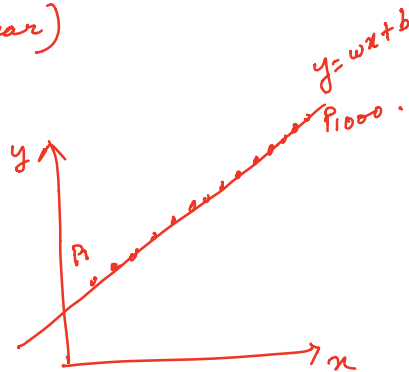
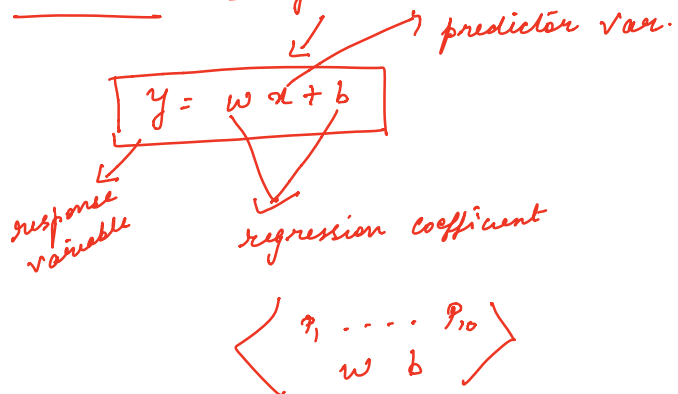


22/01/21 Dec-11

Numerosity redn.
 \swarrow parametric
 \searrow non-parametric

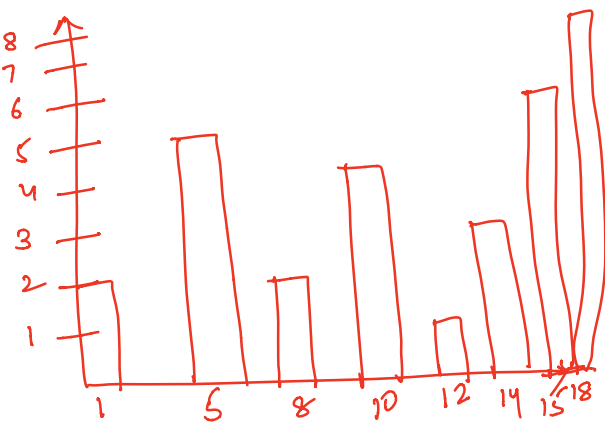
→ parametric (Regression and loglinear)



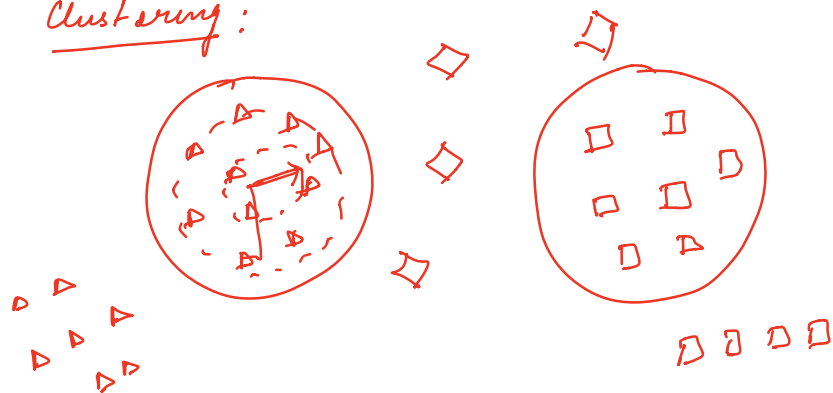
Non-parametric approaches:

1) Histogram:

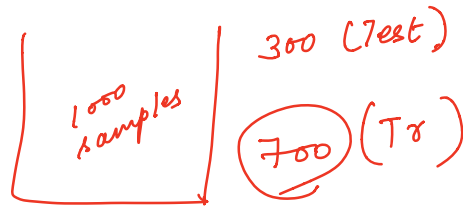
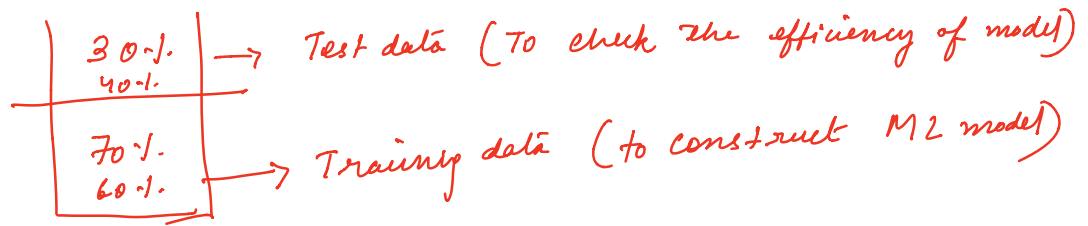
- | | |
|-----------------------|-------------------------|
| $f(1) \rightarrow 2$ | $f(12) \rightarrow 1$ |
| $f(5) \rightarrow 5$ | $f(14) \rightarrow 3$ |
| $f(8) \rightarrow 2$ | $f(15) \rightarrow 0.6$ |
| $f(10) \rightarrow 4$ | $f(18) \rightarrow 0.8$ |



2) Clustering:



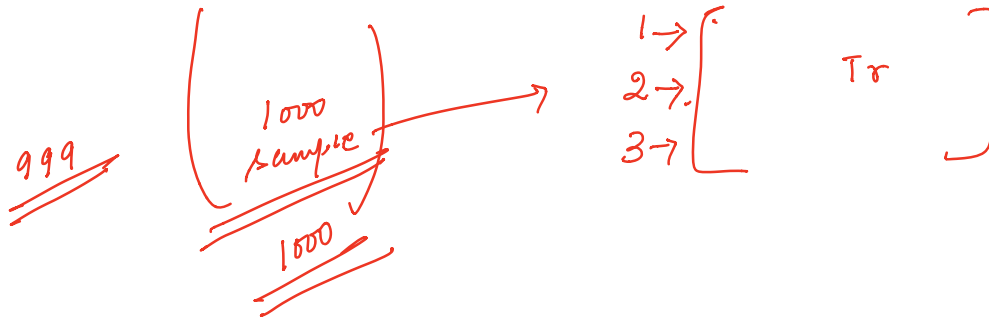
3) Sampling: a) simple random sample without replacement (SRSWOR):



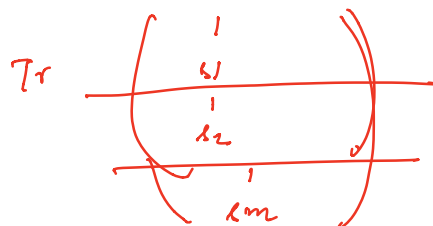
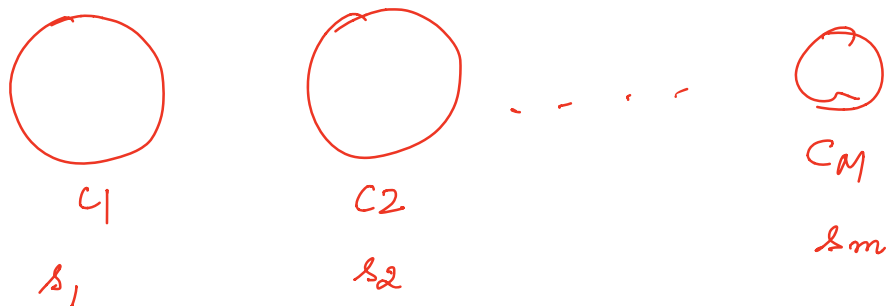
prob. of a tuple getting selected = $\frac{1}{N}$

$= \frac{1}{1000}$

b) SRSWR: (simple random sample with replacement)



c) cluster sample:



d) Stratified sample:

T ₁	Youth
T ₂	"
T ₃	"
T ₄	"
T ₅	Middle aged
T ₆	"
T ₇	"
T ₈	"
T ₉	"
T ₁₀	senior
T ₁₁	"

T ₂	Youth
T ₃	"
T ₇	Middle aged
T ₈	"
T ₁₀	senior
T ₁₁	"

4) data cube aggregation:

Year 2010	
Year 2009	
Year 2008	
Q ₁	\$224
Q ₂	\$408
Q ₃	\$350
Q ₄	\$586

→

Year	sales
2008	\$1568
2009	\$2356
2010	\$3594

Normalization:

The attribute data are scaled so as to fall within a small specified range such as -1 to +1 or 0 to 1.

3 ways for data normalization:

1) min-max normalization:

(performs a linear transformation on the original data)

$$V' = \frac{V - \min_A}{\max_A - \min_A} (\text{new max}_A - \text{new min}_A) + \text{new-min}_A$$

Eg: min-income: \$12000

max-income: \$98000

Map to $[0, 1]$ a income value of \$73,600.

Soln:
$$V^1 = \frac{73,600 - 12,000}{98,000 - 12,000} (1 - 0) + 0$$

$$= \underline{\underline{0.716}}$$

2) Z-score normalization: (or zero mean normalization)
The values for an attribute A are normalized based on the mean and std. deviations of A .

$$V^2 = \frac{V - \bar{A} \rightarrow \text{mean}}{\sigma_A \rightarrow \text{std. devn.}}$$

Eg: Let mean and std. dev. of attribute income are \$54,000 and \$16,000 resp.

Using Z-score normalization, find the normalized value for \$73,600.

Soln:
$$V^1 = \frac{73,600 - 54,000}{16,000} = 1.225$$

3) Decimal scaling:

This normalizes by moving the decimal points of values of attribute A . The no. of decimal points moved depends on the maximum absolute value of A .

$$V^1 = \frac{V}{10^j}$$

where j is the smallest integer
s.t. $\text{Max}(|V^1|) < 1$

Eg: Suppose the recorded values of A range from -986 to 917 .

The max. abs. value of A is 986 .

To normalize by decimal scaling, we divide the value by 10^3 (ie 1000).

\therefore Normalized value: -0.986 and 0.917 .

* Attribute construction:

$$\begin{pmatrix} \text{ht} & \text{width} \\ f_1 & f_2 \end{pmatrix} \quad \begin{matrix} \downarrow \\ f_3 \end{matrix} \rightarrow \text{area}$$

New attributes can be constructed from the given attributes and added encoder to improve the accuracy and understanding of high dimensional data.