

## Chapter - 2

### Data - Preprocessing

→ Raw data contained in databases is unprocessed, incomplete and noisy.

→ ~~Reasons:~~

Data preprocessing is the process of preparing raw data for analysis by cleaning and transforming it into a useful format.

- Real datasets often having missing or incomplete entries.
- Data may contain typos, duplicates or outliers.
- Different features may have different scales not suitable for data mining models.
- Some features are irrelevant, redundant or highly correlated.
- values not consistency with policy or common sense.

→ Steps in Data Preprocessing :-

- Data cleaning
- Data Integration
- Data Transformation
- Data Reduction.

## Data cleaning :-

→ It is the process of identifying and correcting errors, or inconsistencies in the dataset.

→ It involves handling missing values, removing duplicates, and correcting incorrect or outlier data to ensure the dataset is accurate & reliable.

→ Clean data is essential for effective analysis, as it improves the quality of results & enhances the performance of data models.

A. Missing values : when data is absent from a dataset.

→ Remove records with missing values.

→ Fill with mean / median / mode (imputation)

Noisy Data : It refers to irrelevant or incorrect data.

→ use smoothing techniques (Binning, moving avg.)

[Binning] : Data is sorted into equal segments and each segment is smoothed by replacing values with mean or boundary values]

→

1. Ignores the tuple :-

This is usually done when the class label is missing. This method is not very effective, unless the tuple contains several attributes with missing values.

2. Fill in the missing value manually :-

In general, this approach is time-consuming may not be feasible given a large dataset with many missing values.

now  
replace all missing attributes value by the same constant  
such as label like "unknown" or "-∞".

4. use the attribute mean to fill in the missing value.
5. use the " " " for all samples belonging to the same class as the given tuple.  
Use the most probable value to fill in the missing value.

In DT, regression & Bayesian formalism.

### B. Noisy Data

→ Noise is a random error or variance in a measured variable.

- incorrect attribute values may due to
  - faulty data collection instrument.
  - data entry problems
  - data transmission problems
  - technology limitations.
- inconsistency in naming convention.

→ other data problems.

- Duplicate records
- incomplete data
- inconsistent "

## → (i) Binning method

- first sort data and the sorted values are distributed into a number of "buckets", or "bins".

→ When one can smooth by bin means, smooth by bin median, smooth by bin boundaries etc.

### (i) By Bin means

each value in a bin is replaced by the mean value of the bin

### (ii) By Bin median:

→ each bin value is replaced by the bin median.

### (iii) By bin boundaries:

→ minimum & maximum values in a given bin are identified as the bin boundaries.

→ Each bin value is then replaced by the closest boundary value.

→ In general, the larger the width, the greater effect of smoothing. So ~~because~~ bins may be equal-width,

or equal depth (frequency) partitioning: (distance)

### equal width (distance) partitioning:

→ It divides the range into N intervals of equal size

→ If A & B are lowest & highest values of the attribute, the width of intervals will be  $w = (B - A)/N$ .

→ The most straightforward.

→ But outliers may dominate presentation.

## Equal-depth (frequency) Partitioning:

It divides the range into N intervals, each containing approximately same number of samples

- Good data scaling.
- Managing categorical attributes can be tricky.

Q Sorted data for price (in dollars):

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

→ Partition into (equal-depth/equal-frequency) bins:

Bin 1: 4, 8, 9, 15

→ Smoothing by bin median.

Bin 2: 21, 21, 24, 25

Bin 1: 8.5, 8.5, 8.5

Bin 3: 26, 28, 29, 34

Bin 2: 22.5, 22.5, 22.5

Bin 3: 28.5, 28.5, 28.5, 28.5

→ Smoothing by bin mean:

Bin 1: 9, 9, 9, 9

→ Smoothing by bin boundaries:

Bin 2: 23, 23, 23, 23

Bin 1: 4, 4, 4, 15

Bin 3: 29, 29, 29, 29

Bin 2: 21, 21, 25, 25

Bin 3: 26, 26, 26, 34

Q Partition into (equal-width) bins:

lowest value A = 4, highest value B = 34

$$\text{Range} = 34 - 4 = 30$$

$$\text{Width of interval} = \frac{30}{3} = 10$$

→ Hence choosing 3 bins.

S-3 Bin intervals      S-4 Assign Data to Bin

Starting at min = 4

Bin 1: 4, 14)

Bin 2: [14, 24)

Bin 3: [24, 34)

Q want equi-frequency

S-1 Data must be sorted.

S-2 Take 3-bins.

Total value = 12

Each bin should have  $12 \div 3 = 4$  values.

$$\underline{\underline{S-3}}$$
$$\text{Bin 1: } \{4, 8, 9, 15\}$$

$$\text{Bin 2: } \{21, 21, 24, 25\}$$

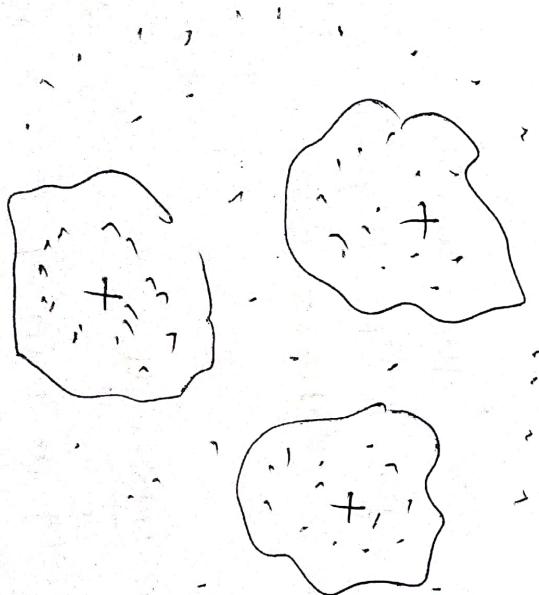
$$\text{Bin 3: } \{26, 28, 29, 34\}$$

2. Regression: Data can be smoothed by fitting the data to a function such as regression.

(1) Linear regression: best line to fit two attributes (Dependent & independent variables)

(2) Multiple linear regression: It is an extension of linear regression where more than two attributes are involved.

3. Clustering:



## Identifying misclassification

Misclassification means that a model predicts the wrong class label compared to the actual (true) class label in your dataset.

<u>Brand</u>	<u>Frequency</u>
USA	1
France	1
US	156
europe	46
Japan	51

(Notice strange about this frequency distribution)

- The frequency distribution shows 5 classes. However, two of the classes USA & France have count of only one each. What is clearly happening here is that two of the records have been inconsistently classified w.r.t. to the origin of manufacture.
- To maintain consistency, the record have been labeled as europe instead of USA & France -

## Graphical methods for identifying outliers

- Outliers are extreme values that go against the trend or remaining data.
- Graphical methods for identifying outliers for numerical variables is to use "histogram" or the variable.

Q 45, 50, 52, 48, 47, 49, 51, 46, 50, 50, 48, 49, 5

S-1 45, 46, 47, 47, 48, 48, 49, 49, 50, 50, 51, 51, 52, 52,

$$\bar{x} = 49, Q_1 = 47, Q_3 = 51$$

$$IQR = Q_3 - Q_1 = 51 - 47 = 4$$

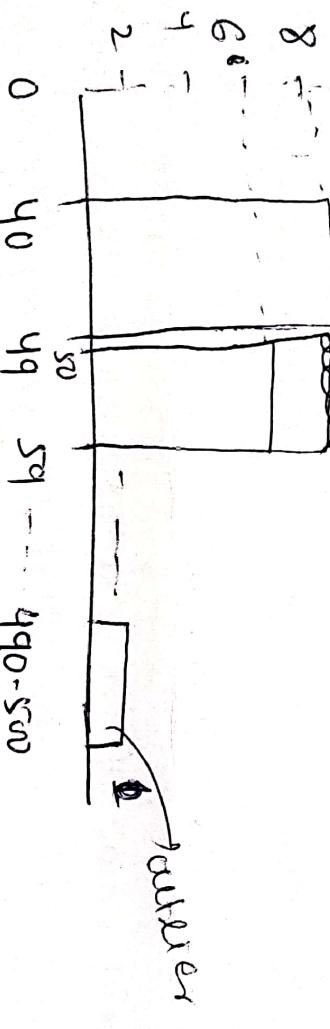
outlier

$$\text{lower Bound} = Q_1 - 1.5 \times IQR = 41$$

$$\text{upper } " = Q_3 + 1.5 \times IQR = 57$$

S-5 anything > 57 → outlier

$$\Rightarrow 50 > 57 \rightarrow \text{no}$$



Outlier

### Scatter plot

Student score

1

45

50

2

52

50

3

48

48

4

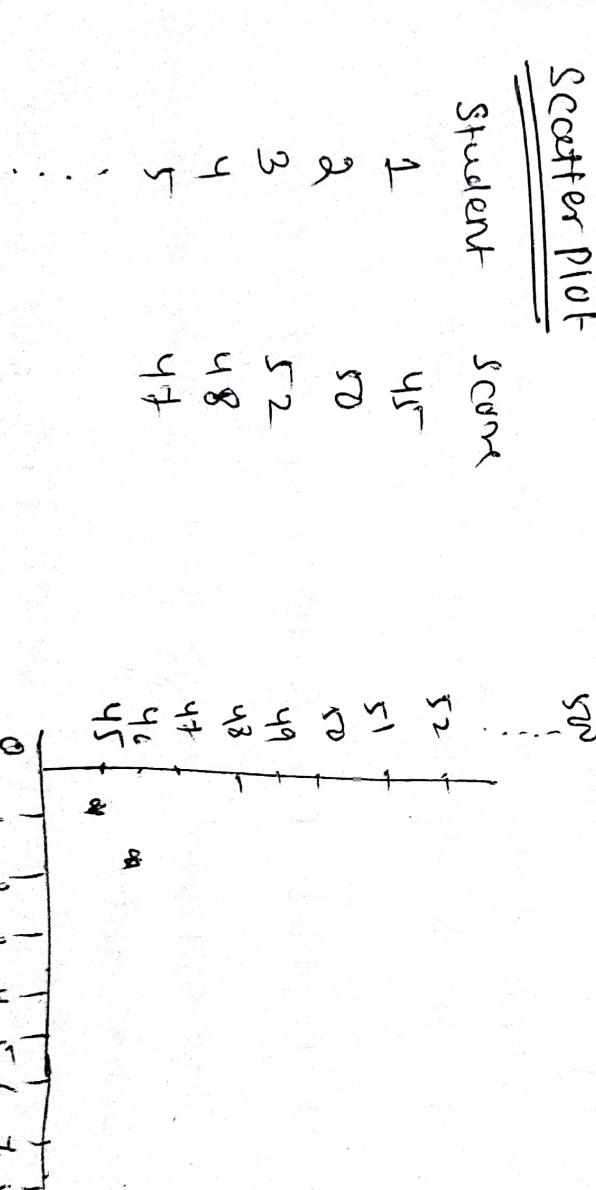
47

47

5

45

45



Measures of center & spread

$$\text{center (mean)} = \bar{x} = \frac{\sum x}{n} \text{ & spread (SD)} = \sqrt{\frac{\sum (x - \bar{x})^2}{n-1}}$$

## 1a Integration

Data Integration - the merging of data from multiple data sources. It helps to reduce and avoid redundancies and inconsistencies in the resulting dataset which improve the accuracy and speed of the subsequent data mining process.

- entity resolution problem →
- It occurs when we try to determine whether two or more records from different data sources refer to the same real-world entity.

Possible → Rule based matching

- probabilistic matching (similarity score)
- Machine learning model.

## Data Transformation

→ It is the process of converting data from one format, structure or value representation to another to ensure compatibility, consistency and better quality for analysis.

### Types

1. Smoothing :-

- Removes noise from data.
- Ex: Binning, regrouping & clustering.

2. Aggregation

Summarizing data.

Ex: Daily sales data → Monthly sales totals.



### 3. Generalization:

Replacing detailed data with higher-level concepts.

Ex 23 years → "Young Adult" or "Young"

middle-aged & senior

### 4. Normalization (scaling)

Rescaling values to a standard range (Ex [0, 1] or

$x_{\text{new}} = \frac{x - \text{min}(x)}{\text{range}(x)}$   $\rightarrow$  Z-score)

### 5. Attribute construction (feature engineering)

Creating new attributes from existing ones.

Ex from DOB create "Age"

### MIN-MAX NORMALIZATION

Min-max normalization works by scaling a linear transformation on the original data.

$$x_{\text{new}} = \frac{x - \text{min}(x)}{\text{range}(x)} = \frac{x - \text{min}(x)}{\text{max}(x) - \text{min}(x)}$$

where  $x = \text{original value}$

$\text{min}(x) = \text{minimum value of}$

the attribute

$\text{max}(x) = \text{maximum value of}$

the attribute

$x' = \text{normalized value (both only)}$



50, 60, 70, 80, 90, 100

$$\min = 50, \max = 100$$

For normalized 80:  $x' = \frac{80 - 50}{100 - 50} = \frac{30}{50} = 0.6$

so the normalized value of 80 is 0.6.

Q Suppose min & max values for the continuous 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 —

$$\frac{\underline{x}}{\text{val}} = \frac{73,600 - 12,000}{98,000 - 12,000} = 0.716$$

Q minimum weight of 1613 pounds & max = 4997 pounds

$$\text{range} = 4997 - 1613 = 3384 \text{ pounds}$$

Q) Find min-max normalization weight for a extra weight vehicle of 1613 pounds.

$$x_{\text{mm}} = \frac{x - \min(x)}{\text{range}(x)} = \frac{1613 - 1613}{3384} = 0$$

(iii) For midrange =  $\frac{\max(x) + \min(x)}{2} = \frac{4997 + 1613}{2} = 3305 \text{ pounds}$

so min-max normalization for midrange is

$$x_{\text{mm}} = \frac{3305 - 1613}{3384} = 0.5$$



Scanned with OKEN Scanner

For heaviest vehicle.

$$X_{\text{mm}} = \frac{4499 - 1613}{3384} = 1$$

Z-score standardization

Already done in appendix

Decimal Scaling :

Decimal Scaling is a normalization technique where the normalized value lies between  $-1 \leq 1$ .

$$x^*_{\text{decimal}} = \frac{x}{10^d}$$

where  $x$  = original value

$x^*$  = Normalized value or no. of digits in  
 $d$  = smallest integer such that  $|x^*| < 1$   
the data value such that  $|x^*| < 1$   
for all values.

$$-987, -120, 56, 300, 850$$

maximum absolute value is 987  
to make it less than 1, divide by  $10^3$  ( $987 < 1000$ )

$$\therefore d = 3$$

$$x^* = \frac{x}{1000}$$

Original      Normalized

$$-987$$

$$-120$$

$$56$$

$$300$$

$$850$$

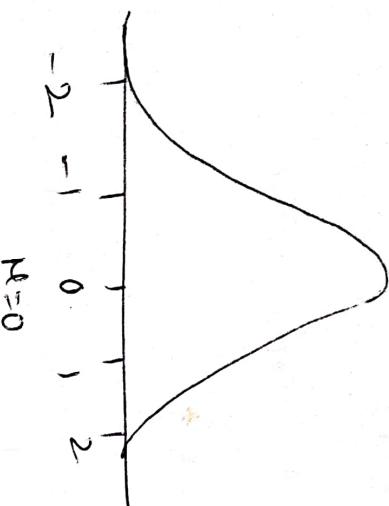
i.e. normalized value lies between  $-1 \leq 1$



Scanned with OKEN Scanner

## Transformation to Achieve Normality:

The normal distribution is a continuous probability distribution commonly known as bell curve, which is symmetric.



→ For standard normal distribution  $z$  has mean  $\mu = 0$  &  $SD(\sigma) = 1$  but distributions may still be skewed.

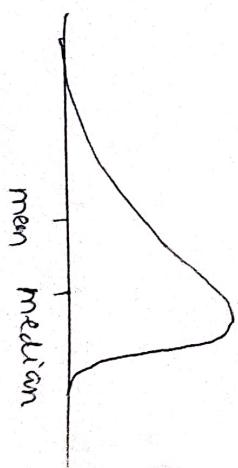
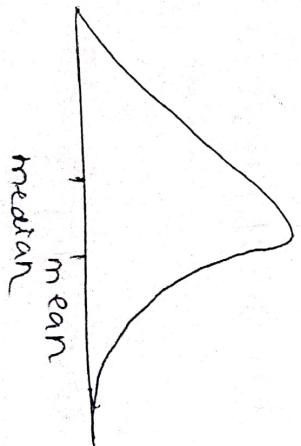
$$\text{Skewness} = \frac{3(\text{mean} - \text{median})}{\text{Standard deviation}}$$

for weight

$$\text{Skewness} = \frac{3(300.490 - 283.5)}{852.646} = 0.6$$

for weight - z

$$\text{Skewness} = \frac{3(0 - (-0.2))}{1} = 0.6$$



mean > median → positive skewness

mean < median → negative skewness

→ To eliminate skewness apply a transformation to the data.

→ common transformations are logarithmic, log transformation.

$$\bar{x} = \log(x+1)$$

(ii) Square root transformation.

$$\bar{x} = \sqrt{x}$$

(iii) inverse square root transformation.

$$\bar{x} = \frac{1}{\sqrt{x}}$$

Flag variables

→ Flag variables (dummy or indicator variable) is a categorical variable taking only two values, 0 & 1.

Ex categorical predict sex. (female & male)

By flag variable as sex-flag.

If sex = female then sex-flag = 0  
If sex = male then sex-flag = 1.

Ex categorical predictor region has K=4 (north, east, south, west)

then

north-flag : If region = north then north-flag = 1, otherwise

north-flag = 0

east-flag : If region = east then east-flag = 1,  
east-flag = 0

south-flag : If region = south " south-flag = 1  
south-flag = 0.

## Classifying categorical variables:

- Reclassifying categorical variables mean combining, merging or transforming categories or variables into fewer, more meaningful groups to improve data analysis or model performance.
- when a categorical variables has too many categories or irrelevant distinctions, we merge or rename them to make analysis easier and more accurate.

### why

- Too many categories exist.
- some categories have very few observations.
- several " " have similar in meaning or behaviour.
- To create more useful or predictive grouping for DM Algorithms.

e.g. 50 states could be classified as the variable region, like Northeast, southeast, North central, southwest & west.

OR can be reclassified as economic - level like richer state, midrange state and poorer state.

## Adding an index field

- An index field is a new variable added to the dataset to uniquely identify each record, it is especially useful when the dataset does not contain a unique ID or key.

### Reason

1. uniquely identifying each record.
2. Data tracking
3. Debugging
4. Joining datasets.

Vedavita Non-discriminatory  
Advoctate "andrea" 14 year old  
S. a 5 - 3 0 - 2 "tinyt" analysis

## DRPA

Removing variables that are not useful

- Some variables in the dataset may not contain useful information for DM task. These variables should be identified and removed to simplify the dataset and improve the performance of the model.

### Reason

- Reduce noise (Relevant variables)
- Improve computation
- Simplify interpretation.
- Avoid redundancy.

#### Ex ① unary variables

- It is a variable that has only one value for all records in a dataset.

an	ID	member	country	salary
1	Male	India	40k	
2	Male	India	30k	
3	Female	India	20k	

#### ② variables that are very nearly unary

- A nearly unary variables that has almost the same value for all records with only a few records differing.

#### Ex

ID	Membership	Age
1	Gold	38
2	Gold	28
3	Gold	45

(4) River 30 → only differ



Karen  
angono

Variables that should probably not be removed  
some variables that appear unimportant at first may  
later turn out to be quite useful.

### Types of variables

- categorical identification that group records meaningfully.
- variables with potential interaction effects.
- variables with missing or rare values.
- 1) that are meaningful to interpretation.

Ex

ID	Name	Department	Salary	Age
1	John	HR	40k	25
2	Ivan	Sales	50k	28
3	Ray	HR	30k	40

