

Introduction to Statistics and Bayesian Networks

Soma Sekhar Dhavala
Iowa State University

what is a statistic

- any function of the data

mean,
max,
absolute
value ...

why statistics

- to quantify uncertainty
- express randomness
- make an informed decision
(however bad it may be :))

(how likely is it that it rains today
given that it rained yesterday and
monsoon hit the western ghats)

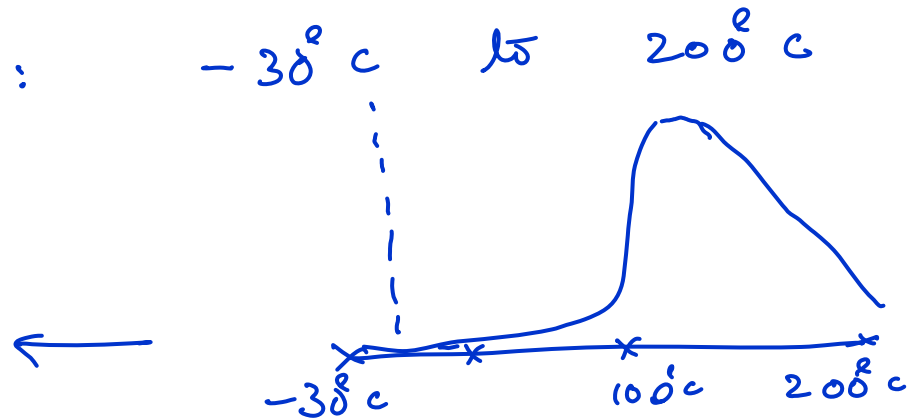
Key terminology :

random variable: . the observations of which over a period of time are possibly different characterized by a distribution

distribution : a function of the random-variable, when evaluated at a possible outcome / realization of a r.v's, tells us how probable that r.v can be observed over a sufficient length of interval

e.g of a r.v
boiler
temperature :
(θ)

$f(\theta)$



e.g.

$$f(99 \leq \theta \leq 102) = 0.5$$

$$f(-38 \leq \theta \leq -10^\circ\text{C}) = 0.0001$$

$$\int_{\theta = -38^\circ\text{C}}^{200^\circ\text{C}} f(\theta) d\theta = 1 \quad \begin{matrix} f(\theta) \neq 0 \\ \geq 0 \end{matrix}$$

laws of probability

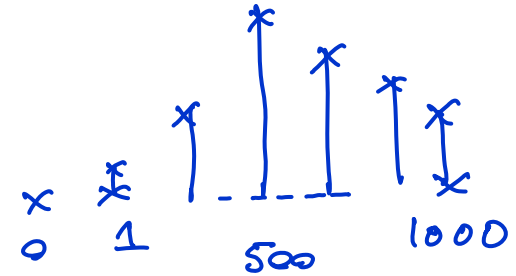
discrete random variables

ATM transactions 0 to 1000

the distribution is now called

probability mass function

or
(pmf)



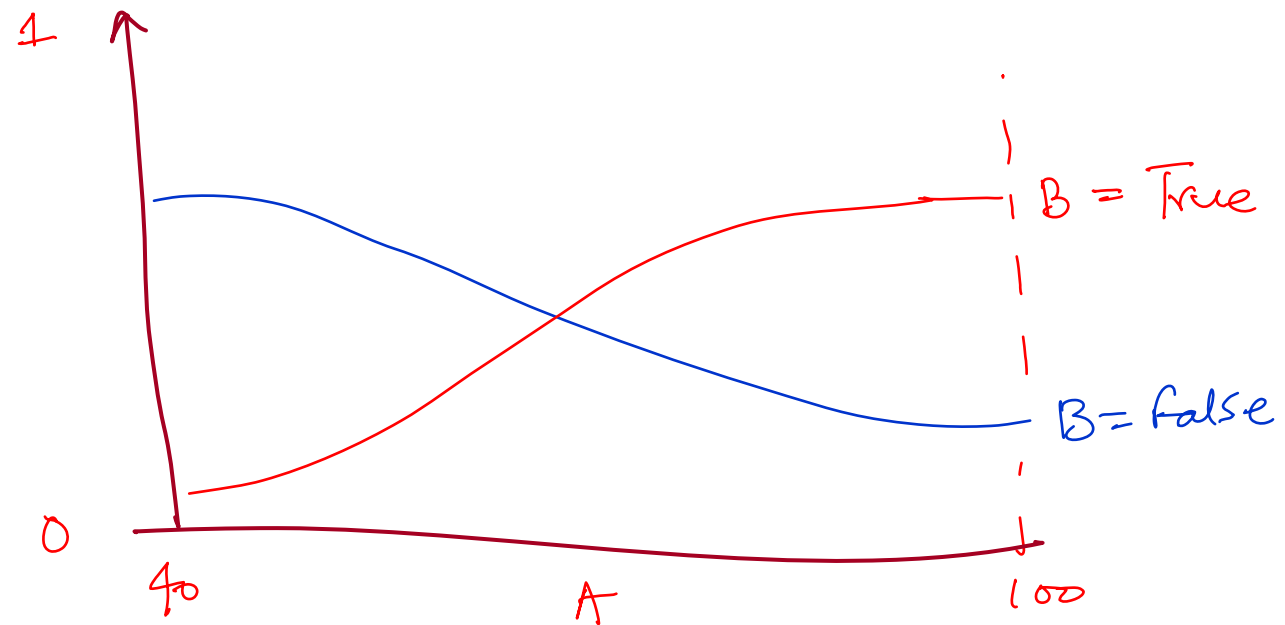
if the r.v is continuous, it is prob density function
(pdf)

$$P(\# \text{ATM} = 100) = 0.01$$

Several r.v.'s (continuous or discrete)

A : price of petrol/liter (Rs 40/- to 100/-)

B : war on terrorism (True or false)



$p(A, B)$ is the joint distribution

If $P(A, B)$ is known to JS,
we can draw several inferences
such as

Given that oil price is Rs 90/- per
liter, which is more likely
was on terrorism true or false
— x —

$$P(B / A = 90 \text{ Rs}) \quad ??$$

$$P(B = T / A = 90) > P(B = F / A = 90)$$

what is a statistical model??

eg. Yield of Corn (is the variable of interest)
in tons/sq. mile

$$Y = \beta_1 X_1 + \beta_2 X_2 + \dots + \underbrace{\epsilon}_{\text{error terms}}$$

$\underbrace{Y}_{\text{observed quantity}}$
 $\underbrace{\beta_1 X_1}_{\text{amount of nitrogen}}$
 $\underbrace{\beta_2 X_2}_{\text{\% rain fall}}$
 $\underbrace{\epsilon}_{\text{error terms}}$
 $\underbrace{\beta_1, \beta_2, \dots}_{\text{are unknowns}}$
 $\underbrace{X_1, X_2, \dots}_{\text{control variables}}$

Given the observations, what are the likely values of β 's which could have possibly produced the observations that we have at hand

Different estimation techniques

- Maximum likelihood
- Maximum a posteriori
- Bayesian methods

Then we can optimize the parameters of the model to yield maximum gains

max γ
 given β 's w.r.t $(x_1, x_2 \dots x_N)$

100 ppm Sucfar ~~with~~ under 10 mm
 rainfall will yield maximum returns

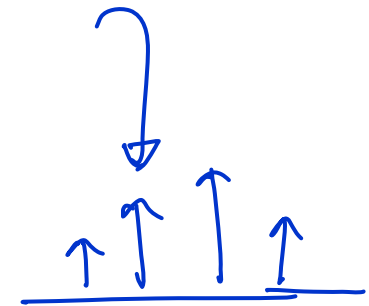
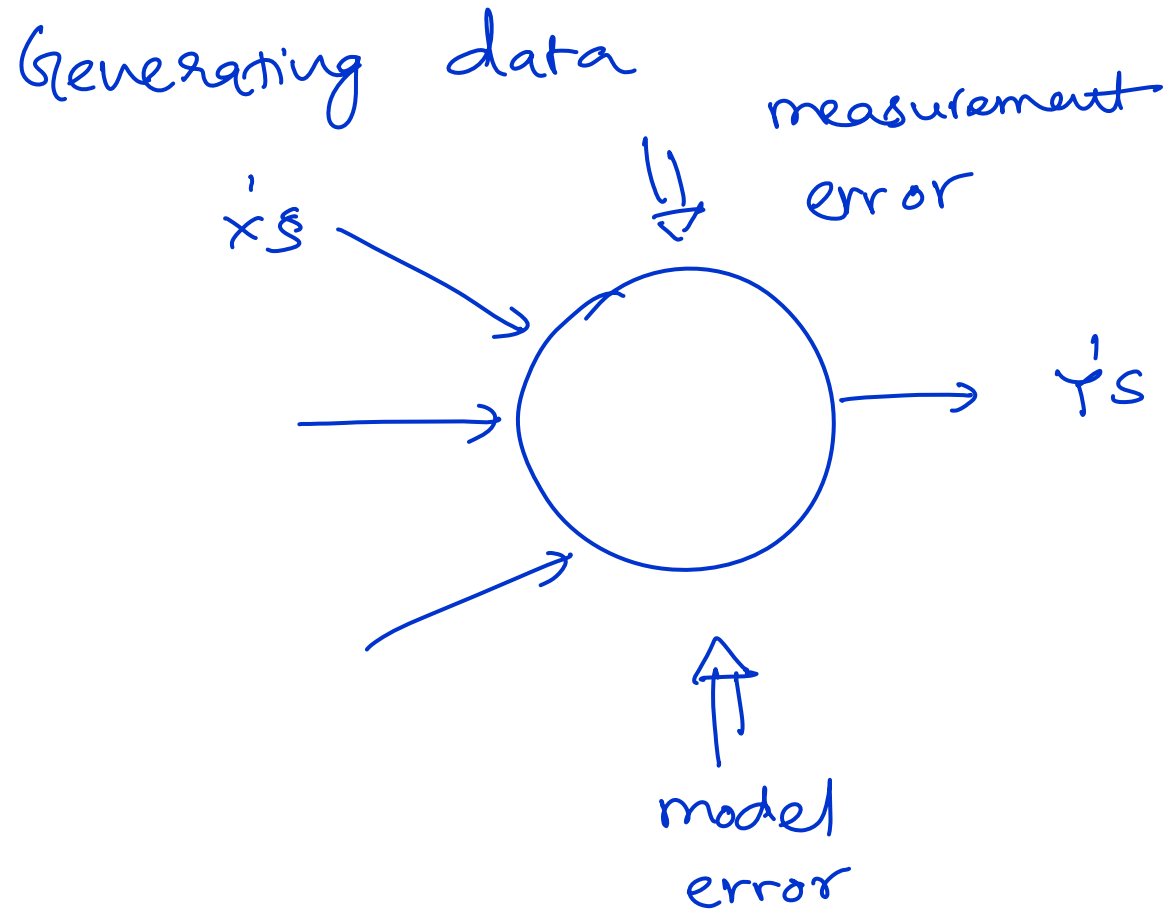
But how do we generate data
 Design of experiments

fix the inputs (controlled / uncontrolled)

set the objective of the experiment

study the effects of varying the
input on the outcome in a
methodological fashion

$f(\text{inputs}) \rightarrow \text{output}$
how many inputs
how many levels / factors each input takes



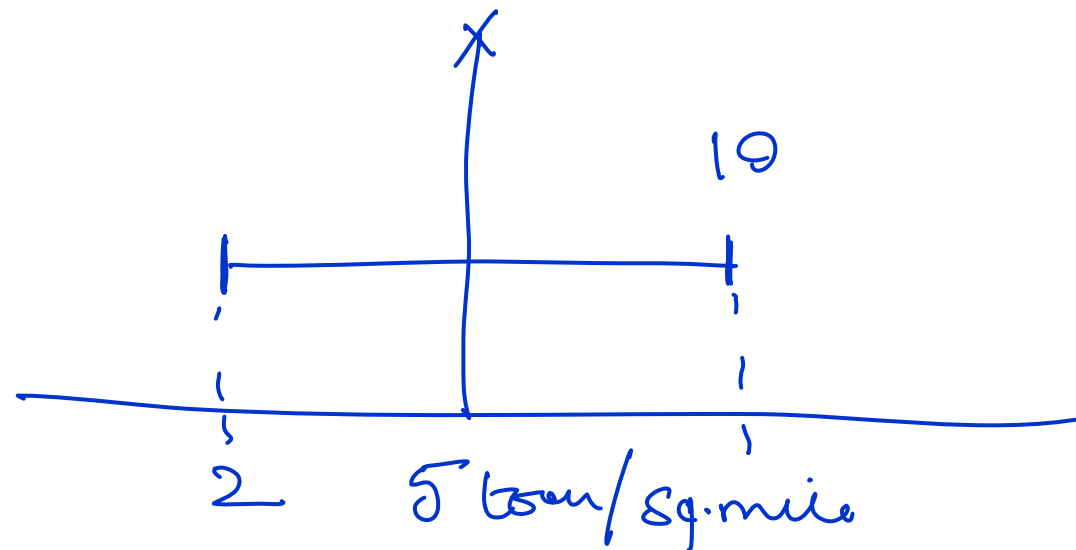
histogram/
distribution
of y 's

Some inferences on the yield

→ avg. yield

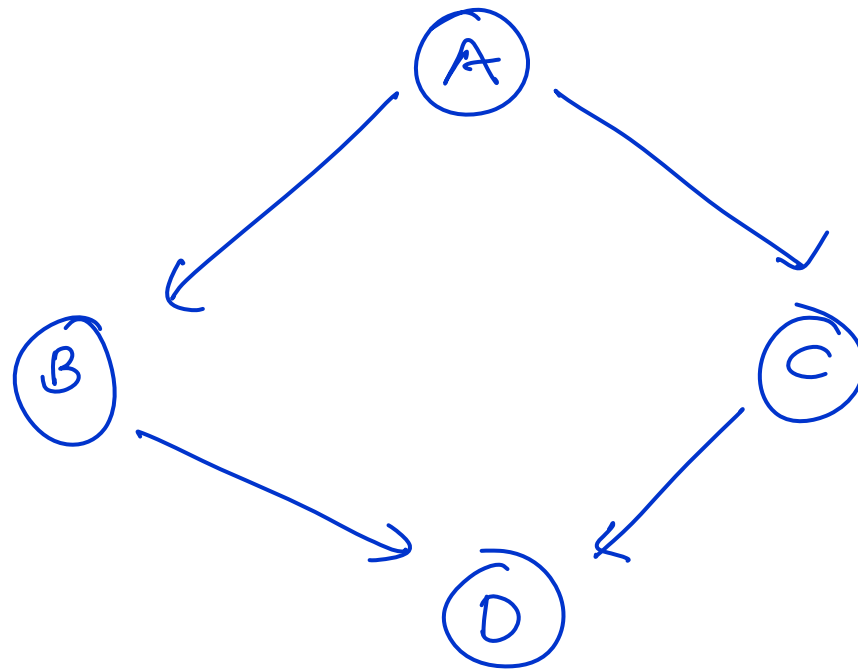
$$\underbrace{E(Y)}$$

and what is the confidence interval



95% confidence level

Bayes Networks



- Parent-child relationship
- nodes are random variables
- Cause — effect analysis possible either ways