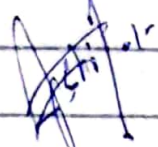


# Lecture 4:- Visualization Model, pipeline & data preparation

• - Visual Mapping must be:-



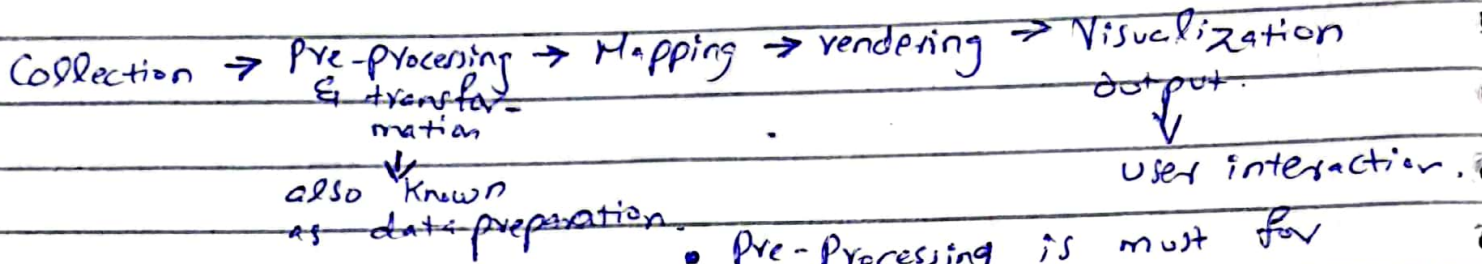
Computable (math)

$$\text{Visual} = f(\text{data})$$

Comprehensible (invertible)  $\text{data} = f^{-1}(\text{Visual})$

Creative

• - Visualization pipeline  $\rightarrow$  Data preparation step within Knowledge discovery process.



- Pre-Processing is must for high data quality which impacts reliability & effectiveness
- Enhance accuracy
- Enhance Completeness (Fill in missing values & ensure all data is present)
- maintains consistency
- Improves Timeliness
- Boosts believability

• - ① Pre-Processing  $\rightarrow$  • Data cleaning • Data reduction  
• Data integration • Data transformation & data discretization

Data Cleaning  $\rightarrow$  maybe • missing data  
• noisy, errors or outliers.  
• inconsistent

• Reason for Missing data  $\rightarrow$  • Equipment malfunction (data collection devices failed, leading to gaps in data)

↓  
If not handled then,  
• incomplete analysis  
• inaccurate models  
• Inconsistent data (data might be deleted)  
• misunderstanding

• Perceived unimportance.

How to handle?

- Ignore tuple
- Fill missing values (use mean (symmetric) median (skewed))

• Most likely value:-  
use bayesian or decision trees  
to predict value

• - Inference-Based → • Bayesian  
methods • decision tree

• - Identify relationship among  
Variables → • Linear regression  
• Multiple regression  
• Non-Linear regression

• - Nearest-Neighbours Estimator  
• Fill missing values with  
most frequent value (categorical  
data)  
or average value (numerical  
data)

• - Noisy data → • Random error or variance in  
measured variable.

• Handle noisy data by Binning, Regression, clustering

• Binning → • Sort the data  
• partition into bins  
• Smooth bins by → mean, median,  
boundary.

Example:-

data = (8, 4, 9, 15, 21, 21, 24, 25, 26, 28, 29, 34)

① Sort → 4, 8, 9, 15, 21, 21, 24, 25, 26, 28, 29, 34

② Bin 1 → [4, 8, 9, 15]  $\frac{8+9}{2} = 8.5$

Bin 2 → [21, 21, 24, 25]

Bin 3 → [26, 28, 29, 34]

Bin Means →

[9, 9, 9, 9], [23, 23, 23, 23],  
[29, 29, 29, 29]

Bin Medians → [8.5, 8.5, 8.5, 8.5]  
(Adding middle  
number and  
divided by 2) [23, 23, 23, 23],  
[28, 28, 28, 28]

↑ identify lowest and  
highest value

Bin Boundaries →

[4, 4, 4, 15]



metadata  $\rightarrow$  name, meaning, datatype, range, null values.

Predicts one attribute based on other.

- Regression  $\rightarrow$  used to smooth data by fitting into regression function

Linear regression

- Finding the best line that fits two attributes. goal is to minimize distance b/w data points and the line.

- Clustering  $\rightarrow$  detect & remove outliers. values outside set of clusters.

- ② Data integration (Combining data from multiple sources into one)  $\rightarrow$  EIP (Entity Identification Problem)

- While integrating, EIP is a great Challenge.

- Schema integration & object matching combining data from different sources which use different formats.

- Metadata can help avoid errors in schema integration and transform data.

- When matching attributes from two databases, structure of data should be checked.

- Handle redundant data  $\rightarrow$  by detecting through correlation analysis and covariance analysis.

- ③ Data reduction  $\rightarrow$  reduced representation of data set (smaller in volume but produces same analytical results)

$\downarrow$  use

b/c more data & complex data analysis take very long time to run.

Strategies

- Dimensionality reduction
- Wavelet transform
- PCA
- Feature selection, creation.

• - Numerosity reduction → • Parametric (Regression, log linear model)

• Non-Parametric (Histogram, clustering sampling)

• data cube aggregation.

• - Data compression

→ lossless (reconstruction, no loss of information)

lossy (reconstruct only approx of original data)

i) Dimensionality reduction (① Wavelet transforms ② PCA ③ supervised & non-linear)

• - curse of dimensionality → • no. of features in data ↑  
data becomes harder to analyze

• difficult to measure distances between points, imp for clustering or outliers.

ii) • - Dimensionality reduction →

• removes noise in data and unnecessary features.

• less memory used & data processing faster.

• easier to visualize data by showing in fewer dimension (3D to 2D)

iii) • - Principle Component analysis (PCA) →

(Reduce no. of features in data while keeping most important information)

• Works for numerical data.

• PCA handles sparse data better than wavelet transforms.

Steps: • Normalize input data.  
$$X' = \frac{X - X_{\min}}{X_{\max} - X_{\min}}$$

(Find for every entry in the table)

• Covariance matrix / Z-score standardization

• Eigenvalue, Eigen vector

• Select principle components

• Transform data.

• Sort components by significance

• reduce data dimensionality.



c) - Numerosity reduction →

• Parametric methods

(e.g. Regression, non-linear model)

→ data follows a specific model.  
• estimate and store parameters and discard data (except possible outliers)

Consider each tuple as a point in an  $n$ -dimensional space.

• Non-parametric

(e.g. histogram, clustering, sampling)

→ do not assume models.

• Regression Analysis (way to understand relationship b/w dependent & independent variables)

• Parameters are estimated to give best fit.

• Best fit is evaluated by using least squares method.

• used for prediction, hypothesis testing & modeling for causal relationships.

• Histogram Analysis (Divide data into buckets and store avg sum of each bucket)

rules → equal width & frequency.

• Clustering (store cluster representation)

• Sampling (obtain small data to represent whole dataset).

may not reduce database I/O

d) - Data Compression →

String Compression (lossless)

Audio/Video Compression (lossy)

Dimensionality & Numerosity may be forms of data compression.

## ④ Data Transformation →

### • Descriptive Statistics

Range, Min/Max

Average (one measure of central location in data set)

Median

Mode (Another measure of central location in data set)

### • Distribution Statistics

How meaningful is chosen average?

Varience

Std (~~Varience~~)

Histogram / Normal distribution

• When data skewed, Mean & SD can be misleading.

• Skewness

( $SK > |1|$ )  
↓  
(non-symmetrical)

• Negatively skewed

(Mean < Median)

• Positively skewed

Mean > Median

• Histogram shows (Center, spread, skewness, outliers, multiple modes in data)

↓  
for small dataset  
Histogram can be misleading.

• Histogram tells nothing about relation among variables

• Normalization → • normalize data b/w  $[-1, 1]$  or  $[0, 1]$

$$Z = \frac{V - \mu}{\sigma}$$

### Discretization

(Divide range of continuous attribute into intervals)

→ • Binning (unsupervised)

• Histogram (un)

• Clustering (un)

• Decision trees (supervised) (use entropy)

• Correlation (un)