

20/2/23
Monday

Module - 2

Data preprocessing

3 types of data

Incomplete

Noisy

Inconsistent

why data preprocessing

• Real world data

• Incomplete - malfunctioning of data.

• Noisy - unwanted data.

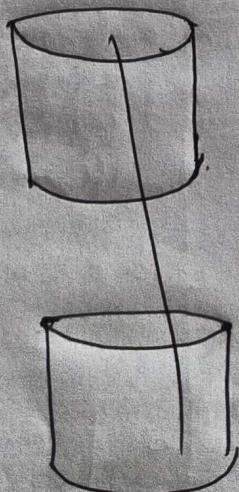
• Inconsistent

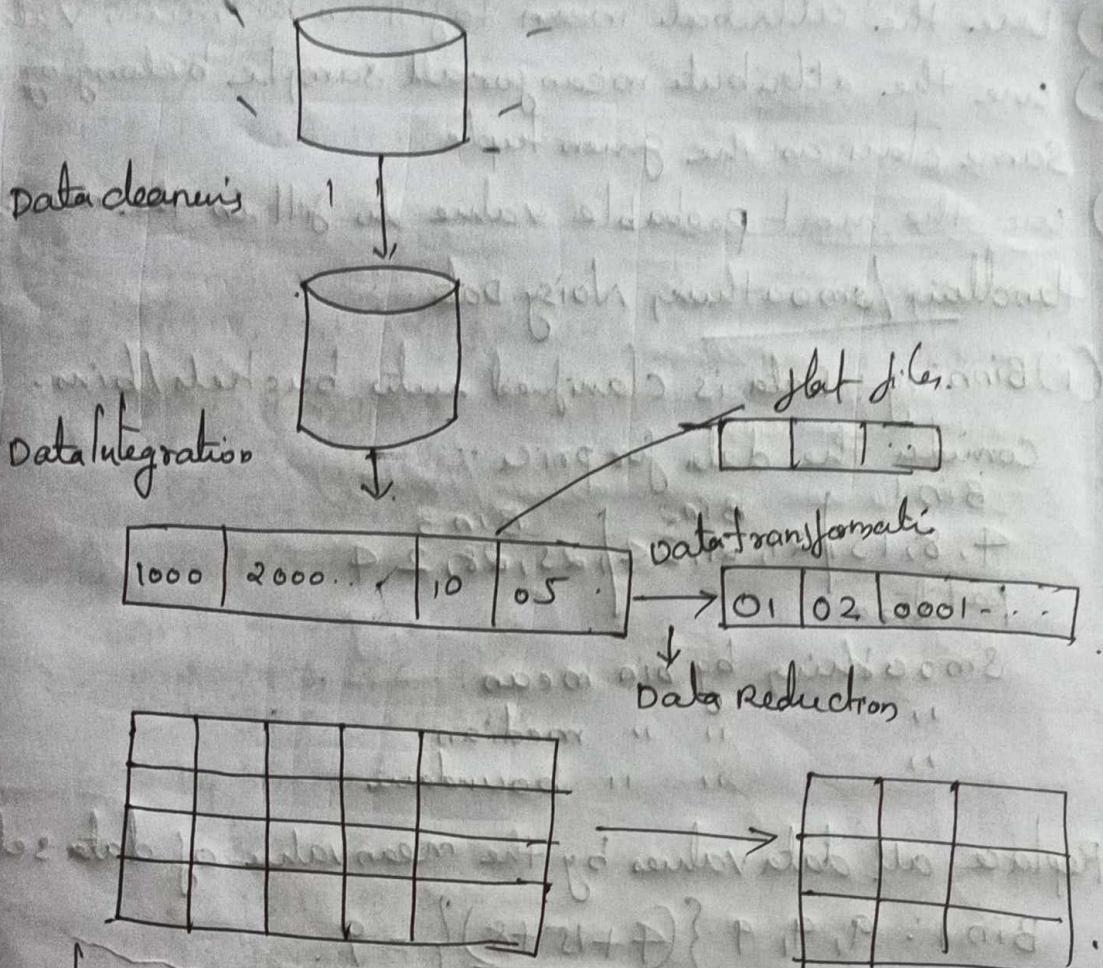
↓
Naming conventions : same data and different column names

Normalization

Major steps in data preprocessing

- (i) Data cleaning
- (ii) Data integration
- (iii) Data transformation
- (iv) Data reduction





Data cleaning : Remove noisy data.
filling missing.

Data transformation : To convert into common format
(normalization).

Data reduction : To reduce its size after data transformation.

23/2/23.

Thursday

Data cleaning : Filling the missing value, smoothing noisy data
identifying or remove outliers, resolving inconsistencies

I

Handling missing values

1) Ignore the tuple.

Method is not very effective

2) Filling is the missing values manually.

This approach is time consuming, may not be feasible.

3) use a global constant to fill is the miss value.

- 4) Use the attribute mean to fill in the missing value
- 5) Use the attribute mean for all samples belonging to H.
- 6) Same class as the given tuple.
- 7) Use the most probable value to fill in the missing value.

II

Handling / smoothing Noisy Data

(i) Binning: Data is classified into buckets/bins.

Consider the data for price.

4, 8, 15	Bin 2	Bin 3
21, 21, 24	25, 28, 24	

Smoothing by bin mean.

" " median,

" " boundaries,

Replace all data values by the mean value of data set each bin.

$$\text{Bin 1: } 9, 9, 9 \left\{ \frac{(4+15+8)}{3} \right\} = 9.$$

$$\text{Bin 2: } 21, 21, 24 = 22.$$

$$22, 22, 22$$

$$\text{Bin 3: } 25, 28, 24.$$

$$29, 29, 29.$$

Smoothing by median.

$$\text{Bin 1: } 8, 8, 8.$$

$$\text{Bin 2: } 21.$$

$$\text{Bin 3: } 28.$$

$$\text{Bin 1: } 8, 8, 8.$$

$$\text{ " 2: } 21, 21, 21.$$

$$\text{ " 3: } 28, 28, 28.$$

Smoothing by bin boundaries.

$$4, 4, 15$$

$$21, 21, 24$$

$$25, 28, 34.$$

Apply smoothing techniques by min mean, median, binning, boundaries on the following data.

13, 7, 24, 50, 3, 9, 11, 12, 34, 26, 17, 23, 1, 8, 30
binsize = 5.

sort

12, 3, 7, 8 / 9, 11, 13, 17, 23, 24, 26, 30, 34, 50.

Bin 1: 1, 2, 3, 7, 8

" 2: 9, 11, 13, 17, 23.

" 3: 24, 26, 30, 34, 50.

smoothing by bin means.

Bin 1: 4.2, 4.2, 4.2.

" 2: 14.6, 14.6, 14.6

" 3: 32.8, 32.8, 32.8

smoothing by bin median.

Bin 1: 3, 3, 3, 3, 3.

" 2: 13, 13, 13, 13, 13

Bin 3: 30, 30, 30, 30, 30.

smoothing by bin boundaries

Bin 1: 1, 1, 1, 8, 8

Bin 2: 9, 9, 9, 23, 23

Bin 3: 24, 24, 24, 50, 50.

Regression (prediction)

It is also used for noise recovery.

Types of regression

Simple linear.

Multiple

Polynomial

Histogram

Data cleaning as a process.

Step 1: Discrepancy detection

2: Data transformation.

Discrepancy

Human error in data entry

25/12/23
Saturday

Avoid discrepancy

1) Using metadata.

Commercial tools.

Data scrubbing tools.

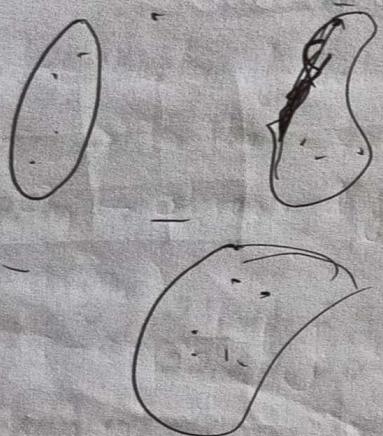
Data auditing tools.

Data transformation

Data Migration tool

ETL (Extraction / Transformation / Loading).

clustering



Data Integration

The data integration is the process of combines data from multiple sources into a coherent data store.

Issues in data integration

1. Entity identification problems

2. Data Redundancy problems.

3. Detection and resolution of data value conflicts.

1. Entity identification problem

Eg: Data analyst or computer cannot ensure that customer-id in one database and cust-number in another refer to the same attribute, due to metadata.

helps to avoid errors in schema entity identifies problems
Data Redundancy (repetitiveness of data)

using normalization to avoid redundancy. Derived attribute
are also examples of redundancies.

Solution to data Redundancy

b) Correlation Analysis

For numerical attributes we can evaluate the correlation b/w
2 attributes A & B. Known as Pearson's product moment coefficient
(named after its inventor, Karl Pearson)

$$\gamma_{A,B} = \frac{\sum_{i=1}^N (a_i - \bar{A})(b_i - \bar{B})}{N \sigma_A \sigma_B}$$

$$= \frac{\sum_{i=1}^N (a_i b_i) - N \bar{A} \bar{B}}{N \sigma_A \sigma_B}$$

N - No. of tuples.

a_i, b_i - respective values of A & B in tuple i

\bar{A} & \bar{B} - respective mean values of A & B

$\sigma_A \sigma_B$ - respective standard deviation.

A, B - attributes.

$\gamma_{A,B}$ - correlation coefficient.

$-1 \leq \gamma_{A,B} \leq 1$.

{ if $\gamma_{A,B} > 0$, A & B are +vely correlated. which means the value
of A increases as the values of B increase.
Highly Correlated
if $\gamma_{A,B} = 0$, A & B are independent & there is no correlation b/w
them.

If $r_{A,B} < 0$

A & B are negatively correlated, where the values of one attribute increase, the values of the other attribute decrease.

If $r_{A,B} \neq 0$ A & B dependent to each other.

Correlation analysis for categorical attribute

A correlation relationship b/w 2 attributes A, B can be discovered by a χ^2 (Chi-square) test. Suppose A has c distinct values, namely a_1, a_2, \dots, a_c . B has r distinct values, namely b_1, b_2, \dots, b_r .

$$\chi^2 = \sum_{i=1}^c \sum_{j=1}^r \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

The data tuples described by A & B can be shown as a contingency table, with the c values of A

O_{ij} - observed frequency (actual count) of the joint event (A_i, B_j)

E_{ij} - expected frequency of (A_i, B_j)

The sum is computed over all of the

$$E_{ij} = \text{count}(A=a_i) \times \text{count}(B=b_j)$$

Problem of categorical analysis

- 1) Suppose that a group of 1,500 people was surveyed. The gender of each person was noted. Each person was polled as to whether their preferred type of reading material was fiction or non-fiction. Thus, we have 2 attributes, gender and preferred reading. The observed frequency (or count) of each possible joint events is summarized in the Contingency, where the numbers in parentheses are the expected frequencies.

		Male	Female	
b ₁	Fiction	250 (1,1)	200 (2,1)	450
b ₂	Nonfiction	50 (1,2)	100 (2,2)	1050
		300	1200	1500

Hypothesis

Gender and preferred reading are independent each other.

B_{ij} = Total cost

(0 + 0 + 0001)

$$e_{ij} = \frac{(a_{ij} \times b_{ij})}{N}$$

$$e_{11} = \frac{a_1 \times b_1}{N}$$

$$= \frac{300 \times 450}{1500} = \underline{\underline{90}}$$

$$e_{21} = \frac{1200 \times 360}{1500} = \underline{\underline{360}}$$

$$e_{12} = \frac{a_1 \times b_2}{N}$$

$$= \frac{300 \times 1050}{1500} = \underline{\underline{210}}$$

$$e_{22} = \frac{a_2 \times b_2}{N} = \frac{1200 \times 1050}{1500} = \underline{\underline{840}}$$

	a ₁	a ₂	Total	
	Male	Female		
b ₁	Fiction	250 (1,1) (90)	200 (2,1) (360)	450
b ₂	Nonfiction	50 (1,2) (210)	100 (2,2) (840)	1050

$$\chi^2 = \sum_{i=1}^r \sum_{j=1}^c \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

$$= \frac{(250-90)^2}{90} + \frac{(200-360)^2}{360} + \frac{(50-210)^2}{210} + \frac{(1000-840)^2}{840}$$

$$= 284.44 + 71.11 + [21.90 + 30.48]$$

$$= \underline{\underline{507.93}}$$

$$\text{degrees of freedom} = (r-1)(c-1)$$

$$= (2-1)(2-1) = 1.$$

If the critical value from the chi-square distribution table is checked against the chi-square, if calculated chi-square value < critical value, hypothesis can be accepted.
Otherwise hypothesis rejected.

Critical value corresponding to $\text{df}((r-1)(c-1))$ with maximum significant levels of 0.001 is 10.828.

$507.93 > 10.828$ hypothesis can be rejected. $\therefore A \& B$ are strongly correlated i.e., preferred reading and gender are correlated.

Detection and resolution of data value conflicts.

To avoid this, attributes functional dependencies and referential constraints in the source system match those in the target system.

28/10/2023
Tuesday

Data transformation

The data are transformed into forms appropriate for mining.

Smoothing which works to remove noise from the data. such techniques include binning, regression, clustering

Aggregation : where summary or aggregation are applied to the data.

e.g. the daily sales data may be aggregated so as to compute

monthly and annual total amounts. This step is typically used in constructing a data cube for analysis of the data at multiple

generalization of the data

Generalization of the data, where low-level data are replaced by higher level concepts through the use of concept reaches.

e.g.: categorical attributes, like street, can be generalized to higher level concepts like city or country.

Attribute construction (feature construction), where new attributes are constructed and added from eg: we may be wish to add the attribute area based on the attribute height & width. the given set of attributes to help the mining process.

Normalization: where the attribute data are scaled so as to fall within a small specified range.

Methods for data normalization:

Min max normalization.

Z-score "

Min-max normalization: performs a linear transformation on the original data.

Let \min_A and \max_A are the minimum and max values of an attribute A.

Maps a value, v , of A to v' in the range [new_min_A, new_max_A].

$$v' = \frac{v - \min_A}{\max_A - \min_A} (new_{max_A} - new_{min_A}) + new_{min_A}$$

It will encounter an "out of bound" error if a future input value for normalization fall outside of the original data range for A.

Q: Suppose that the minimum and max values for the attribute income are \$12,000 & \$98,000 respectively. we would like to map income to the range [0.0, 1.0]. By min-max normalization, a value of \$73,600 for income is transformed to

$$v' = \frac{v - \min_A}{\max_A - \min_A} (\text{new } \max_A - \text{new } \min_A) + \text{new } \min_A$$

$$\min_A = 12000$$

disadvantage

$$\max_A = 98000$$

$$\text{new } \min_A = 0.0$$

$$\text{new } \max_A = 1.0$$

$$\text{Given } v = 73600$$

$$v' = \left[\frac{73600 - 12000}{98000 - 12000} (1.0 - 0.0) \right] + 0.0$$

$$= \frac{61600}{86000} \times 1 + 0.0 = \underline{\underline{0.716}}$$

~~Z-score or Z-mean normalization.~~

~~shorthand~~ For eg: consider the data 99200.

$$V' = \frac{99200 - 12000}{98000 - 12000}$$

Z-score or Z-mean normalization

Values for an attribute, A, are normalized based on the mean and standard deviation of A.

Maps a value V , of A to V' by computing

$$V' = \frac{V - \bar{A}}{\sigma_A}$$

where \bar{A} and σ_A are the mean & sd, respectively, of attribute A.

(useful when the actual min & max of attributes are unknown, or when there are outliers that dominate the min-max normalization.)

Suppose that the mean & sd. of the values for the attribute income are \$54,000 & \$16,000 respectively. With z-score normalization, a value of 73.600 for income is too large.

$$V = 73600$$

$$\bar{A} = 54000 \quad \sigma_A = 16,000$$

$$v' = \frac{v - \bar{A}}{\sigma_A}$$

$$= \frac{73600 - 54000}{16000}$$

$$= \frac{19600}{16000} = \underline{\underline{1.225}}$$

Normalization by decimal scaling

Normalizes by moving the decimal point of values of attribute A

The no: of decimal points moved depends on the max absolute value of A

A value of v_i of A is normalized to v'_i by computing.

$$v' = \frac{v}{10^j},$$

where j is the smallest integer such that $\text{Max}(|v'|) < 1$.

Q: Suppose that the recorded values of A range from -986 to 917. The max absolute value of A is 986. To normalize by decimal scaling, we therefore

divide each value by £1,000.8 (i.e., $j=3$).
 So that -986 normalizes to -0.986 and 917
 normalizes to 0.917.

$v' =$ the following
 use the 2 methods to normalize data.

200, 300, 400, 600, 1000

(A) Min-max normalization by setting $\min_A = 0$
 $\max_A = 1$.

(B) Z-score normalization.

$$v' = \frac{v - \min_A}{\max_A - \min_A} [\frac{\text{new max}_A - \text{new min}_A}{\text{old max}_A - \text{old min}_A}] + \frac{\text{new min}_A}{\text{old min}_A}$$

$$\text{new max}_A = 1$$

$$\text{new min}_A = 0$$

$$v = 200$$

$$\min_A = 200, \max_A = 1000$$

$$v' = \frac{200 - 200}{1000 - 200} [1 - 0] + 0 = \underline{\underline{0}}$$

$$v = 300$$

$$v' = \frac{300 - 200}{1000 - 200} [1 - 0] + 0 = \frac{-100}{800} = -0.125$$

$$V' = \frac{V - 200}{1000 - 200} [1 - 0] + 0$$

$$= \frac{200 - 200}{800} = \frac{-200}{800} = +0.25$$

$$V' = \frac{V - 200}{1000 - 200} [1 - 0] + 0$$

$$= \frac{200 - 200}{800} = \frac{-200}{800} = +0.25$$

$$V' = \frac{1000 - 200}{1000 - 200} = \frac{200 - 1000}{800} = \frac{-800}{800} = -1$$

Normalized data = 0, 0.125, 0.25, 0.5, 1

B) Z-score normalization.

$$V' = \frac{V - \bar{A}}{\sigma_A}$$

$$\bar{A} = \frac{200 + 300 + 400 + 600 + 1000}{5}$$

$$= \frac{2500}{5} = 500$$

$$\sigma^2 \text{ variance} = \frac{\sum (x - \bar{x})^2}{n}$$

$$SD = \sqrt{\frac{\sum (x - \bar{x})^2}{n}}$$

$$\sigma^2 \text{ variance} = \frac{1}{5} \left[(200 - 500)^2 + (300 - 500)^2 + (400 - 500)^2 + (600 - 500)^2 + (1000 - 500)^2 \right]$$

$$\begin{aligned}\sigma^2 \text{ variance} &= \frac{1}{5} \left[(200 - 500)^2 + (300 - 500)^2 + (400 - 500)^2 + (600 - 500)^2 + (1000 - 500)^2 \right] \\ &= \frac{1}{5} \left[90000 + 40000 + 10000 + 10000 + 250000 \right] \\ &= 80000\end{aligned}$$

$$SD = \sqrt{80000} = 282.84$$

$$V' = \frac{200 - 500}{282.84} = -1.06$$

$$V' = \frac{300 - 500}{282.84} = -0.7$$

$$V' = \frac{400 - 500}{282.84} = -0.35$$

$$V' = \frac{600 - 500}{282.84} = 0.35$$

$$V = \frac{1000 - 500}{282.84} = 1.76$$

Data Reduction

with summarization - Data The process of retaining relevant attributes from the data set.

Techniques used for data Reduction

- Data cube Aggregation
- Attribute Subset Selection
- Dimensionality reduction
- Numerosity reduction
- Discretization & Concept hierarchy generation

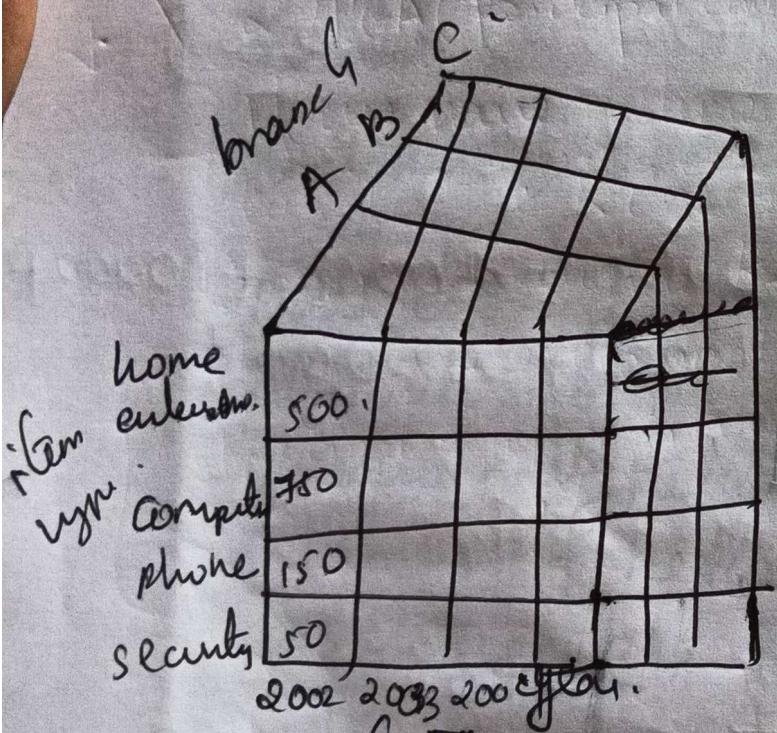
(i) Data cube Aggregation

Year 2002	
Quarter	Sales

year 2004
year 2003
year 2002

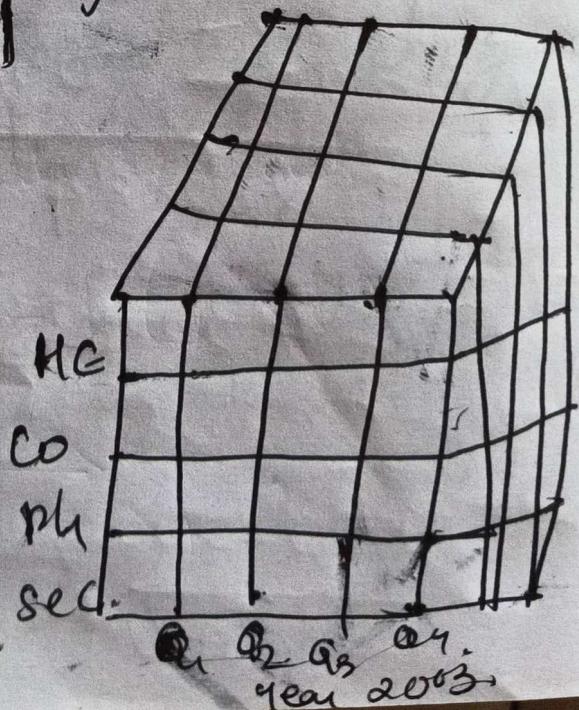
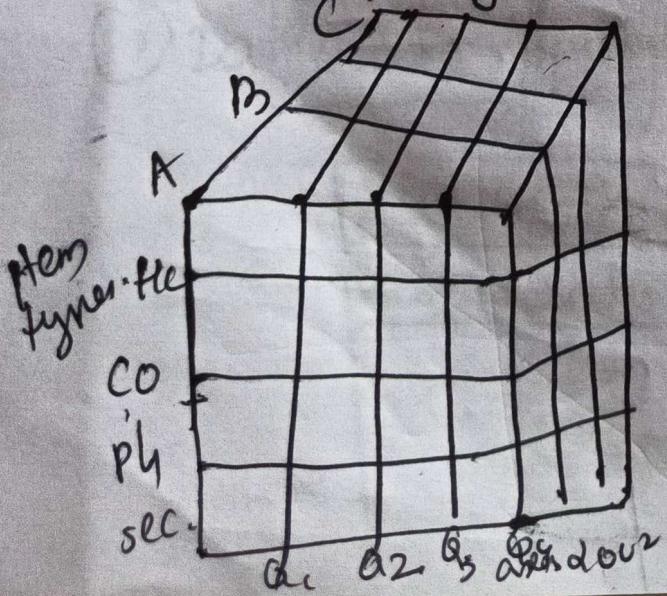
Quarter	Sales
Q1	100
Q2	100
Q3	100
Q4	100

Year	Sales
2002	
2003	
2004	



* Bar cube, * Area cube

showing level abstract
level represented:
highest abstract
level.



13/23 i.f. Dimensionality Reduction

lossless
lossy
(approx 1/ ρ)
dont rec
original

1) DWT (Discrete wavelet transform)

Original data is signal processing efficient.
Replaced by wavelet coefficients.

lossless : If the original data can be reconstructed from the reduced O/P - Then it is known as lossless reduction.

lossy : If the original data can be lost the reconstructed only the approximate i/p.

2) PCA (principal component Analysis).

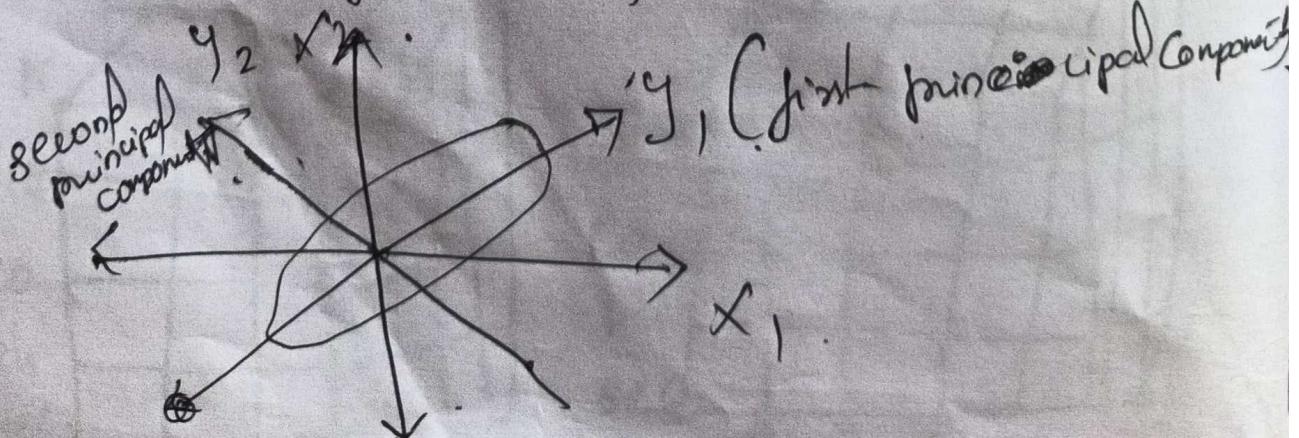
If we have n dimensional data set - we have to calculate k (n-dimensional) principal components where $k \leq n$.

They are orthogonal vectors \perp to original data. Principal components are unit vectors calculated from original data i/p.

$$z^n = \sum$$

$$\frac{1}{|z|}$$

It can be diagrammatically represented as.



Numerosity Red?

to reduce the no. of counts.

(1) parametric methods.

- Regressions.

• log linear models → probability concept

strictly based on parametric
uses some parametric models.

(2) Non-parametric methods.

- Histograms

- clustering.

- sampling.

Histogram

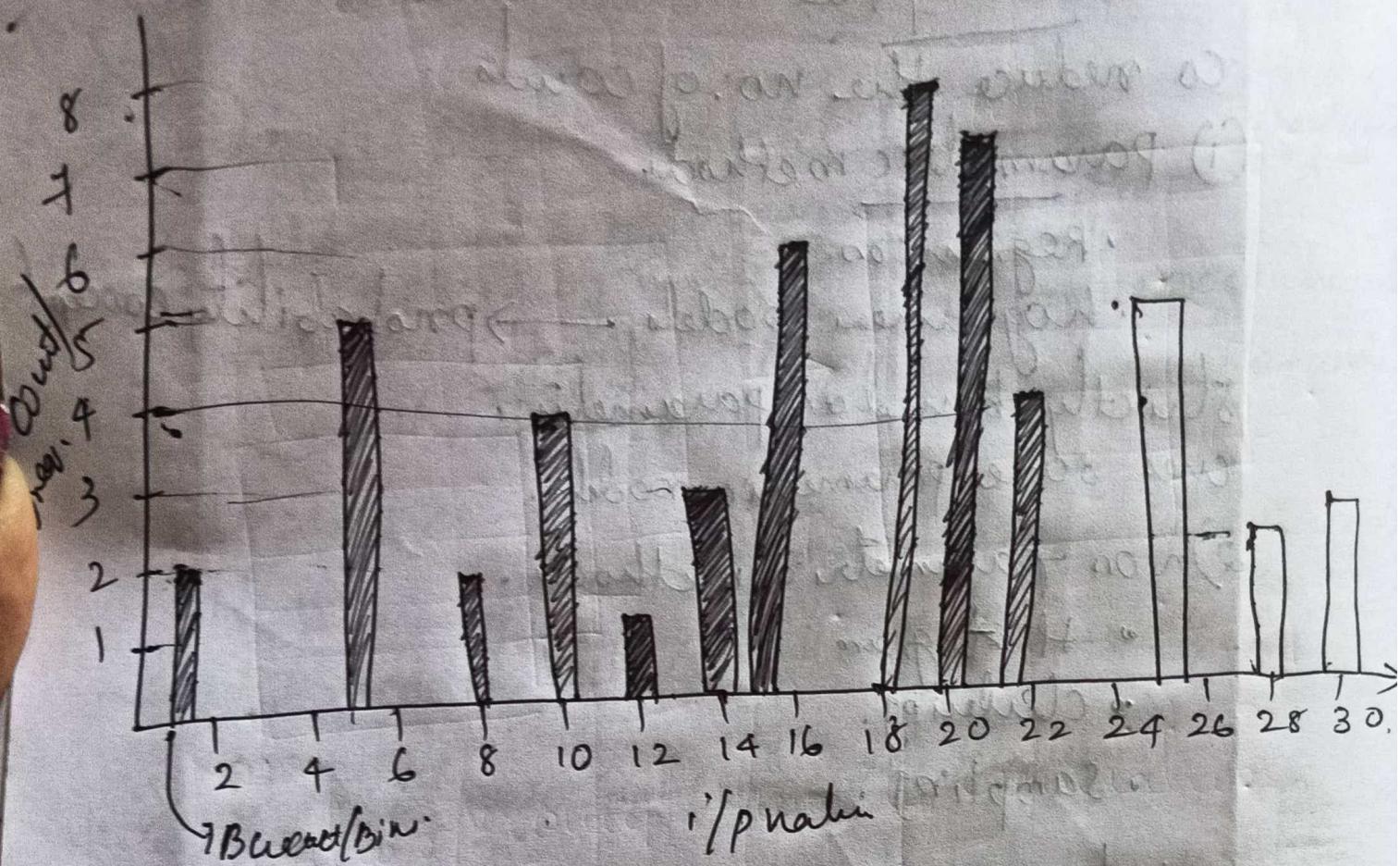
convert the data to a histogram

1, 1, 5, 5, 5, 5, 5, 5, 8, 8, 10, 10, 10, 10, 10, 12, 14, 14, 14

15, 15, 15, 15, 15, 15, 18, 18, 18, 18, 18, 18, 18, 18, 18, 18, 18, 20,

20, 20, 20, 20, 20, 20, 21, 21, 21, 21, 21, 25, 25, 25, 25

28, 28, 30, 30, 30



Process of partitioning histograms:

(1) Equal width.

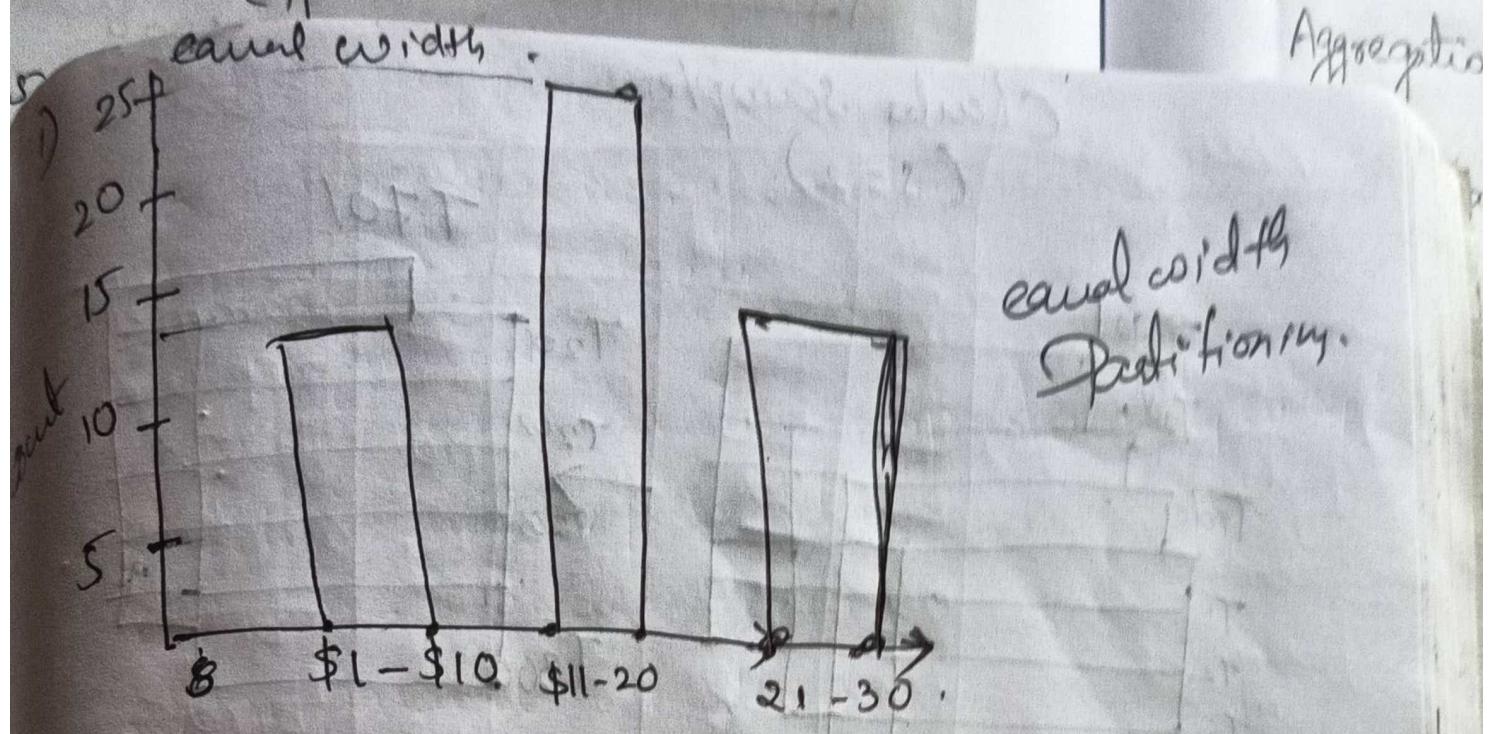
↳ width of each bucket is uniform.

(2) Equal frequency.

frequency of each bucket is a constant
each bucket contains same no. of elements

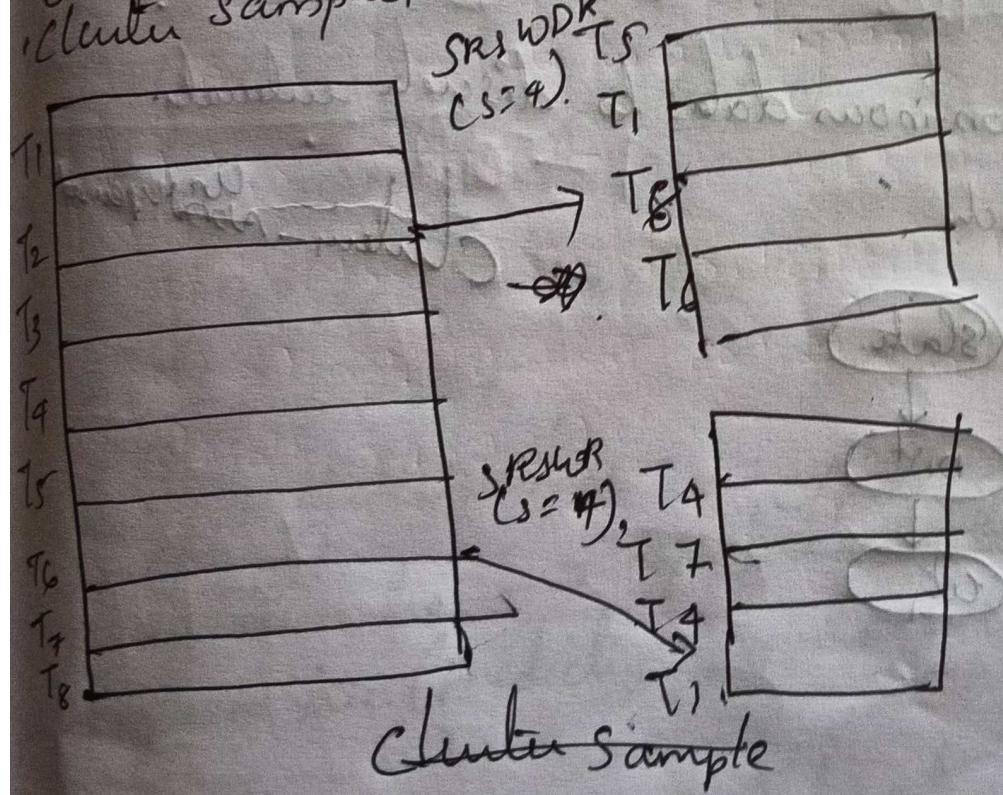
(3) V-optimal → Bucket with least variance

(4) Max Diff - Diff b/w each partition of adj values.

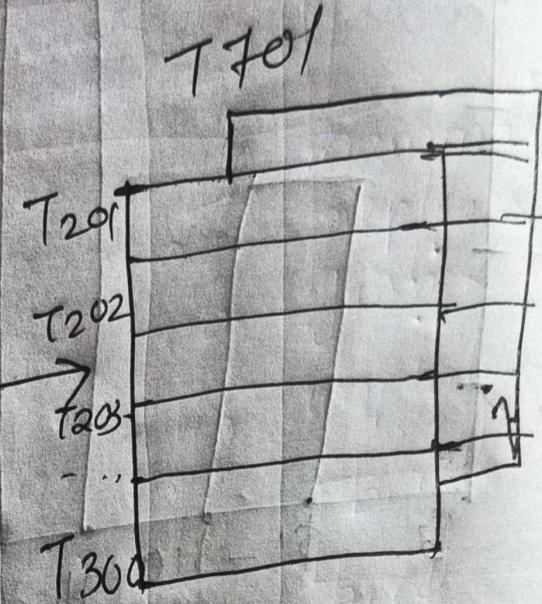
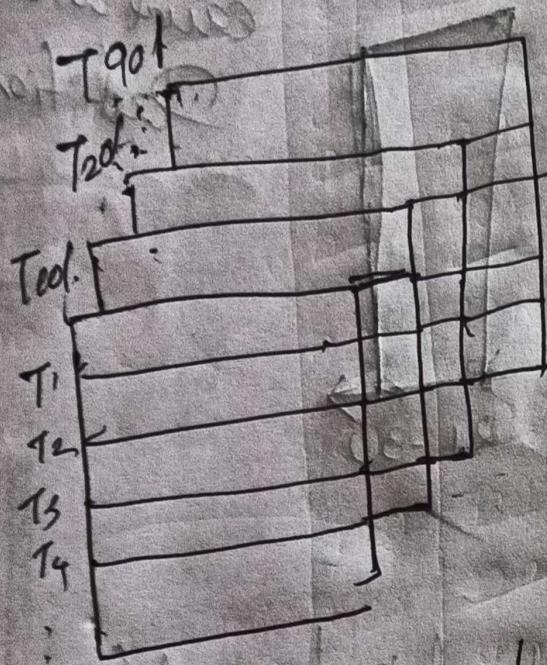


3/23 Sampling

- Monday
- Simple Random Sampling without Replacement (SRSWOR)
 - Simple Random Sampling with Replacing (SRSWR)
 - Cluster Sampling



Clustering sample
($s=2$).



Discretization & concept hierarchy generations

Discretization
converting continuous data to discrete elements.
concept hierarchy.

Clustering - Unsupervised

