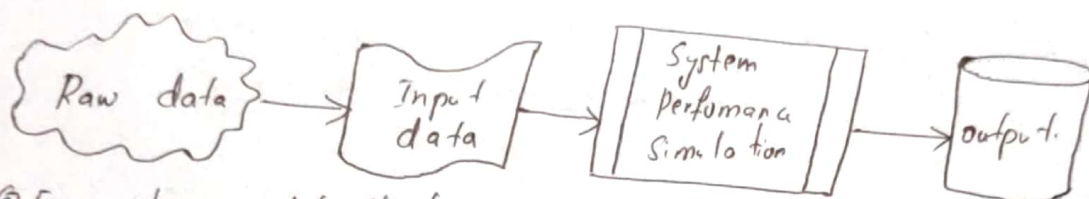


Write a note on following topic:

1) Data collection:

- ① Data collection is one of the biggest task in solving real world problem
- ② It is the most ^{important} and difficult problem 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 "GIGO" or 'Garbage-in, Garbage Out' is a basic component in computer science & 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: pretesting session.

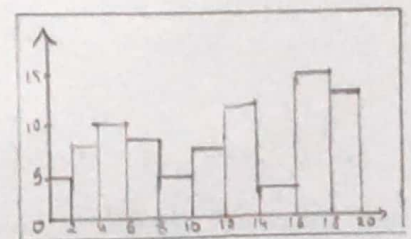
② Analyze data being collected ③ Check for variable relationship ④ check for autocorrelation.

2) Identifying the Distribution with Data:

→ ① A frequency distribution or histogram is useful in identifying the shape of a distribution.

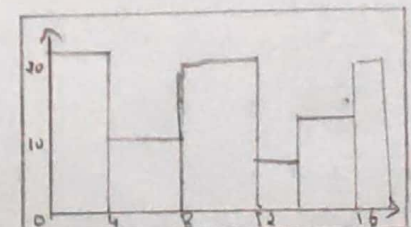
② A histogram is constructed as follows.

③ Divide the range of the data into intervals. Usually of equal width; however unequal width may be used if height of frequencies are adjusted.



④ Label the horizontal axis to conform to the intervals selected.

⑤ Determine the frequencies of occurrences within each interval.

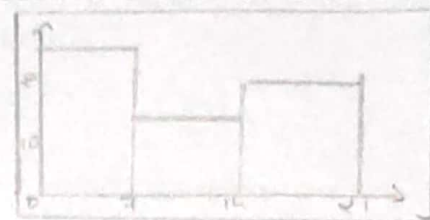


⑥ Label the vertical axis so that occurrences can be plotted for each interval.

⑦ plot the frequencies so that total occurrences on the vertical axis.

⑥ The number of class interval depends on

- The no. of observation
- The dispersion of data.
- suggested No. of intervals.



⑦ For continuous data:

- corresponds to probability density function of theoretical distribution.

⑧ For discrete data:

- corresponds to probability mass function.

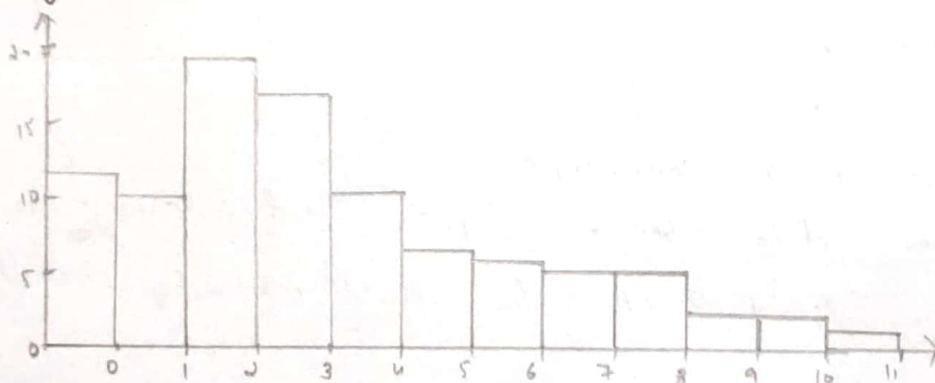
⑨ if few data points are available.

- combine adjacent cell to eliminate the legal appearance of diagram

Figure above shows some data with different interval size.

Ex: The No. of vehicles arriving at the northeast corner of intersection in 5 min b/w 7 am & 7.05 am. was monitored for 5 days over 10 week period. Table shows resulting data. The 1st entry in table indicates 1st 5 min period during which vehicles arrived & so on. The no. of automobiles is a discrete variable & there are sample data, so the histogram may have a cell for each possible value in range. The resulting histogram is shown below.

Arrival per period	Frequency
0	12
1	10
2	17
3	17
4	10
5	7
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:

① The sample mean is $\bar{x} = \frac{\sum_{i=1}^n x_i}{n}$

② The sample variance is $s^2 = \frac{\sum_{i=1}^n x_i^2}{n} - n\bar{x}^2$

③ If data are discrete in a frequency distribution, then we can re-write the equation as

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

where k is no. of distinct value of x and f_j is observed for value x_j of x .

④ If the data are continuous we "discretize" then estimate the mean $\bar{x} = \frac{\sum_{j=1}^p f_j m_j}{n}$

and variance $s^2 = \frac{\sum_{j=1}^p f_j m_j^2}{n} - n\bar{x}^2$

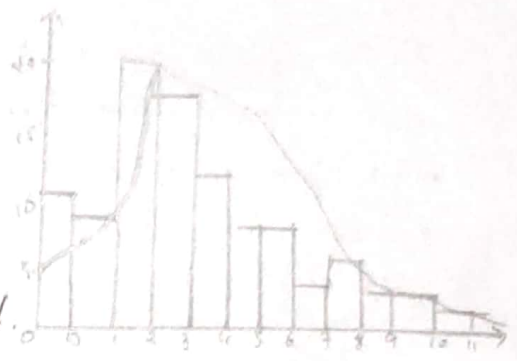
where f_j is observed frequency of j th case interval, m_j is point of j th interval & C is no. of class intervals

④ A parameter is an unknown constant, but an estimator in statistic.
 Eg: Vehicle arrival example. Table in histogram of vehicles example can be analysed to
 $n_x = 100$, $f_i = 12$, $x_i = 0$, $f_2 = 10$, $x_2 = 1$, & $\sum_{i=1}^h f_i x_i = 364$ and $\sum_{j=1}^h f_j x_j^2 = 2080$

The Sample mean & 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.

⑥ However, most samples 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.

4) Goodness - of fit tests:

→ ① previously helpful guidance for evaluating suitability of a input models

② There is no single correct distribution in usual application exists.

③ If very little data are available, it is unlikely to reject all candidate distribution.

④ Conduct hypothesis testing on input data distribution using:

→ chi-square test & → kolmogorov-Smirnov Test.

* Chi-Square test:

one procedure for testing hypothesis random size n of random variable x follows specific distribution form in random variable x follow specific distribution form in
 chi-square, the test procedure begins by arranging n -observation into sets of k class intervals or statistics is given by $\chi^2 = \sum_{i=1}^k \frac{(O_i - E_i)^2}{E_i}$ where $O_i \rightarrow$ observed frequency
 $E_i \rightarrow$ Expected frequency

The expected frequency for each class intervals is computed as $E_i = n P_i$
 where $P_i \rightarrow$ is theoretical hypothesis probability associated with i th class interval.

The hypothesis are following:

H_0 : The random variable x , conforms the distributional assumption with parameters given by parameter estimate.

H_1 : the random variable x does not conform.

⑤ Each value of random variable should be class interval unless combining is necessary as $P = p(x_i) = p(x = x_i)$

⑥ for continuous case with assuming pdf $f(x)$, or assumed cdf $F(x)$, p_i can be computed by, $p_i = \int_{a_{i-1}}^{a_i} f(x) dx = F(a_i) - F(a_{i-1})$

where a_{i-1} $a_i \rightarrow$ endpoints of i th class interval

$f(x)$ assumed pdf

$F(x)$ assumed cdf

Table below are made to aid in determining the no. of class interval for continuous data.

Sample size n	Number of class interval k
20	Do not use chi-square
50	5 to 10
100	10 to 20
100	\sqrt{n} to $n/5$

Table: Recommendation for no. of class-interval for continuous data.

* Chi-Square Test with equal probabilities:

- ① if continuous assumption is being tested class interval are equal in probability rather than equal in width of interval should be used.
- ② Unfortunately there is not method for determining the probability associated with each interval that maximize the power of that of given size.

$E_i = n p_i \geq 5$
Substituting for p_i yields $n/k \geq 5$ & Solving for k yields $k \leq n/5$

* Kolmogorov-Smirnov Goodness-of-fit-tests:

This test is particularly useful when sample size is small & when no parameter have been estimated from data and Ex: Suppose 50 interarrival times are collected over the following 100 min interval.

0.44, 0.53, 2.04, 2.74, 2.00, 0.30, 2.45, 0.52, 2.02, 1.89, 1.53, 0.21, 2.80, 0.04, 1.35, 8.32, 2.34, 1.95, 0.10, 1.42, 0.46, 0.071, 1.09, 0.76, 5.55, 3.95, 1.07, 2.26, 2.88, 0.67, 1.21, 6.26, 6.57, 5.37, 0.12, 3.19, 1.63, 1.46, 1.08, 2.06, 0.85, 0.83, 2.44, 2.11, 3.15, 2.90, 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 dT_1, T_2, \dots } time exponential, the arrival times distributed on interval $(0, T)$. The arrival

$T_1, T_1+T_2, T_1+T_2+T_3, \dots, T_1+\dots+T_{50}$ are obtained by adding times on $\alpha(0,1)$ interval the points will be $[T_1/T, (T_1+T_2)/T, \dots, (T_1+\dots+T_{50})/T]$

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	0.1161	0.3033	0.3866	0.3508
0.3535	0.3561	0.3670	0.3746	0.3746	0.4694	0.4796	0.5225	0.5315	0.5382
0.5494	0.5520	0.5977	0.6514	0.6514	0.6845	0.7009	0.7154	0.7202	0.7468
0.7553	0.7636	0.7870	0.7982	0.7982	0.8417	0.8732	0.9022	0.9660	0.9944

Following the procedure of D^+ of 0.1054 & n of 0.0080 therefore this statistic is $D = \max(D^+, D^-) = \max(0.1054, 0.0080) = 0.1054$. The critical value of D for level of significance of $\alpha = 0.05$ & $n = 50$ & 0.005 , $50 = 1.36 / \sqrt{n} = 0.1923$; but $D = 0.1054$, So the hypothesis that the interarrival times are exponentially observed cannot be rejected.

* P-values & Best fit:

P-value for test statistics - the significance level at which we would reject H_0 for the given statistics.

A measure of fit, the larger the value

① large - P-value: good fit

② Small p-value: poor fit.

5) Fitting non stationary poisson process

- ① Fitting a NSPP data is difficult, possible approaches.
- ② Fit a very flexible model with lots of parameters.
- ③ Approximate constant arrival over some basic interval time but vary it from time interval.
- ④ Suppose we need to model arrival over time $[0, T]$ our approach is most appropriate.
- ⑤ observe the time period repeatedly.
- ⑥ Count arrival / record arrival time.
- ⑦ Divide the time period into k equal interval of length $\Delta t = T/k$
- ⑧ Over n period of observation let i_j no. of arrivals during i th interval in j th period.

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

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

n = no. of observation periods

Δt = time interval length

a_{ij} = no. of arrival during time interval of j th observation period.

Ex: Divide 10 hours business day (8 am, 6 pm) interval $k=60$ whose length $\Delta t = 1/6$ observe over $n=3$ days.

Time Period	No of Arrival			Estimate Arrival Rate (arrivals/hr)
	Day 1	Day 2	Day 3	
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 intervals,
 $1/3(0.5) * (0.3 + 2.6 + 3.2) = 54$ arrival/hr.

6) Selecting input models without data.

→ * if data is not available, some possible sources to obtain information about process are:

① Engineering Data: often product or process has performance ratings provided by manufacturers, or company rules specific time or production standard.

② Expert opinion: pessimistic and most-likely times and they may know the variability.

③ Physical or conventional limitations: physical limit on performance, limit or bounds that narrow the range of input process.

Ex: Production planning simulation.

④ Input of sales volume of various product is required, salesperson of product XYZ says that: → No fewer than 1000 units & no more than 2000 units, at 0.5% chance of selling more than 4500 units.

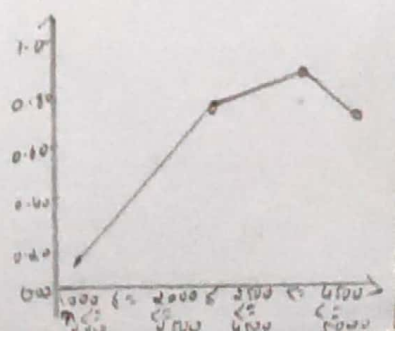
- Translating these information into a cumulative probability or being less than or equal to those goals for simulation

input:

① The nature of the process

② The uniform triangular & beta distribution are often used as input models

i	interval (Sales)	PDF	Cumulative Frequency
1	$1000 \leq x < 2500$	0.1	0.10
2	$2000 \leq x < 2500$	0.65	0.75
3	$2500 \leq x < 4500$	0.24	0.99
4	$4500 \leq x < 5000$	0.01	1.00



7) Multivariate & Time Series input models.

- ⊖ The random variable discussed until now we considered to be independent of any other variables within the context of the problem.
 - However, variable may be related
 - If they appear as input, the relationship should be investigated & taken into consideration

* Multivariate input model:

- ⊖ Fixed, finite numbers of random variable x_1, x_2, \dots, x_n
- ⊖ Ex: lead time & annual demand for an inventory model.
- ⊖ An increase in demand results in lead time increase, hence variables are dependent

* Time Series input Model:

- ⊖ Infinite sequence of random variables, ex: x_1, x_2, x_3, \dots
- ⊖ Ex: time between arrival of orders to buy & sell stocks.
- ⊖ Buy & sell orders tends to arrive in bursts; hence, time between arrivals are dependent.

8) Types of simulation with respect to output analysis.

- ⊖ Terminating versus non-terminating simulation.

⊖ Terminating Simulation:

- Runs for some duration of time T_ϵ , where ϵ is a specified unit that stops the simulation.
- Starts at time 0 under well-specified initial conditions.
- Ex: opens at 8:30 am (time 0) with no customers present at 8 of the or 11 teller working. (initial condition) & closes at 4:30 pm (Time $T=80$)
- The simulation analyst chooses to consider it a terminating system because the object of interest is one day operation.

9) Stochastic nature of output data:

- ⊖ Model output consists of one or more random variables (x, v) because the model is an input-output model transformation & input variables are r.v.s.

⊖ M/M/1 queuing eg:

- poisson arrival rate = 0.1 per minute, service time $\sim N(m=9.5, s=1.75)$

- System performance: long-term mean queue length $L_q(t)$.

- Suppose we seen a single simulation for a total of 5000 minutes.

- Divide the time interval $[0, 5000]$ into 5 equal subinterval of 1000

- Average no. of customers in queue from time $(j-1) 1000$ to $j(1000)$ is y_j

• M/M/1 queuing eg (cont:)

- Batched Average queue length for 3 independent replication

Batching interval (minutes)	Batch j	REPLICATIONS		
		1, y_{1j}	2, y_{2j}	3, y_{3j}
0, 1000	1	3.61	2.91	2.67
1000, 2000	2	3.21	9.00	19.53
2000, 3000	3	2.18	16.15	20.36
3000, 4000	4	6.92	24.53	8.11
4000, 5000	5	2.82	25.19	12.62
0, 5000		3.75	15.56	13.66

① Inherent variability in stochastic simulation both within a single replication and across different application.

② The average across 3 replications can be regarded as independent observations, but average within a replication y_{1j}, \dots, y_{3j} are not.

(10) Absolute measure of performance and their estimation.

→ ① Consider the estimation of a performance parameter, g (or F) of simulated system.

→ Discrete time data: $[y_1, y_2, \dots, y_n]$ with ordinary mean: \bar{g}

→ Continuous time data: $\{y(t) \mid 0 \leq t \leq t_e\}$ with time weighted mean: \bar{f}

② Point estimation of discrete time data.

→ The point estimator: $\hat{g} = \frac{1}{n} \sum_{i=1}^n y_i$, is unbiased if its expected value is g .
i.e. if: $E(\hat{g}) = g$, is biased if:

* point estimator:

• point estimator for continuous time data

→ The point estimator: $\hat{f} = \frac{1}{t_e} \int_0^{t_e} y(t) dt$

→ Is biased in general where: An unbiased or low-bias estimator is desired.

• Usually, system performance measures can be put into the common framework of g or f .

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

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

③ Performance measure that doesn't fit: quantile or percentile:

→ estimating quantile: the increase of problem of estimating a proportion or probability $\text{Pr}\{y \leq \alpha\} = P$.

→ Consider a histogram of observed value y : Find such that 100p% of the histogram is to the left of y .

* Confidence-Interval Estimation:

① To understand this fully, it is important to distinguish b/w measures of error, risk, eg: confidence interval versus prediction interval.

② Suppose the model is normal distribution with mean μ , variance σ^2 (both unknown).

→ let y_i be avg cycle time for parts produced on i th replication of simulation.

→ Avg cycle time will vary from day to day, but over the long-term the avg of the averages will be close to μ .

→ Sample variance across R replication: $s^2 = \frac{1}{R-1} \sum_{i=1}^R (y_i - \bar{y})^2$

* Confidence - Interval estimation

→ Confidence interval (CI):

• A measure of error.

• Where y_i are normal distributed: $\bar{y} \pm t_{\alpha/2, R-1} \frac{s}{\sqrt{R}}$

• 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 the error.

• The more replication we make, less error occurs.

→ Prediction interval (PI):

• A measure of risk.

• A good guess for average cycle time on a particular day is over 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:

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

• The length of PI will not go to 0 as R increases because we can never simulate away risk.

• Pls limit is: $\bar{y} \pm Z_{\alpha/2} \sigma$

ii) Output analysis for terminating

→ A terminating simulation: runs over simulated time interval $[0, T_E]$.

A common goal is estimate:

ⓐ $\bar{\theta} = E\left[\frac{1}{n} \sum_{i=1}^n y_i\right]$, for discrete output.

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

In general, 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.

(i) 2 performance measures of system: cycle time for part and work in (wp).

(ii) Across replication data are formed by summarizing within replication data.

ⓐ Overall sample avg and interval replication sample avg are always unbiased estimates of expected daily avg cycle time or daily avg wp.

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

ii) 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.

$$\Theta = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n Y_i \text{ , for discrete measure.}$$

$$\Phi = \lim_{T \rightarrow \infty} \frac{1}{T} \int_0^T y(t) dt, \text{ continuous measure.}$$

• The sample size is design choice, with several consideration in mind.

(i) Initialization Bias: Method to reduce the point-estimator bias caused by using artificial & unrealistic.

→ Intelligent initialization

• Initialize the simulation in 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.

→ Divide simulation into 2 phases:

• An initialization phase from 0 to time T_0 .

(ii) Error Estimation: ⊙ if Y_1, \dots, Y_n are not statistically independent then \bar{Y} is biased estimator of true variance.

⊙ Suppose point estimator \bar{Y} is sample mean $\bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i$

⊙ Variance of \bar{Y} is very hard to estimate.

⊙ For system with steady state, produce an output process that is approximately stationary.

(iii) Replication Method: ⊙ Approach make R replication initializing & deleting from each one the same way.

⊙ Important to do through if obj of investigation condition bias:

• Basic row output data $\{Y_{rj}, r=1, R, j=1, \dots, n\}$ is derived

• Individual observation from with 'r'.

• Proch mean from within replication 'r' of some of no. of discrete time observation.

13) Model building, verification & validation.

• ⊙ The first step in model building consists of observing real system & interaction among its various components & collecting data on its behaviour.

⊙ Operators technicians, repairs & maintenance personnel, engineer & managers under certain aspects of system which may be unfamiliar to others.

⊙ As model development proceeds, new question may arises & model developers will returns to this step of learning true system structure & behaviors.

⊙ The second step in model building is the construction of a conceptual model.

→ a collection of assumption on components & structure of system, plan hypothesis & the value of model into parameters, illustrated by following figure.

② The third step is translation of the operational model into computer recognizable form - the computerized model.

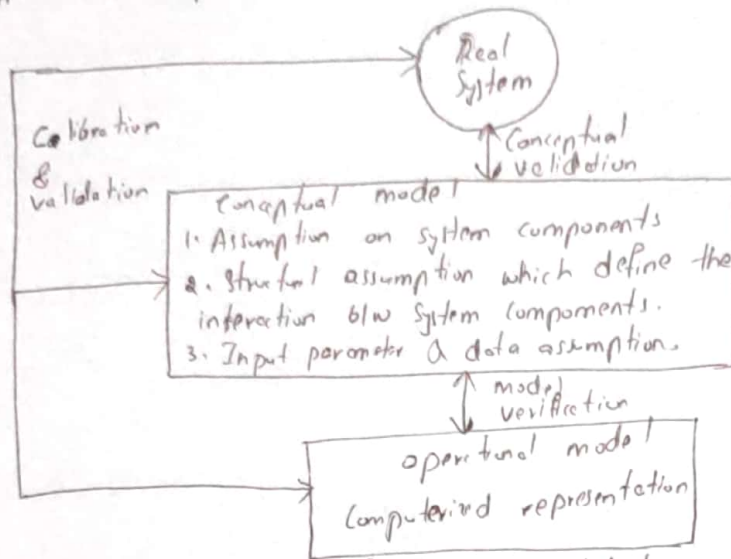


Fig: Model Building, verification & validation.

14) Verification of simulation model.

→ The purpose of model verification is to assume that conceptual model is reflected accurately in computerized representation:

② Make someone else check the model.

② Make a flow diagram that indicates each logically possible action a system can take during an event.

② Closely examine the model output for reasonableness under variety of input parameters.

② Make the operational model or self documentary as possible.

② The interactive run controller (IRC) assist in debugging in the following ways:

② Simulation can be focused on a particular line of logic.

② Attention can be focused on a particular line of logic.

② Values of selected model components can be observed.

② The simulation can be temporarily suspended or paused, not only to view information but also to reassign value or redirect entities.

② Graphical interfaces are recommended for accomplishing verification & validation.

15) Calibration & validation for models, optimization via simulation.

→ ② verification & validation although are conceptually distinct, usually are conducted simultaneously by the models.

② Validation is the overall process of comparing the model to real system & its behaviour.

② Calibration is the interactive process of comparing the model to real system,

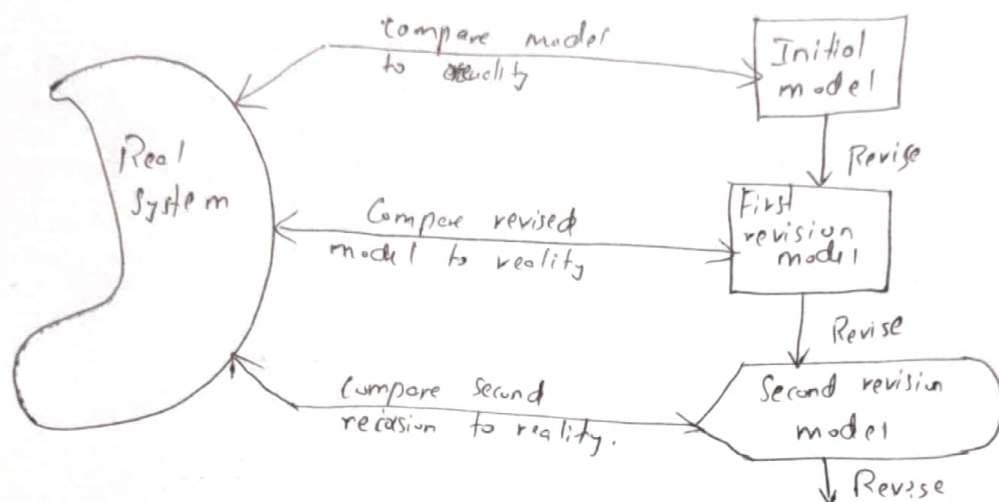
making adjustments to the model, comparing again & so on.

② Comparison of model to real system.

③ Subjective tests - Requires data on real systems behaviours & the output of the model.

④ Subj validation is not an proposition - no model is ever totally representative of system under study.

⑤ In addition, each revision of model, as in fig. involves some cost, time and efforts.



⑥ As an aid in validation process, Naylor & Tinger 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 input-output transformation for the real system.