Data Mining and Warehousing

Storing Procedure

Store date entreit into useful
information Data mining

↙

Knowledge Data Discovery

⇒ Clean (Redundancy Removed) ~~data~~
↓

Duplication

Same data in Same table or multiple
table

⇒ Data Integration
↓ data
Irrelevant paragraph / Ambiguity

2) Selection (Normalization)
Or Transformation step by step
(i) Normalization     Several Techniques
(ii) selection     ↙ اس سے Data relevant ہے.
chor dain(3 کہ ہے رہی rd

Applications
Data

4 5 6 7 8 9 10
11 12 13 14 15 16 17
18 19 20 21 22 23 24
25 26 27 28 29 30

May **20**
Ramazan 1439 **18**

# Data Pre Processing

Pre - Before        task کو حل کرنے سے پہلے, understand کرنا

Post - After        flawful data اور     . 2 / Perform

Processing          error data

Pre - Before                Post after

fee before                  Salary after month

a) Cleaning

=) Integration         [DataBase]   (File)   meaningful
                                              form

        Different   Resources (2,3,4 سے
                                کوئی آٹا)

2) Transformation ( Repitation )
                    Removal

Normalization

Don't Data   repeat

⇒ Reduction ( volume reduce but meaning Not
                                    changed ).

        → Discritization

        → clarify

        alphabets / Digit / Numeric
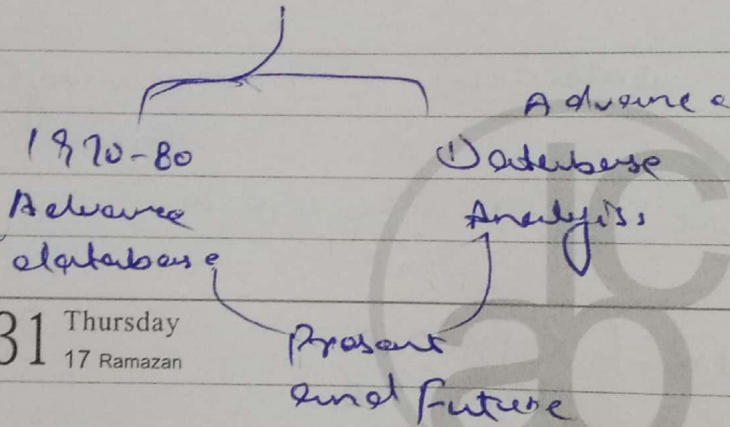
**20 / 18 May**
Ramazan 1439

MAY
M T W T F S S
1 2 3 4 5 6
7 8 9 10 11 12 13
14 15 16 17 18 19 20
21 22 23 24 25 26 27
28 29 30 31

**30 Wednesday**
16 Ramazan

Age of Data
Bulk of data
Large Collection of Knowledge
Future Prediction through datamining

1980 Earlier
↓

1970 → Database

1970-80          Advance
Advance          Database
database          Analysis

**31 Thursday**
17 Ramazan

Present          Customer and
and Future          Consumer

Prevewias Data & Future data in Prediction
Which can be mined?
→ Characterization                    entri
                                → Classes / attri
2) Association and Correlation
→ Clarification  and  Regeression   entity → T
→ Cluster Analysis
→ outlier Analysis                Table entries
                                        ↑
                    → Data Characterization
① Characterization —
                    → Data (Discriminization)
                                        ↓
Data / elaborate اور explain کریں ہم بتائیں    comparison

compare features کی خصوصیت کے ساتھ → features کی خاصیت اس

Classification oh the basis of

Friday
18 Ramazan   1

funct.

process of modeling
function

إلى قيا س في (6

Sala/iep top



Age          Income

s/w ?

Neural Network

① middle layer Network          input

② Deep layer Network          output

Lay e.

 feed forward

G o/

=> Regression .

↓

, Numeric Analysis

① Identify data

Taxonomy
↓

=> cluster -> Same group of data
NO dataset boundaries          Observation

=> outlier -> Noise / fraud
chance of error existence

**20 18** **June**
Ramazan 1439

| | | | | 1 | 2 | 3 |
|4|5|6|7|8|9|10|
|11|12|13|14|15|16|17|
|18|19|20|21|22|23|24|
|25|26|27|28|29|30| |

**4 Monday**
21 Ramazan

1) Data can be mined?

→ Class / Concept Discrimination

→ mining frequent Pattern → Single or multi
   Classification / Regression        dimensional
                    ↓              ↳  Pattern
                                        Frequency
        Supervised learning         Mathematically
              ↕                          ↓
           Label data               Calculation
                                    Presentation is
                                    called regression

→ Cluster                          1) Data B Regression
          Analysis                           so
      ↓
   unlabeled data              } Internet
                                        ↓

**5 Tuesday**
22 Ramazan

→ outlier Analysis              minimize Similar
                                    win

→ outlier Analysis              Intranet → group
      ↓                                ↓
   Fraud Detection              maximize Similar
   / Unusual activity
      → Purchasing
      → Location

① Information retrieval
   Reward or not Complex
   unstructured data

M T W T F S S
30 31         1
2 3 4 5 6 7 8
9 10 11 12 13 14 15
16 17 18 19 20 21 22
23 24 25 26 27 28 29

June **20**
Ramazan 1439 **18**

adhoc

9 Supervised Learning
Unsupervised Learning
Semi Supervised learning
Active learning

→ Business Intelligence → memuel task
                                perform

    → web Search engine
             ↳
        ↳ Directories

## Issues

→ Mining Methodogy →
→ User Interaction

Interdiscipilary

→ Efficieney and Scalability
→ Diversity of Data type type     uncertanity
→ Data mining and Society     patlern evelutvon
    User Interaction           Detection
       ↓

Friendly user interface     heindling of
    or useble            complex dates

efficieney relate to       type
performence
   Do accurate work

Analysis Daily

life

**June**
Ramazan 1439

4 5 6 7 1 2 3
11 12 13 14 8 9 10
18 19 20 21 15 16 17
25 26 27 28 22 23 24
29 30

**8** Friday
25 Ramazan

Getting to know your data    combination of

1. Attribute ⌐→ Characteristics of data features

    ↓

Major Two Types
① Numeric → Digits
② Categorical → Qualitative data

Numeric
Discrete ⌃ continuous → Age, Time
(finite) boundary              continuous
    0 — 9

**9** Saturday    / **10** Sunday
26 Ramazan        27 Ramazan

Feature Vector

First N ⎫ Bivariance          low    High
Last n ⎭                    intensity

                    Binary A (0 — 1)

Categorical                      ↓    ↓
    ⌃                            off  on
Nominal    Ordinal
  ↓          ↓
Names    size of drive ⎫ dependent
colors      drive size ⎭

Center of Tendency
Multiple observation
    Center ⎫ middle
            ⎭ Point

JULY
M T W T F S S
30 31       1
2 3 4 5 6 7 8
9 10 11 12 13 14 15
16 17 18 19 20 21 22
23 24 25 26 27 28 29

June **20**
Ramazan 1439 **18**

Monday **11**
28 Ramazan

Mean

Median

Mode → most frequent occur observation

Range → Difference between max/minimum

Quantile

Bell shape graph      each line $Q_3$    4

100 lines draw    Quantized

Percentile

Standard Deviation      Interquantile

$$\sigma^2 = N^2 \rightarrow variance$$    IQR

$$S'D = \sqrt{N^2} = W$$    $Q_3 - Q_1$ → Maximum − minimum

Variance $\delta^2 = W^2$

$\sigma = S.D = N$

variance

$$\sigma^2 = \sum_{x=1}^{N} \left( \frac{x_i - \bar{x}}{n} \right)^2$$

1, 2, 3, 4, 5, 6, 7, 8,

$$= \left( \frac{1 - 4.5}{8} \right)^2 + \left( \frac{2 - 4.5}{8} \right)^2 \cdots$$

$$= \frac{(9 - 4.5)^2}{8}$$

$$= \left( \frac{(1 - 4.5) + (2 - 4.5)}{8} \right)$$

**13** Wednesday
1 Shawal

Data visualization/ Representation

Know reading
Univariant
Single, more

Q: Plot 2-dimension

Histogram

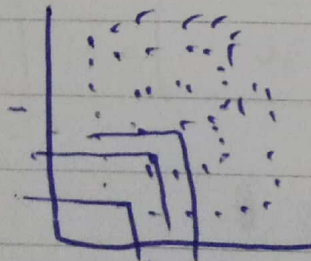**14** Thursday
2 Shawal

Bar graph
(Frequeney graph/ Chart)

Line of code Loc

Standard Plot Scattered graph

524

Pixel Oriented data
↓ dot

ے یہ ورنگ color جہ

0% availability مگر کچھ ورنگ

Pixel are Scattered no data are noise

Geometric Level → ① Different Shapes

awareness  different data representation
③ Color



Scattered plot Matrix $\begin{bmatrix} 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{bmatrix}$

Icon based Visualization

Chernoff

~~Her Hierachie~~

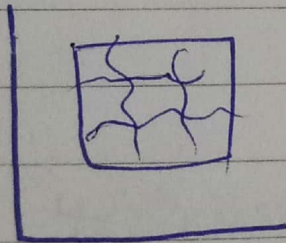Hierarical  → Step-by-study of

data

Measuring  of data Similarity  and
Dissimilarities Matrix

0.0....  Dissimiler  Two attributes  w
0-5...1 Similar ,  $(x_i, y_i)$
one item attributes compare to
another item attribute is called
Similarity and dissimilarities matrix

Name, Color, Age,

Proximity Measures Attribute of Nominal attribute

$$d_{i,j} = \frac{(P - M)}{P}$$

→ Total number of matching attribute

↑ Total number of Attributes

Dil

Ordinal Attribute

Dissimilarities

$$d(i,j) = \lim_{h \to \infty} \left( \sum_{f \to 1}^{P} \left| x_{if} - x_i \right|^h \right)$$

$$L_\infty = P_{max} \left| x_{if} - x_{if} \right|$$

Uniform norm

Cosine Similarity

اپنا پیمانہ دوسرے سے ملا کے دیکھ لے پھر بتا

Total comparison

M T W T F S S
30 31        1
2 3 4 5 6 7 8
9 10 11 12 13 14 15
16 17 18 19 20 21 22
23 24 25 26 27 28 29

June **20**
Shawal 1439
**18**

Wednesday **20**
8 Shawal

# Chapter #3

Data Processing
نتائج کو بنانا accurate اور لیکن
easy to understand    and environment
or User friendly → e.g windows
<u>Time line</u>
Completeness
Consistency → Flow

→ unusual Data

Data <u>Cleaning</u> → Noise (Unwanted effect)
Unwanted ڈیٹا جو کہ ہمارے لیے helpful نہیں ہوتا اسے ہم نکال دیتے ہیں.
Noise کو دور کرنا

<u>Data Integration</u>
Dimension Reduction → Features or Attributes
<u>Reduction</u> ( Volume reduced but meaning
         not changed )

Missing Values

Transformation

<u>Binning</u>
   Bin       Size of bucket in each cycle
   Recyle    Binning
   equally distribution   by mean

$$\frac{8 + 4 + 15}{3} = \frac{27}{3} = 9$$

Equally

Binning by boundary (Median)

Regression

Linear Regression          Two → one to one de
Multiple Regression        one to many

Clustering → unlabel Data

Discripeny <u>Detection</u> → user don't response

Data is more important to

you but you don't response

↔ Resumblance

→ <u>Data alacer</u>

Metadata → Data about data
Relevand Data → Metadata are Sem

Field
<s>fact</s> <u>**Overlooding**</u>                unique
Unique rule → eg id or attributes
Specific person

JULY

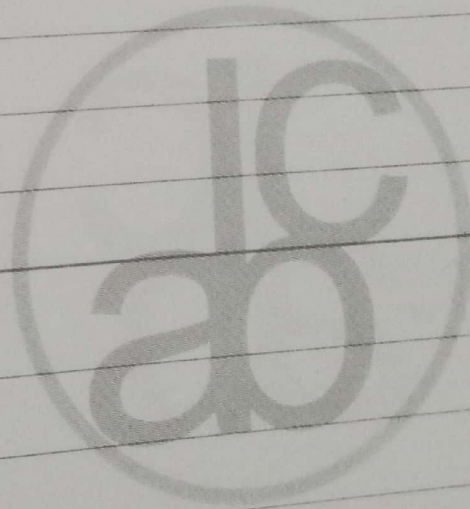| M | T | W | T | F | S | S |
|---|---|---|---|---|---|---|
| 30 | 31 | | | | | 1 |
| 2 | 3 | 4 | 5 | 6 | 7 | 8 |
| 9 | 10 | 11 | 12 | 13 | 14 | 15 |
| 16 | 17 | 18 | 19 | 20 | 21 | 22 |
| 23 | 24 | 25 | 26 | 27 | 28 | 29 |

June **20**
Shawal 1439 **18**

Monday 25
13 Shawal

→ Consective rule

→ Datta discovery Analysis

Datta Scubbing rule and
Address