

Chapter-1

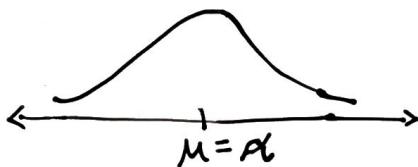
Measures of location

Mean → gives central tendency of data.
 (average) → determined by adding all the data points in a population & then dividing total by the no. of points.

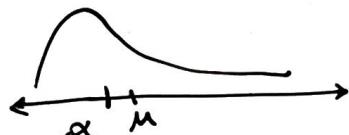
$$\bar{x} = \frac{\sum_{i=1}^n x_i}{n}$$

Median → middle value once the data is sorted (ordered) from small to largest.
 (middle)

$$\tilde{x} = \begin{cases} \text{single mid value} & \text{(odd)} \\ \text{avg of } \left(\frac{n}{2}^{\text{th}} \text{ & } \frac{n+1}{2}^{\text{th}} \right) & \text{(even)} \end{cases}$$

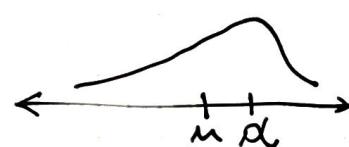


Symmetric



positive skew
 $\alpha < \mu$

μ = mean median = \tilde{x} or α (here)



negative skew
 $\mu < \alpha$

Other measures : Quartiles, Percentiles & trimmed means

10% trimmed mean for eg.
 is mean calc. by eliminating smallest 10% & largest 10%.

Measures of variability

→ gives deviation from above computed central values.
 we need to sometime find out the variability in data.

We cannot simply add the difference of all values from mean since,

$$\sum_{i=1}^n (x_i - \bar{x}) = 0$$

∴ Other ways are:

do absolute value $|x_i - \bar{x}|$ summation

But this isn't a good idea theoretically

square summation then root of result

↳ Pros: solves our problem as well as many theoretical problems

Cons → Works crap with really skew data.

* So, sample variance (s^2) is

$$s^2 = \frac{\sum (x_i - \bar{x})^2}{n-1} = \frac{s_{xx}}{n-1}$$

* Sample standard deviation $= \sqrt{s^2}$

* Standard deviation is σ^2 for variance, & σ for standard deviation

$$\sigma^2 = \frac{\sum (x_i - \mu)^2}{N}$$

Note

($n-1$ is statistical correction to compensate if data is extremely non symmetrically skewed.)

If observation is within 1.5 fs from closest fourth it's mild

further than 1.5 fs → outlier

further than 2 fs → extreme.

Ch2 - Probability

→ finding out likelihood of an outcome amidst multiple outcomes of an event

Sample Space (S)

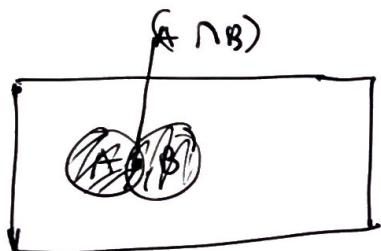
→ Set of all possible outcomes of an experiment

Event

→ any collection of outcomes contained in the sample space S

Simple if exactly 1 outcome

Compound if multiple outcomes that comprise.



$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

* Product Rule:

1st element of ordered pair can be chosen in n_1 ways
& for each of these n_1 ways, second element can be chosen in n_2 ways, then no. of pairs is $n_1 * n_2$

* Permutation & Combination

↓
ordered ↓
unorderd.

$P \rightarrow$ no. of permutations of size k that can be formed from n individuals or objects

$$P_{(n,r)} = \frac{n!}{(n-r)!}$$

subset is called combination.

$$C_r = \frac{n!}{r!(n-r)!}$$

Conditional Probability:

↪ probability of A, given event B has occurred

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

$$P(A \cap B) = P(A|B) \cdot P(B)$$

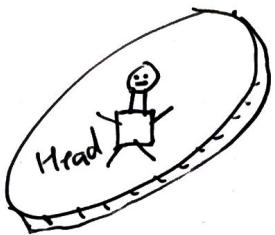
Bayes theorem:

Let A_1, \dots, A_n be mutually exclusive events.

for any event B, $P(B) = P(B|A_1)P(A_1) + \dots + P(B|A_n)P(A_n)$

$$P(B) = \sum_{i=1}^n P(B|A_i)P(A_i)$$

Binomial Variable



lets consider an unfair coin,

$$P(H) = 0.6$$

$$P(T) = 1 - 0.6 = 0.4$$

$X = \# \text{ of heads after } 10 \text{ flips of my coin.}$

(*) These sort of problems can be categorized if they follow certain rules:

- 1) The outcome of each trial can be classified either as a success or failure
- 2) Each trial is independent of the others
- 3) There is a fixed number of trials
- 4) The probability of success on each trial remains constant.

If these conditions are fulfilled we can employ a Binomial variable to solve our problem.

* How to Recognize a binomial Variable you ask?

Ans: Well the process is simple. You ask intelligent questions about the conditions in exhibit A.

Illustration If we are keeping track of # of Kings in a deck by checking if the card is King or not,

i) if we check the card then remove it from deck, we can say its ~~not~~ we can't employ a binomial variable.

(since after each trial, we change the no. of cards in deck, & hence affect (change) the probability of each new card about to be drawn. ~~the~~ each experiment (trial) is now dependent on previous outcome)

trial 1

$$P(\text{success}) = \frac{a}{52}$$

trial 2

$$P(\text{success} | \text{success in } T_1) = \frac{b}{51}$$

* Why we use binomial variables??.

We can say the distribution is binomial,
useful inferences.

i) [sum of a bunch of independent trials]

* 10% Rule?

if sample \leq 10% of population
 \approx independence

lets consider an example

$X = \# \text{ of } H \text{ from flipping coin } 5 \text{ times}$

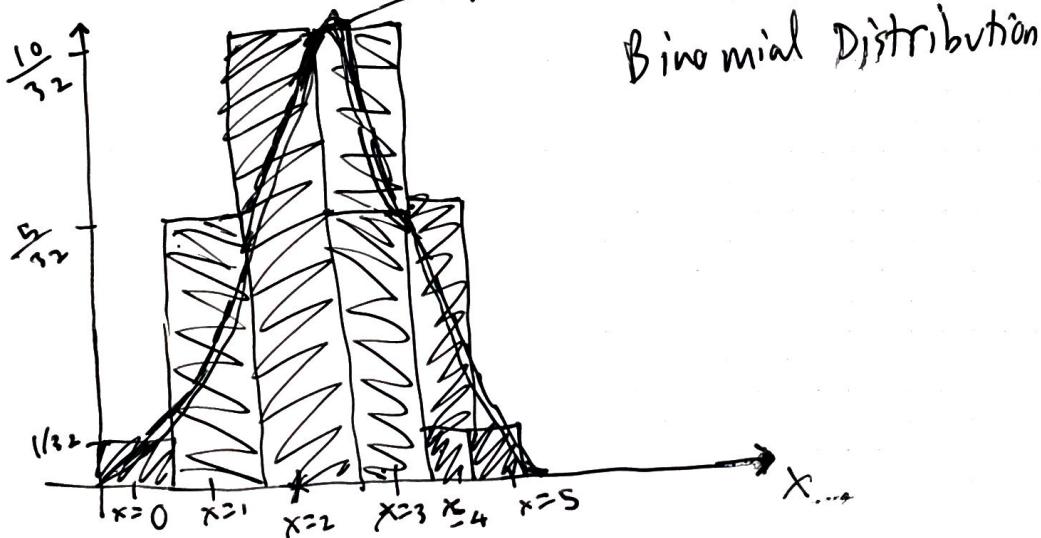
possible no. of outcome from 5 flips $= 2 \cdot 2 \cdot 2 \cdot 2 \cdot 2$
 $= 2^5 = 32$

$$P(X=0) = \frac{1}{32} = \frac{^5C_0}{32} \quad P(X=2) = \frac{^5C_2}{32} = \frac{10}{32}$$

$$P(X=1) = \frac{5}{32} = \frac{^5C_1}{32} \quad P(X=3) = \frac{^5C_3}{32} = \frac{10}{32}$$

$$P(X=4) = \frac{^5C_4}{32} = \frac{5}{32} \quad P(X=5) = \frac{^5C_5}{32} = \frac{1}{32}$$

If we plot this
normal distribution : Bell curve



- * Binomial distribution is discrete version of normal distribution
 - more observations will give out closer & closer discrete versions of the normal distributions
 - i.e. more x values, smoother the bar graph.

Calculating Probability

If I'm playing basketball,

$$\text{prob(score)} = 70\% \text{ or } 0.7$$

$$\text{prob(miss)} = 30\% \text{ or } 0.3$$

$P(\text{exactly 2 scores in 6 attempts})$

$$SS \quad MM \quad M \quad M \\ (0.7)(0.7) (0.3)(0.3)(0.3)(0.3) = (0.7)^2 (0.3)^4$$

→ this combination prob

Similarly,

$$MS \quad MS \quad MM \\ (0.3)(0.7)(0.3)(0.7)(0.3)(0.3)$$

$= (0.7)^2 (0.3)^4$ → this combination prob

We can see many ways one could miss but it's

$$(0.7)^2 (0.3)^4$$

nonetheless

added to however many times it happens

$$\binom{6}{2}$$

∴ total prob of $P(\text{exactly 2 scores in 6 attempts})$

$$= \binom{6}{2} 0.7^2 0.3^4$$

or

$$\left(\frac{\text{total attempts}}{\text{favorable outcomes}} \right) \times \left(\frac{\text{chance of favorable}}{\text{chance of fail}} \right)$$

Hypergeometric Distribution:

e.g.: 6 docs, 19 nurses attend a small conference.

all 25 names in a list, 5 names picked randomly
without replacement

what's $P(4 \text{ docs}, 1 \text{ nurse selected})$?

sol

trials not independent

Binomial dist + not relevant

let's see 😊

$P(4 \text{ docs}, 1 \text{ nurse})$

$$\frac{\binom{6}{4} \binom{19}{1}}{\binom{25}{5}} = \frac{15 \cdot 19}{53130} = 0.00536$$

not easy to do with ~~too~~ complex stuff.

* Had the sampling been done with replacement, using binomial distribution, would be appropriate.

$$P(4 \text{ docs}) = \binom{5}{4} \left(\frac{6}{25}\right)^4 \left(1 - \frac{6}{25}\right)^{5-4}$$

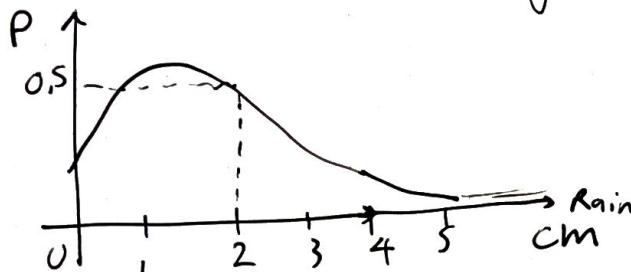
$$\boxed{\frac{= 0.0126}{\text{with replacement!}}}$$

$$\boxed{\begin{array}{l} \text{Without replacement} \\ 0.00536 \end{array}}$$

Probability density function Section 4.1

Continuous RVs need these to express data.

Eg $Y = \text{exact amount of rain tomorrow (cm)}$



$$P(Y=2) = ?$$

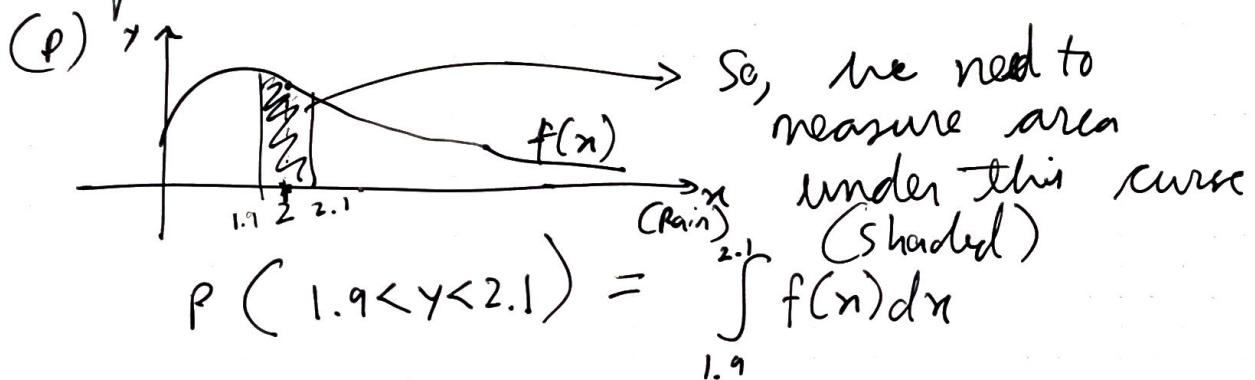
~~$P(Y=2) = 0.5$~~

Wrong!

We can't be that accurate!

So, we make bounds for Y 's value like

~~$\star P(|Y-2| < 0.1) \Rightarrow P(1.9 < Y < 2.1)$~~



~~\star also, $\int_0^\infty f(n) dn = 1$~~

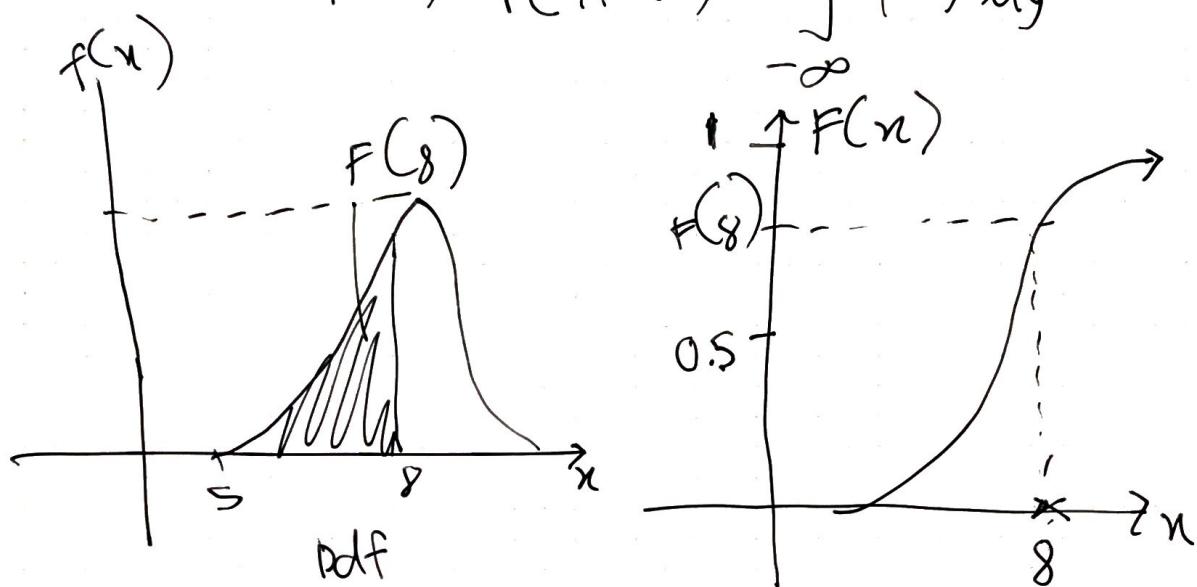
DEF: A continuous RV X is said to have a uniform distribution on the interval $[A, B]$ if pdf of X is

$$f(n; A, B) = \begin{cases} \frac{1}{B-A} & A \leq n \leq B \\ 0 & \text{otherwise} \end{cases}$$

4.2 Cumulative Distribution Functions & EV

CDF is just $F(x)$

$$F(x) = P(X \leq x) = \int_{-\infty}^x f(y) dy$$



~~*~~ $P(X > a) = 1 - F(a)$

~~*~~ $P(a \leq X \leq b) = F(b) - F(a)$

~~*~~
$$\boxed{F'(x) = f(x)}$$

Percentiles

$$P = F(\eta(p)) = \int_{-\infty}^{p(1)} f(y) dy$$

$(100 \cdot p)^{\text{th}}$ percentile $\leadsto \eta(p)$

$$M_x = E(x) = \int_{-\infty}^{\infty} x \cdot f(x) dx$$

* $E[h(x)] = M_{h(x)} = \int_{-\infty}^{\infty} h(x) \cdot f(x) dx$

* $\sigma_x^2 = V(x) = \int_{-\infty}^{\infty} (x - \mu)^2 \cdot f(x) dx$

$$\sigma_x = \sqrt{V(x)} = \sqrt{E((x - \mu)^2)}$$

* $V(x) = E(x^2) - [E(x)]^2$

Normal Distribution

$$\text{if } f(x; \mu, \sigma) = \frac{1}{\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

if $\mu=0, \sigma=1$, \Rightarrow standard normal distribution

& we have standard normal RV $\Rightarrow Z$

$$\text{pdf of } Z = f(z; 0, 1) = \frac{1}{\sqrt{2\pi}} e^{-\frac{z^2}{2}}$$

$$\text{cdf of } Z \quad P(Z \leq z) = \int_{-\infty}^z f(y; 0, 1) dy$$

which is denoted by

$$\Phi(z)$$

$$P\{a \leq Z \leq b\} = \Phi(b) - \Phi(a)$$

$$\& \Phi(\infty) = 1$$

$$Z = \frac{x - \mu}{\sigma} \quad \text{by} \quad P(Z \leq z) = \int_{-\infty}^z \frac{1}{\sqrt{2\pi}} e^{-\frac{y^2}{2}} dy = \Phi(z)$$

$$E[X] = \mu \quad \& \quad V[X] = \sigma^2$$

Gamma Distribution & Exponential Distribution

Expo for ~~the~~ param $\lambda > 0$
if pdf is

$$f(x; \lambda) = \begin{cases} \lambda e^{-\lambda x} & \text{if } x \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

$$E(X) = \frac{1}{\lambda} \quad V(X) = \frac{1}{\lambda^2}$$

Gamma

for $\alpha > 0$

$$\Gamma(\alpha) = \int_0^\infty x^{\alpha-1} e^{-x} dx$$

Properties

$$1) \alpha > 0 \quad \Gamma(\alpha) = (\alpha - 1) \Gamma(\alpha - 1)$$

$$2) \text{any int } n \quad \Gamma(n) = (n - 1)!$$

$$3) \Gamma\left(\frac{1}{2}\right) = \sqrt{\pi}$$

$$\text{if } \alpha > 0, \beta > 0 \text{ if pdf}(x) = f(x; \alpha, \beta) = \frac{1}{\beta^\alpha \Gamma(\alpha)} x^{\alpha-1} e^{-\frac{x}{\beta}}$$

$f(x; \alpha, \beta)$ is standard gamma if $x > 0$, otherwise

$$E(x) = \mu = \alpha\beta$$

$$V(x) = \sigma^2 = \alpha\beta^2$$

in standard gamma dist,

$$F(x; \alpha) = \frac{\int_0^x y^{\alpha-1} e^{-y} dy}{r(\alpha)}$$

X be a gamma dist ,

$$P(X \leq n) = F(n; \alpha, \beta) = F\left(\frac{n}{\beta}; \alpha\right)$$

CHI SQ distribution

$v = \text{int, greater } 0$

if pdf of x is gamma density with $\alpha = v/2$ & $\beta = 2$

pdf of chi-sq rv is

$$f(x; v) = \begin{cases} \frac{1}{2^{v/2} r(v/2)} x^{v/2-1} e^{-x/2} & \text{if } x \\ 0 & \text{otherwise} \end{cases}$$

WEIBULL DIST

if $\alpha > 0, \beta > 0$ if pdf of x is

$$f(x; \alpha, \beta) = \begin{cases} \frac{\alpha}{\beta^\alpha} x^{\alpha-1} e^{-(x/\beta)^\alpha} & \text{if } x \geq 0 \\ 0 & \text{otherwise.} \end{cases}$$

we get weibull dist

cdf of weibull is

$$F(x; \alpha, \beta) = \begin{cases} 1 - e^{(-x/\beta)^\alpha} & \text{if } x \geq 0 \\ 0 & \text{otherwise.} \end{cases}$$

LOG NORMAL DIST

$$f(x; \mu, \sigma) = \begin{cases} \frac{1}{x \sqrt{2\pi}} e^{-\frac{(\ln x - \mu)^2}{2\sigma^2}} & \text{if } x > 0 \\ 0 & \text{otherwise.} \end{cases}$$

we have lognormal dist.

$$E(x) = e^{\mu + \frac{\sigma^2}{2}} \quad V(x) = e^{\mu + \sigma^2} (e^{\sigma^2} - 1)$$

$$F(x; \mu, \sigma) = P(X \leq x) = P(Y \leq y) = P(Z \leq z)$$

$$= P\left(Z \leq \frac{\ln x - \mu}{\sigma}\right) = \Phi\left(\frac{\ln x - \mu}{\sigma}\right)$$