Poisson's distribution assumptions → ① $n \rightarrow \infty$
② events are independent
③ p can be very, very less
sometime p = 0 (p does not matter)

classmate
Date _____
Page _____

→ **POISSON'S DISTRIBUTION :**

IF $n \rightarrow \infty$ (very large)

$$P(X = x) = \frac{e^{-\lambda} \cdot \lambda^x}{x!} \qquad e = 2.72 \text{ (constant)}$$

$\lambda$ = no. of objects per $\begin{cases} \text{time} \\ \text{length} \\ \text{area} \\ \text{volume} \end{cases}$

→ Scales of measurement

① nominal $\Big\}$ categorical data → Shown in bar/pie chart.
② ordinal
③ interval and ratio.

gives Freq

↳ gives rel. Freq.

① nominal : cannot apply $\overset{\text{(mathematical - quantity)}}{\text{sort, add, sub, multiply}}$ division. operations

ex: texts, phone nos.
can be numeric as well (ex : phone no.) → cannot perform any operations on it.

② ordinal : ordered
• ex : grades (A, B, C), quality (high, low), designation
• only sort applicable
(mathematical ops)
• $+, -, *, /$ → not applicable
• exception : Likert scale
In case of reviews, we take average reviews

③ interval and ratio
• can apply all mathematical/arithmetic operations
• it is a scale ← ratio
• all numeric data

numerical $\longrightarrow$ discrete (count)
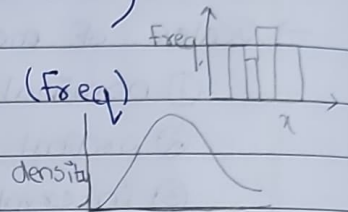$\searrow$ continous (measure)

① discrete : count, variable
ex: no. of people
whole no. /integers

② continous : can be decimal also

if discrete repeats - it can be converted into
categorical for better understanding

numerical data $\longrightarrow$ ① histogram (freq)
② density plot
③ box plot - to understand
5 point summary

$\rightarrow$ ① univariate $\rightarrow$ bar, histogram, pie
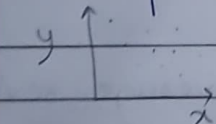② bivariate

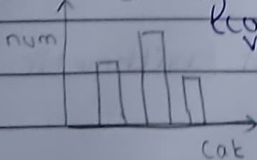② bi-variate : a) numerical & numerical (vs)
b) num & cat. (vs)
c) c & cat (vs)
d) time stamp & num (vs)
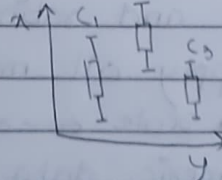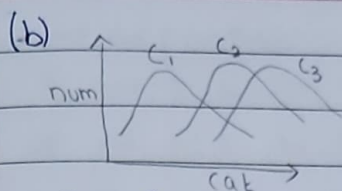• scatterplot $\rightarrow$ covariance or correlation

(a) BAR by aggregate
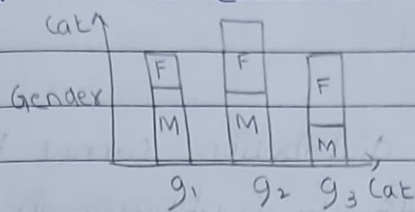(count, sum, var)

(b) BOXPLOT by group
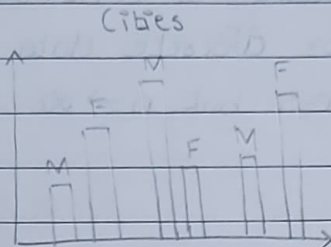
(b)



Density plot by group.

(c) cat. vs cat.



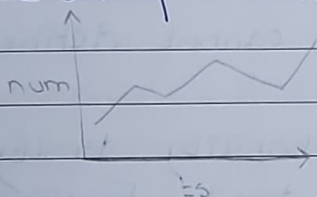stacked bar chart.

F, M → gender
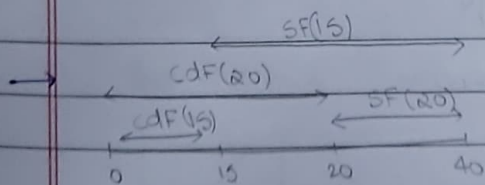
Cities



bar chart by group

d) time stamp vs. numeric.



line chart.

metric : noise
       seasonality
       trends.

→ Data analysis $\longrightarrow$ ① descriptive - pop. data available.
                 $\hookrightarrow$ ② inferential data - when population
                    data is not available
                    - we consider sample data
                    for analysis



3 ways ① $cdF(20) - cdF(15)$
       ② $SF(15) - SF(20)$
       ③ $1 - [cdF\ 15 + SF\ 20]$

to find
probs.
b/w 15 & 20

→ Poisson's distribution:

$$P(x = \lambda) = \frac{e^{-\lambda} \cdot \lambda^{\lambda}}{\lambda!}$$

$$= F(\lambda)$$

pmf → st. poisson. pmf

cdf → st. poisson. cdf
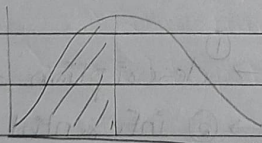
sf → st. poisson. sf

→ Cases when poisson's distribution cannot be applied:
① when random variable is continous in nature
② applicable only for discrete data (count)
③ when p is not less but n → ∞

⇒ Continous prob. dist.
· when random variable is continous
· since data is continous, cannot define probability at a given point
· But you can define probability density using pdf → probility density function $\left( \frac{P(x)}{\lambda} \right)$
· since continous, no difference between
   $P(x < 35)$ and $P(x \leq 35)$ → $\frac{50}{60}$ we always
   use $P(x \leq 35)$



P (X > 35) and P (X ≥ 35) are same

Formula for bell curve:

$$F(x) = \frac{1}{\sigma \cdot \sqrt{2\pi}} e^{-\frac{1}{2} \left( \frac{x - \mu}{\sigma} \right)^2}$$

↑
pdf

$e \rightarrow$ exponent

$\pi \rightarrow 3.14$

$\sigma \rightarrow$ std. deviation

$\mu \rightarrow$ mean

$\rightarrow$ normal distribution:

- never touches $x$-axis.
- $P(x \le a) = \int_{-\infty}^{a} f(x) \cdot dx \rightarrow$ st. norm. cdf $(a, \mu, \sigma)$

- $P(x \ge a) = \int_{a}^{+\infty} f(x) \cdot dx \rightarrow$ st. norm. SF $(a, \mu, \sigma)$

- $P(a \le x \le b) = \int_{a}^{b} f(x) \cdot dx$

3 ways : ① st. norm. cdf $(b, \mu, \sigma)$
- st. norm. cdf $(a, \mu, \sigma)$

② st. norm. SF $(a, \mu, \sigma)$ - st. norm. SF $(b, \mu, \sigma)$

③ $1 - [$ st. norm. cdf $(a, \mu, \sigma) +$ st. norm. SF $(b, \mu, \sigma)]$

$z:$
$\rightarrow f(x) = \dfrac{1}{\sigma \sqrt{2\pi}} \, e^{-\frac{1}{2} \left( \frac{x - \mu}{\sigma} \right)^2}$
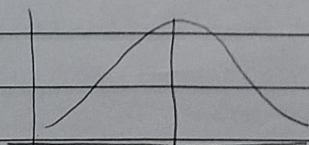
- $z = \dfrac{x - \mu}{\sigma}$

- no units, units cancelled
- if normal distribution followed. $\rightarrow -4 \le z \le +4$
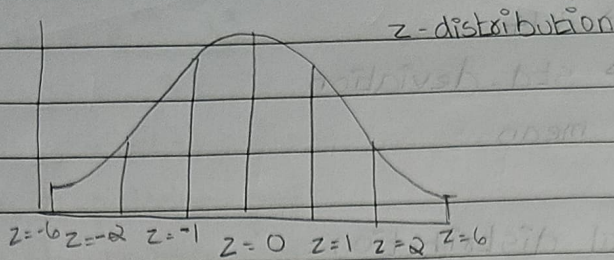- standard normal distribution.

$z = \dfrac{x - \mu}{\sigma}$

For $z = x$ when $\boxed{\mu = 0, \sigma = 1}$

mean $= 0$

$z = 0$

z/≠1



z-distribution

z=-6 z=-2 z=-1  z=0  z=1  z=2 z=6

(if variable follow normal dist.)

$-1 \leq z \leq +1 \Rightarrow 68.2\%$ (area shaded)

$-2 \leq z \leq +2 \Rightarrow 95.4\%$

$-3 \leq z \leq +3 \Rightarrow 99.72\%$      area not shaded
$3.14 \times 10^{-6}$

$-6 \leq z \leq +6 \Rightarrow 99.99999\,6\%$     $= 0.00000314$

if  mean = median = 0. } to determine if it
     k = 0              } is normal dist.
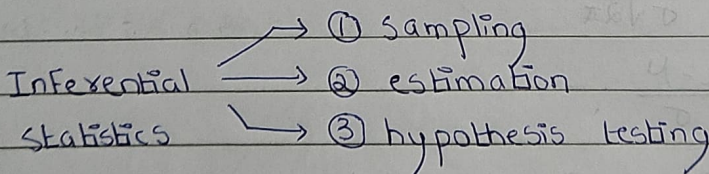
$= 1 - 0.00000314$

= area of graph
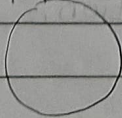
→ IMP functions to use:

① cdf - cumulative, probability - ascending
                                    (lowest to given pt)

② SF - cum. prob. - descending (given pt to highest)

③ isf - inverse survival function (inverse of SF)

④ ppf - point percentile function (inverse of cdf)

⇒ Inferential Statistics:

                    → ① sampling
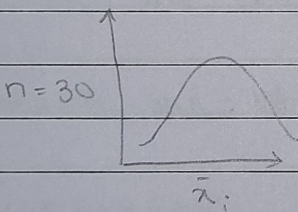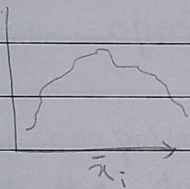   Inferential  ⟶ ② estimation
   Statistics   ↳ ③ hypothesis testing

① Sampling
                population
                      $S_1$ samples
                    → O
                    → O $S_2$

population → parameters
sample → statistic

|  | parameter | statistic |
|---|---|---|
| numerical $\Big\{$ mean | $\mu$ | $\bar{x}$ |
| variance | $\sigma^2$ | $s^2$ |
| SD | $\sigma$ | $s$ |
| categorical $\Big\{$ proportion ① | $P$ | $\hat{p}$ |

$\mu$ and $\bar{x}$ are not same
for all above parameter ≠ statistic

→ suppose, 100 samples
n = sample size
observations
= 15 object





(after n = 30
normal dist. observed)
n ≥ 30 approx. normal
dist.

CLT (Central Limit Theorem)
aka law of large numbers



CLT :

i) $n \geq 30$ - approx. normal. dist.

ii) $\bar{x} \approx \mu$

iii) S.D. of sample mean $\approx \dfrac{\sigma}{\sqrt{n}}$ → std. error

iv)

## Sampling

**Probabilistic / Stochastic**

① Scientific

a) simple random sampling

b) Stratified random sampling

c) Cluster sampling

d) systematic sampling

**non - Probabilistic (not used)**

① non - scientific

a) comfortable and convinient

b) quota sampling

c) judgemental sampling biased

d) snowball sampling