

Statistics

29 AUG 23

The main purposes of statistics are classified under three headings.

- Description
- Analysis
- Prediction.

Individual data are not important but they are considered as a means to measure a certain physical property of interest, the test of hypothesis on the prediction of future occurrences under given conditions.

Whatever the final objective of the experiment statistical methods of "inductive inference" in which a particular set (or sets) of data - "the so-called realization of the samples" is used to draw inferences of general nature on a "population" under study.

Inductive inference drawn from incomplete information may be wrong even if the original information is not. In the field of statistics, the positivity is often related to the process of data collection on the other hand to the fact that we can only make probabilistic statement predictions. It is evident that insufficient or biased data or/and the failure to consider an important influencing factor in the experiment may lead to incorrect conclusion.

5 essential steps to solve problems: (in industry)

1. State the problem on question.
2. Collect and Analyze data
3. Interpret data and make decision.
4. Implement and verify decision.
5. Plan ~~and~~ next action.

examples

Ex: 18

Data \rightarrow Qualitative (categorical)

\downarrow Quantitative (numerical)

① Data in the form of numerical measurements or counts are referred to as quantitative.

② Data in the form of classification into different groups or categories are referred to as qualitative data.

Quantitative \neq Qualitative: A convert अच्छा थाएँ, लेटे 70-80% रुपए Good weather. After Quantitative \neq Qualitative & convert अच्छा

Representation: ① Graphical

②

Measure of centre:

Mean / Average / ~~Average~~: frequent, ~~average~~ distribution

Median: मध्यक: ascending order अनुसारे middle point वर्तमान median

$$l = \frac{(n+1)}{2} \quad (\text{the dataset should be in ordered})$$

$$\bar{x}_m = \frac{n_1 \bar{x}_1 + n_2 \bar{x}_2}{n_1 + n_2}$$

Ex: 18 (continued)

Relative Humidity of 30 days of a

mean: 81
median: 80.5
Mode: 78

(05 Sep 23)

⑥

60	63	64	71	67	73	79	80	83	81
86	90	96	98	98	99	85	80	77	78
71	79	76	84	85	82	90	78	79	79

Taking 2nd Table:

X	f _i
77	1
78	1
80	1
86	1
89	1
90	1
96	1
98	2
98	1
99	1

range = (99 - 77) = 22

Mean = $\frac{\sum f_i x_i}{n}$

= $\frac{10}{10} = 10$.

Median = $\frac{89 + 91}{2} = 90$

Mode = 98.

X	f _i
77	1
78	1
80	1
86	1
89	1
90	1
96	1
98	2
99	1

$$\text{Mean} = \frac{\sum f_i x_i}{\sum f_i}$$

Range & Highest - Lowest
class interval → ✓

$$\text{Variance} = \frac{\sum (x_i - \bar{x})^2}{n-1}, \text{ or, } \sigma^2, \text{ or, } S^2, \rightarrow \text{static Deviation.}$$

Grouped
Data

class	frequency	midpoint	$f_i x_i$	Cumulative frequency	$(x_i - \bar{x})^2$	$f_i (x_i - \bar{x})^2$
60-69	3	62	186	3	330.15	990.45
65-69	1	67	67	4	173.45	173.45
70-79	1	72	288	5	66.92	256.68
75-79	7	77	539	12	9.92	69.44
80-89	6	82	492	18	3.35	20.1
85-89	3	87	261	21	16.65	139.95
90-91	2	92	184	23	139.95	279.9
95-99	1	97	388	30	283.25	11.33
	$\eta = 30$		$\sum f_i x_i = 2905$			$\sum f_i (x_i - \bar{x})^2 = 3062.97$

Mean. $\bar{x} = 8.17$, Mode 79, Median 80,

$$(\bar{x} \text{ or } \mu) \text{ Mean} = \frac{\sum_{i=1}^n f_i x_i}{n}$$

$$\bar{x} \text{ or } \mu = \frac{2105}{30} = 80.167 \approx 80.17$$

Modal class $\rightarrow (75-79)$ कार्यकरण में class व max frequency 7।

\therefore Mode : 79 (frequency 4).

$$\text{Estimated Median} = L + \left[\frac{\frac{n}{2} - Cf}{f} \right] w.$$

L = lower class boundary of the group

Cf = cumulative frequency

f = frequency

w = width

n = total no. of values

$$75 + \left[\frac{\frac{30}{2} - 8}{7} \right] \times 5$$

$$= 80$$

Estimated Mode = L +

$$\frac{(f_m - f_{m-1})}{(f_m - f_{m+1}) + (f_m - f_{m-1})} \times w$$

$$= (75 + \frac{7-4}{(7-4)+(7-6)} \times 5)$$

w = group width

L = lower boundary of the modal class.

f_m = frequency of the group

f_{m+1} = frequency of the group after modal group.

f_{m-1} =

before n =

$$L = 75$$

$$f_m = 7$$

$$f_{m+1} = 4$$

$$f_{m-1} = 9$$

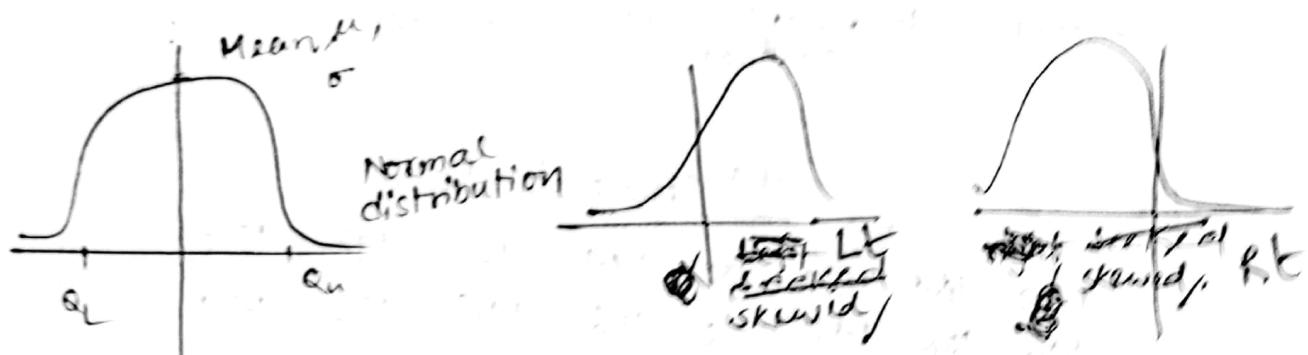
Q3 (probable) & Q1 (probable), Q1-Q3 = 7.50

mean square deviation.

(12 SEP 23)

$$\text{Variance, } \sigma^2 = \frac{\sum_{i=1}^n f_i(x_i - \bar{x})^2}{(n-1)} = \frac{3062.97}{30-1} = 103.62$$

$$\text{Standard Deviation: } \sqrt{s^2} = \sigma = 10.2$$



$$\bar{y} = a\bar{x} + b \quad \text{and} \quad s_y = \sigma_x s_a$$

(13 SEP 23)

~~Empirical Rule:~~

For any mound shaped, nearly symmetric distribution of data

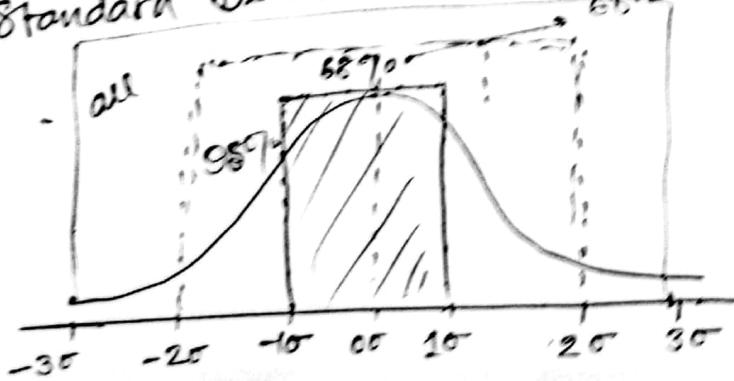
interval

- The $\bar{x} \pm s$ contains 68% data points. $\rightarrow 68.2607\%$
- The interval $\bar{x} \pm 2s$ contains 95% data $\rightarrow 95.449\%$
- $\sim \sim \sim \bar{x} \pm 3s \rightarrow$ all data $\rightarrow 99.72\%$

Standard Normal Distribution:

Mean $\rightarrow \mu = 0$

Standard Deviation $= 1$ (1)



(19 Sep 23)

Standard Values and Z-scores:

$$z\text{-score} = \frac{\text{measurement} - \text{mean}}{\text{standard deviation}}$$

$$z\text{-score} = \frac{189.966.99 - 37298.03}{84.369} = 1.81$$

40 students Data

138	169	180	132	144	125	199	157
146	158	190	147	136	148	152	149
168	126	138	176	163	119	159	165
146	173	192	197	135	153	190	135
161	195	135	142	150	156	195	128

Largest weight = 176 lb

Smallest \rightarrow = 119 lb

$$\text{Range} = (176 - 119) = 57 \text{ lb.}$$

$$\frac{57}{5} \approx 11 \rightarrow \text{(total division)}$$

interval

$$\frac{57}{9} \approx 7$$



(if q)

f_i

~~118 - 122~~

123 - 127

128 - 132

133 - 137

138 - 142

143 - 147

152 - 156

162 - 166

162 - 166

(f_i) frequency

x_i Midpoint

118 - 122 → 3
127 - 135 → 5
136 - 144 → 9
145 - 153 → 12
154 - 162 → 5
163 - 171 → 1
172 - 180 → 2
 $\frac{70}{70}$

118 - 122	1
123 - 127	2
128 - 132	2
133 - 137	1
138 - 142	6
143 - 147	8
148 - 152	9
153 - 157	1
158 - 162	2
163 - 167	3
168 - 172	1
173 - 178	2

120

125

130

135

140

145

150

155

160

165

170

175

Ex: 1.1 (Pg. 51 of 83)

#

$$\textcircled{a} \quad y_i = x_i + 50; \quad \bar{y}_i = \bar{x}_i + 50. = 550 \text{ \$}$$

standard deviation unchanged.

$$\textcircled{b} \quad y_i = 1.10x_i; \quad \bar{y} = 1.10\bar{x} = 1.10(500) = 550 \text{ \$}$$

standard deviation, $S_y = \sqrt{1.10^2 S_x^2}$

$$= \sqrt{(1.10)^2 (125)^2}$$

$$= 137.5 \text{ } \underline{125}$$

Given $f = 500$
 $f = 125$
 $\textcircled{1} \quad \text{increase price of each by 50\$}$
 $\textcircled{2} \quad \text{each by 100\$}$
 $\textcircled{3} \quad \text{increase by 50\$}$

std variance! $\sqrt{s^2} = \sigma$

variance $s^2 = \sigma^2$

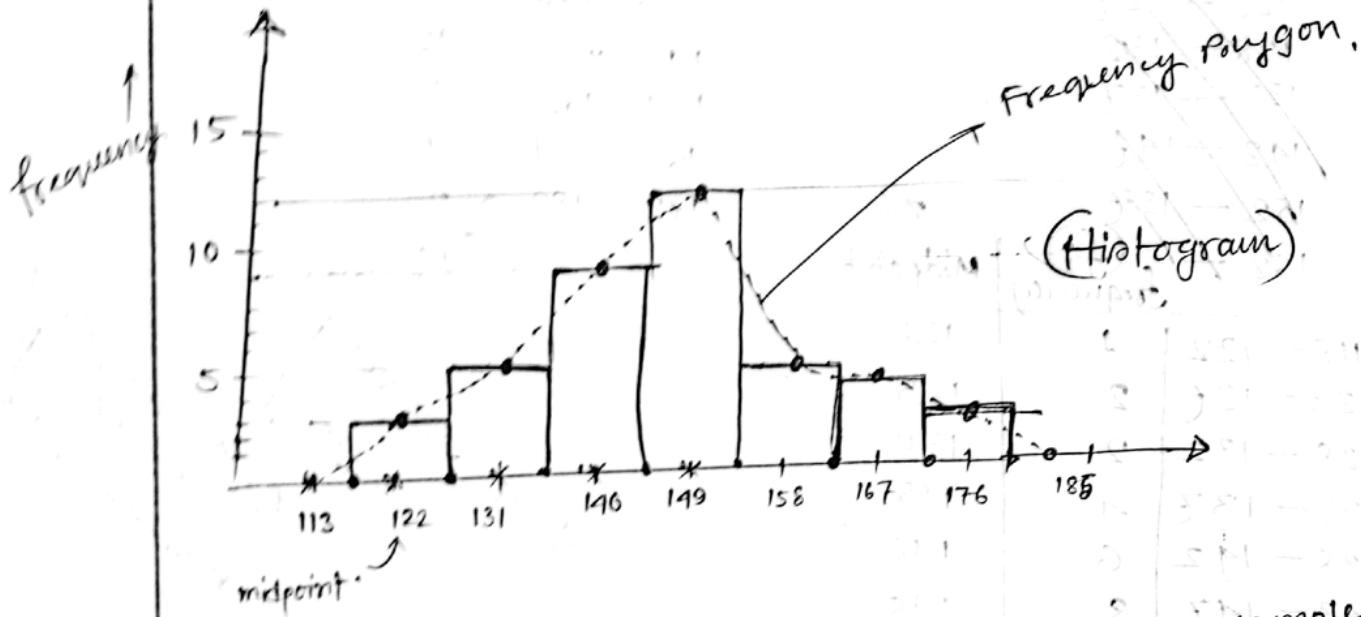
$\bar{y} = a\bar{x} + b$
 $s_y^2 = a^2 s_x^2$

s_x^2, a^2

Histogram & First & Last interval

(26 Sep 23)

Frequency Polygon: Midpoint Required
Histogram : First & Last point.



Chebyshev's Inequality

Let \bar{x} and s be the sample mean and standard deviation of data set, assuming that $s > 0$, Chebyshev's inequality states that for any value of $K \geq 1$, greater than $100(1 - \frac{1}{K^2})$ percent of the data lies within the interval from $\bar{x} - ks$ to $\bar{x} + ks$.

$(\bar{x} - 1.5s, \bar{x} + 1.5s)$ interval \approx ? % data

$$\text{Suppose } \frac{13}{2} = K$$

$$100 - \left(1 - \frac{1}{9}\right) =$$

$$= 100 \left(\frac{5}{9}\right)$$

= ≈ 55.56 percent of data lie within $(\bar{x} - 1.5s, \bar{x} + 1.5s)$

Freq
Value

Let $S_K = \sum f_i k_i^2$

$$\frac{N(S_K)}{n} > 1 - \frac{n-1}{nK^2} > 1 - \frac{1}{K^2}$$

$N(K) = \text{Number of } i : |x_i - \bar{x}| > Ks$

$$\frac{N_K}{n} < \frac{1}{K^2}$$

$$|x_i - \bar{x}| < Ks$$

$$\bar{x} - Ks < x_i < \bar{x} + Ks$$

(3 OCT 23)

Three basic statistics used in inferential statistics

Variance

$$S^2 = \frac{\sum f_i (x_i - \bar{x})^2}{n-1}$$

পথন অসমৰ্গত data
এবং frequency ।
এবং data আলাদা

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

CT
Exam
Question Review

Estimated

$$n = 64$$

$$\therefore \frac{n}{2} = \frac{64}{2} = 32$$

\therefore Median = $\frac{32^{\text{nd}} + 33^{\text{rd}} \text{ entry}}{2}$

$$= \frac{5+5}{2} = 5.$$

$$\text{class width} = 2000$$

$$2 \rightarrow 3 = \frac{2}{2} = 1$$

$$2 \rightarrow 1 = \frac{2}{4} = 0.5$$

32nd & 33rd entry

32nd & 33rd position

32nd & 33rd position

B1

The following data represents the lifetimes (in hrs) of a sample of 10 transistors.

112	121	126	108	111	109	136	139
121	118	143	116	108	122	127	140
113	117	126	130	134	120	131	133
118	125	151	147	137	140	132	119
110	129	132	152	135	130	136	128

(a) Determine the sample mean, median and mode.

(b) Draw cumulative relative frequency plot.

(c) Are the data apparently normal?

(d) What percentage of data fall within $\bar{x} \pm 1.5\sigma$?

(e) Compare the result in (d) to that given by empirical rule.

(f) Use Chebychev's inequality.

Soln: Mean = 127.125

Median = 127.5

Mode = 108

Standard Deviation = 11.873

Range = 152 - 109 = 48



$$\text{formula} = 100 \left(1 - \frac{1}{K^2}\right)$$

$$68\% \cdot (\bar{x} - ks, \bar{x} + ks)$$

	freq	relative freq	cumulative freq
109 - 116	7	$\frac{7}{40}$	7
117 - 128	19	$\frac{19}{40}$	20
129 - 140	19	$\frac{19}{40}$	39
141 - 152	5	$\frac{5}{40}$	40
	$n=40$		

chebychev's inequality yields that $100 \left(1 - \frac{1}{(\frac{3}{2})^2}\right)$
 $\Rightarrow 100 \times \frac{5}{9} = 55.55$

percent of data lies in the interval

$$(\bar{x} - 1.5s, \bar{x} + 1.5s) = \{(127.425 - 1.5 \times 11.873), (127.425 + 1.5 \times 11.873)\} \\ = (109.615, 145.235)$$

standard deviation

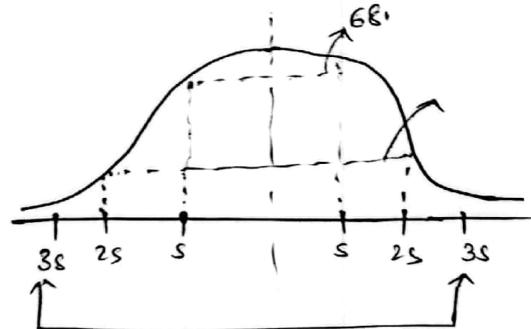
\xrightarrow{s} Chebychev's band:

$$\frac{N(s_{1.5})}{40} \geq 1 - \frac{n-1}{nK^2} \geq 1 - \frac{1}{K^2} = \frac{5}{9}$$

$\frac{N(s_k)}{n} \geq 1 - \frac{n-1}{nK^2} > 1 - \frac{1}{K^2}$

$$1 - \frac{40-1}{40 \cdot (\frac{3}{2})^2} \geq \frac{5}{9}$$

$$\therefore \frac{51}{90} > \frac{5}{9}.$$



95.44 %