimation

→ Confidence Interval (CI):

- A Measure of Error.
- Where y_r are normal distributed : $y_r \pm t_{\alpha/2, R-1} \frac{S}{\sqrt{R}}$
- We cannot know for certain how far is \bar{y} but CI attempts to bound that error.
- A CI, 95% tells us how we can trust Interval actually bound that error.
- The More replication we Make, less error occurs.

→ Prediction Interval (PI):

- A Measure of Error.
- ~~where y_r~~ A good guess for Average cycle time on a particular day is our estimator but it is unlikely to be exactly right.
- PI is designed to be wide enough to contain the Actual Average cycle time on any particular day with high probability.
- Normal - Theory Prediction Interval:

$$y = \bar{y} \pm t_{\alpha/2, R-1} \sqrt{\frac{S^2}{R} + \frac{1}{R}}$$

- The length of PI will not go to 0 as R increase because we can never simulate away risk.
- PI's limit is : $\theta \pm t_{\alpha/2, \infty} S$.

⑪ Output Analysis for terminating Simulation.

⇒ A Terminating Simulation: runs over simulated time interval $[0, T_E]$

A Common goal is estimate:

$$a) \Theta = E \left[\frac{1}{n} \sum_{s=1}^n y_s \right], \text{ for discrete output.}$$

$$b) \phi = E \left[\frac{1}{T_E} \int_0^{T_E} y(t) dt \right] \text{ for Continuous Obj } y(t), 0 \leq t \leq T_E$$

In General, ~~Inputs~~ Independent replication are used, each run using a different random no Stream & Independently chosen initial Conditions.

* Statistical Background.

- Important to distinguish Within replication data from Across-replication data.
- For eg, Simulation of Manufacturing System.

② Performance Measure of system: cycle time for Part & Work in process (wip).

③ Let y_{ij} be cycle time for j^{th} part production in i^{th} replication

iii) Across Replication Data are formed by summarizing Within Replication data.

- Overall Sample avg and interval replication Sample Avg are always unbiased estimators of Expected daily avg cycle time or daily avg wip

- Across replication data are independent and identically distributed, but within-replication data do not have these properties.

12) Output Analysis for Steady-time Simulations.

→ Consider a single run of Simulation Model to estimate a steady state long-run characteristics of system

- The Performance Measure

$$\bar{\theta} = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n y_i, \text{ for discrete Measure.}$$

$$\bar{\psi} = \lim_{T_E \rightarrow \infty} \frac{1}{T} \int_0^{T_E} y(t) dt, \text{ Continuous Measure.}$$

- The Sample size is design choice, with several consideration in mind.

⑨ Initialization Bias

Method to reduce the point-estimator bias caused by using artificial & unrealistic

- Intelligent Initialization

- Initialize the simulation in a state that is more representative of long run conditions.
- If the system exists, collect data on it and use these data to specify more nearly typical initial conditions.

⑩ Error Estimation

- If $\{y_1, \dots, y_n\}$ are not statistically independent then s^2/n is biased estimator of true variance
- Suppose point estimator $\hat{\theta}$ is Sample Mean

$$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$$

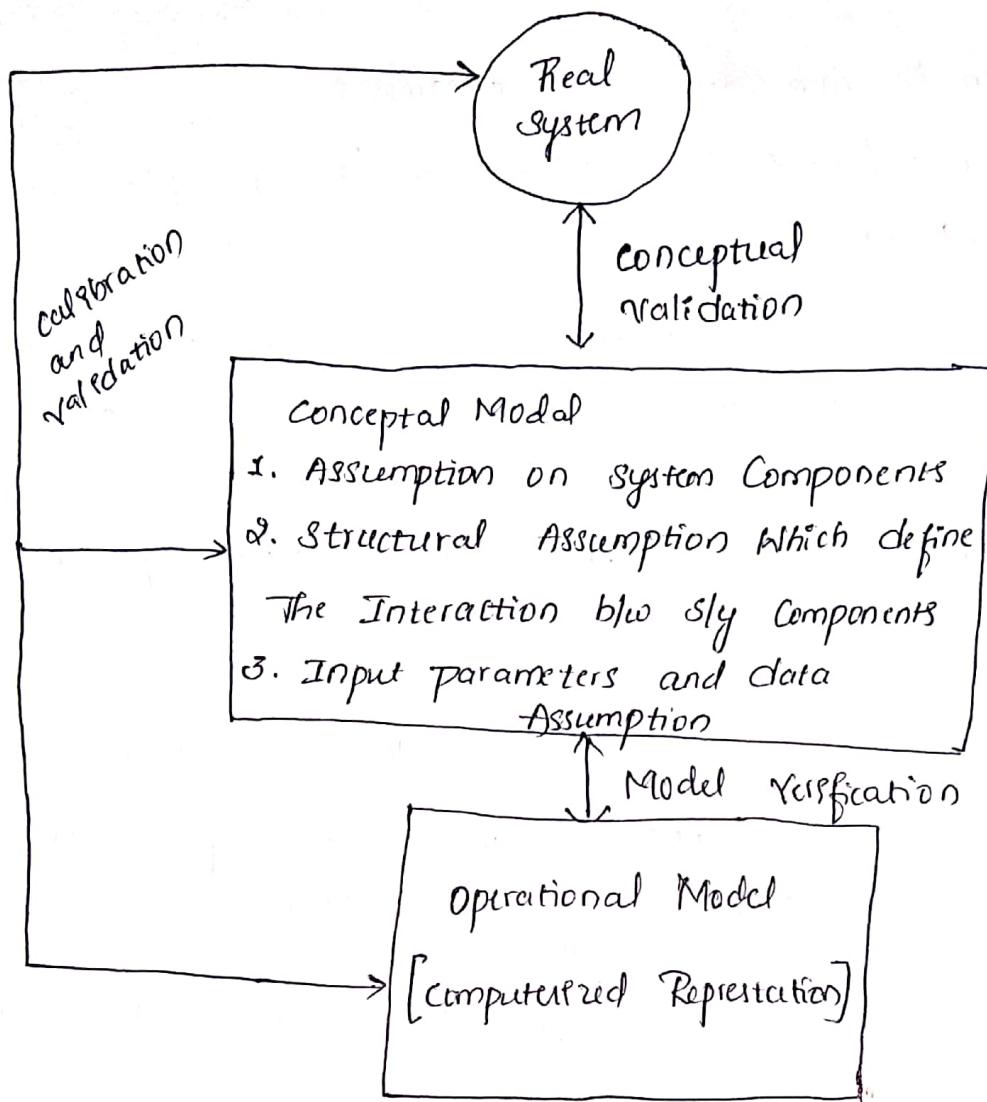
- Variance of \bar{y} is very hard to estimate.
- for SLE with steady state, produce an off process that is approximately convergent stationary.

⑫ Replication Method

- Approach make R replication initializing & deleting from each one that same way.
- Important to do through job of investigation Condition bias:
- Basic Raw o/p data $\{y_{ij}, i=1, R, j=1, \dots, n\}$ is derived
 - Individual Observation from within 'x'
 - Patch Mean from within replication 'r' of some no. of discrete-time Observation.

⑬ Model building, Verification and Validation.

- The first step in Model building consists of Observing real system and Interactions among its various components & collecting data on its behavior.
- Operators, technicians, repairs and Maintenance Personnel, Engineer and Managers under certain aspect of System which may be unfamiliar to Other.
- As Model development proceeds, new question may arise & Model development will return to this step of learning true System structure & behaviour.
- The second step in Model building is the Construction of a Conceptual Model \rightarrow a Collection of Assumption on Components & Structure of system, plan Hypothesis on The Value of Model into Parameters, Illustrate by following figure.
- ④ The Third step is Translation of the Operational Model into Computer recognizable form - The Computerized Model.



(ii) Verification of Simulation Model

⇒ The purpose of Model verification is to assure that Conceptual Model is reflected accurately in Computerized representation.

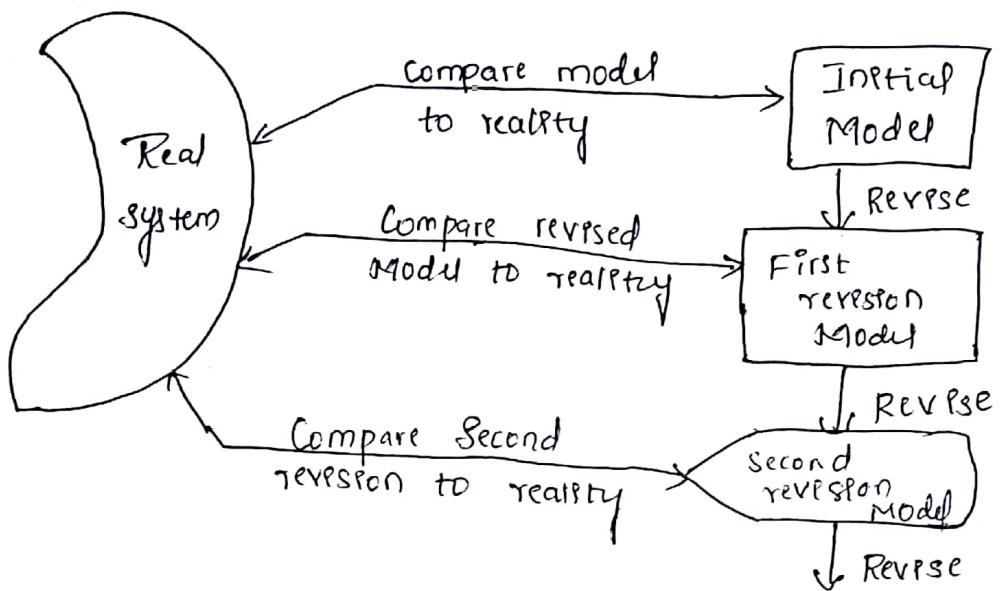
- Make someone else check the Model.
- Make a flow diagram that includes each logically possible action a system can take during an event.
- Closely examine the model o/p for reasonable under a variety of Input Parameter setting.
- Make the Operational Model as self documentary as possible.
- The Interactive run Controller (IRC) assist in debugging in the following ways:

- (e) Simulation can be monitored as it progresses.
- (f) Attention can be focused on a particular line of logic.
- (g) Values of selected Model Component can be Observed
- (h) The simulation can be temporarily suspended or paused, not only to view information but also to reassgn values or redirect entities.
- graphical Interfaces are recommended for accomplishing Verification and Validation.

(15) Calibration and Validation for Models, Optimization via simulation.

⇒ Verification and validation although are conceptually distinct, usually are conducted simultaneously by the Models.

- Validation is the overall process of comparing the Model to real system and its behaviour.
- Calibration is the interactive process of comparing the Model to real system, making adjustment to the model, comparing again & so on.
- Comparison of Model to real system.
- Subjective Tests - people who are knowledgeable about the system.
- Objective Tests - Requires data on real system's behaviour & the output of the Model.
- Validation is not an proposition - no model is ever totally representative of system under study.
- In addition, each revision of Model, as in fig. involves some cost, time and efforts.



As an aid in Validation Process, Naylor and finger formulated a 3 step approach which been followed:

- 1) Build a Model that has high face validity.
- 2) Validate Model Assumption.
- 3) Compare The Model Input-output transformation to Corresponding I/O transformation for the real sly.