

Ch-2 Data Preprocessing

The data are incomplete \rightarrow lacking attribute values
Ex: salary = ??
occupation = ??
 \rightarrow noisy \vdash containing errors \rightarrow data
or outliers values Ex: salary = -10

\rightarrow inconsistent - containing discrepancies
Age = 29, Birthday = 5/04/1985

Data cleaning \rightarrow clean the data by filling
missing values,

smoothing noisy data

identifying or removing outliers
and resolving inconsistencies

Data integration \rightarrow integrating multiple databases,
data cubes or files

Data transformation \rightarrow normalization and aggregation

Data reduction \rightarrow a reduced representation
of the data set that is much
smaller in volume, yet produces
the same analytical results.

(1) Data aggregation \rightarrow building a
data cube.

(2) Attribute subset selection \rightarrow
removing irrelevant attributes
through correlation analysis.

(3) Dimensionality reduction.
Ex - DFT, PCA, DCT, DWT

(4) Numerosity reduction \rightarrow
replacing the data by
smaller alternatives
such as clustering or
parametric models.

Data characteristics can be obtained from

- (1) central tendency \rightarrow mean, median, mode and midrange.
- (2) Dispersion of data \rightarrow quartiles, interquartile range (IQR) and variance.

There are 3 measures for central tendency.

\rightarrow Distributive measure \rightarrow sum(), count()

Algebraic measure \rightarrow mean() = $\frac{\text{sum}()}{\text{count}()}$
weighted mean.

Holistic measure.

Distributive measure \rightarrow can be computed for a given data set by dividing the data into smaller subsets, then computing the measure and then merging the results.

Algebraic measure \rightarrow computed by applying an algebraic function to one or more distributive measures.

Ex \rightarrow mean(),

weighted mean = $x = \frac{w_1x_1 + w_2x_2 + \dots + w_Nx_N}{w_1 + w_2 + \dots + w_N}$.

Trimmed mean \rightarrow mean obtained after chopping off values at the high and low extremes.

median \rightarrow mid value if N is odd.
average of two middle values, if N is even.

Holistic measure \rightarrow must be computed on the

entire data set as a whole.

Ex: approximate median = $L_1 + \left(\frac{N/2 - (\text{freq})_L}{\text{freq median}} \right) \text{width}$

L_1 = lower boundary of the median interval.

N = total no. of values in the entire dataset.

f_{median} = frequency of the median interval

width = width of the median interval.



Mode \rightarrow The value that occurs most frequently in the set.

Data sets with one, two or 3 modes are called unimodal, bimodal and trimodal.

A data set with 2 or more modes is called multimodal.

Midrange \rightarrow Average of the largest and smallest values in the set.

Range, Quartiles, Outliers and Boxplots.

Let x_1, x_2, \dots, x_N be a set of observations.

Kth percentile of a set of data is the value x_i having the property that K percent of the data entries lie at or below x_i .

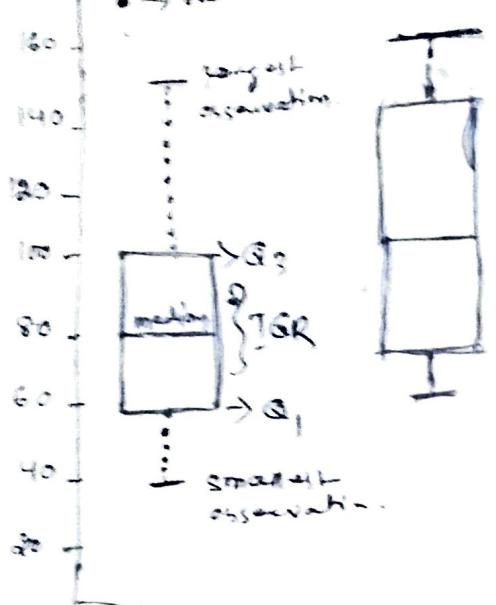
1st quartile (Q_1) \rightarrow 25% IQR \rightarrow interquartile range
 3rd " (Q_3) \rightarrow 75%

$$\text{range} = Q_3 - Q_1$$

(3)

Outliers → are values falling at least 1.5 IQR times above the 3rd quartile or below the 1st quartile.

Boxplots → is a popular way of visualizing a distribution.



For branch 1,
median price of items sold is \$80
Q₁ is \$60
Q₃ is \$100

minimum → 40
max → 150

outliers → 175 and 202
marked as
as the value dot
greater than ($1.5 \times \text{IQR}$)

Branch 1 Branch 2 Branch 3 Branch 4

(Boxplot for unit price data for items sold at four branches)

Variance and Standard Deviation

$$\sigma^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2 \rightarrow \text{Variance of } N \text{ observations.}$$

x_1, x_2, \dots, x_N

Standard deviation, $\sigma = \sqrt{\sigma^2}$

Graphical representation of Data

(3)

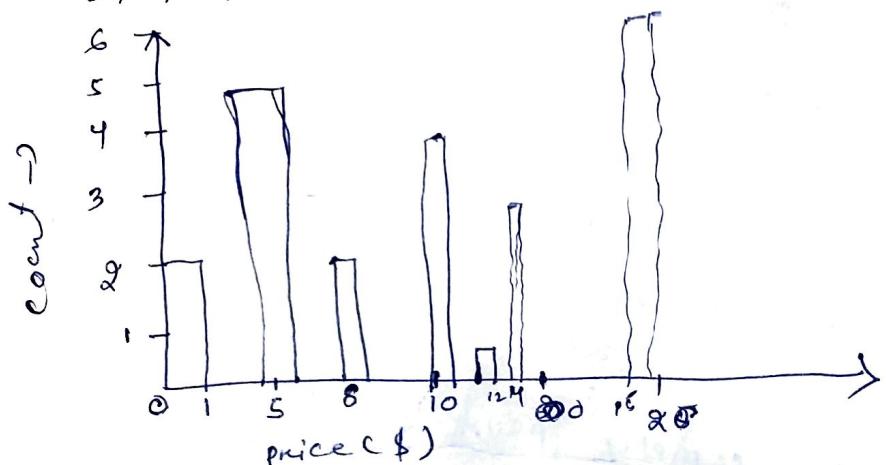
To display data distribution, we can use

(1) Histogram on Box chart.

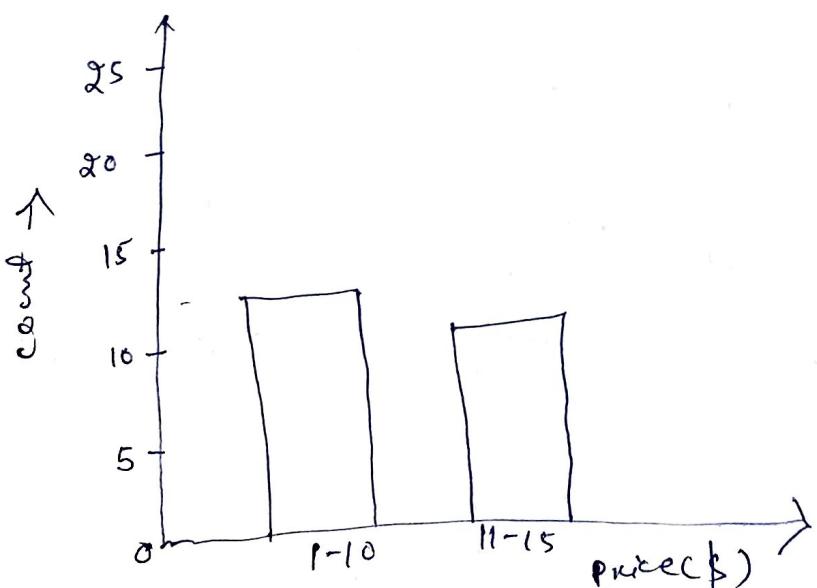
(2) Quantile plot.

Ex The following data are a list of prices commonly sold items. The numbers have been sorted.

1, 1, 5, 5, 5, 5, 5, 8, 8, 10, 10, 10, 10, 12, 14, 14, 14,
15, 15, 15, 15, 15, 15, 18, 18, 18, 18, 18, 18, 18, 18, 20,
20, 20, 20, 20, 20, 21, 21, 21, 21, 25, 25, 25, 25, 25, 25, 28, 28, 30, 30, 30.



(Each bucket represents one price-value)

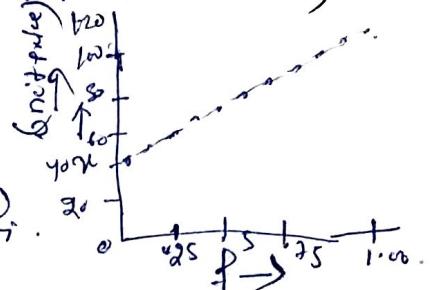


Equal-width histogram for price.

Quantile plot

(4)

- Let x_i , for $i = 1 \dots N$, be the data sorted in increasing order so that x_1 is the smallest observation and x_N is the largest.
- Each observation, x_i is paired with a percentage, f_i which indicates that approximately $100f_i\%$ of the data are below or equal to the value x_i .
- $$f_i = \frac{i - 0.5}{N}$$
- x_i is plotted against f_i .



Quantile-quantile plot (Q-Q plot)

Let $x_1 \dots x_N$ be the data from the 1st branch and $y_1 \dots y_M$ " " " " and branch where each data set is sorted in increasing order.

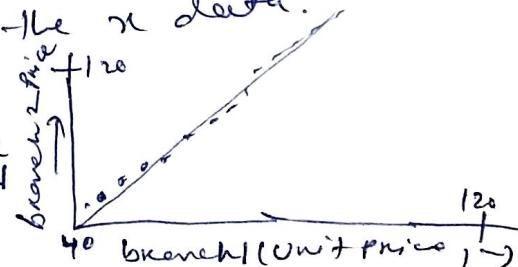
If $M = N$, simply plot y_i vs x_i , where y_i are the $(i - 0.5)/N$ quantiles of their respective data sets.

If $M < N$, here y_i is the $(i - 0.5)/M$ quantile of the y data, which is plotted against the x data.

The $(i - 0.5)/M$ quantile of the y data.

Each point corresponds to the same quantile for each dataset, shows the price of items sold at Branches, branched at the quantiles.

Scatter plot



Hence, each pair of values is treated as a pair of coordinates in an algebraic sense and plotted as points on the plane.

Data cleaning

How to handle missing values?

(1) ignore the tuple where there is a missing value
But not an efficient way.

(2) Fill it manually. (tedious, infeasible)

(3) Fill it with global constant (e.g. unknown or 0)

Not an efficient method.

(4) By the attribute mean.

(5) By most probable value (this is)
determined by regression or Bayesian formula on decision tree)

How to handle noisy data?

Noise is a random error on variance of a measured variable.

Happened due to → fault data collection
data entry problems.
technology limitation

It can be handled by the following techniques:-

(1) Binning → 1st sort the data, then partition into bins or buckets of equal length and then one can smooth by 3 ways.

(2) Regression → Data can be smoothed by fitting the data to a function, such as regression.

(3) clustering → one can detect the outliers by clustering and remove it.

Binning →

(1) Smoothing by bin means \rightarrow each value \leftrightarrow
-the bin is replaced
by -the mean value
of the bin.

(2) Smoothing by bin medians \rightarrow replaced by bin median.

(3) Smoothing by bin boundaries \rightarrow the minimum and max. values are identified as the bin boundaries. Each bin value is then replaced by the closest boundary value.

replaced by
ex 4, 8, 15, 21, 21, 24, 25, 28, 34 (price data)
on \$

partition into (equal-frequency) bins:

Born 1: 4, 8, 15

Bin 2: 21, 21, 24

Bin 3: 25, 28, 34

smoothing by bin means

seen 1 ! 9, 9, 9

Bin 2: 22, 22, 22.

Bin 3: 29, 29, 29 ..

Smoothing by bin boundary

Ben 1: 4, 4, 15

Bin 2: 21, 21, 24

Bin 3 : 25, 25, 34

Ex
 13, 15, 16, 16, 19, 20, 20, 21, 22, 22, 25, 25, 25, 25,
 30, 33, 33, 35, 35, 35, 36, 40, 45, 46, 52, 70.
 we smoothing by means , using a bin depth of 3 .

Data integration -

Various Issues in Data Integration -

Schema
integration

Object
matching

Redundancy

curl_ID on one
database >
curl_NO on another
database.

(When integrating
problem arises)

Brian Lara
= Lara Dutta
or

Bill Clinton =
William Clinton

An attribute is
redundant if
it can be derived
from another
attribute or a
set of attributes.
Ex: Annual revenue.

✓ Redundancies can be detected by
correlation analysis

i.e. by analyzing how strongly one
attribute influences the other.
(i.e. causality effect)

For Numeric attributes, the correlation between
two attributes can be evaluated by computing
the correlation coefficient

↓
also called Pearson's product-
moment coefficient

⇒ The correlation doesn't imply causality.
i.e. if A, B are correlated, this doesn't
necessarily imply that A causes B or
B causes A

Ex:- we may find that attributes
representing the no. of hospitals
and the no. of car thefts in a city are
correlated. This does not mean that one
causes the other.

Chi-Square test (χ^2) -

→ The correlation relationship between two attributes^(A & B), can be evaluated by χ^2 test.
 (of categorical data) Ex: Discrete attribute: ZIPcode, profession
 i.e discrete continuous " " " form, height, weight, weight

1: Assume a null hypothesis.

2: Compute value of $\chi^2(P)$ → Degree of freedom

3: calculate value of χ^2 & DOF

+: write the rule to prove the hypothesis.

5: else write the summary based on decision.

$$\chi^2 = \sum_{i=1}^c \sum_{j=1}^r \frac{(O_{ij} - E_{ij})^2}{E_{ij}} \quad \begin{array}{l} (A \text{ has } c \text{ distinct values}) \\ (B \text{ has } r \text{ " }) \end{array}$$

or

$$\chi^2 = \sum_{i=1}^P \sum_{j=1}^Q \frac{(O_{ij} - E_{ij})^2}{E_{ij}} \quad \begin{array}{l} O \rightarrow \text{observed value (frequency)} \\ E \rightarrow \text{expected value (frequency)} \end{array}$$

$$E = \frac{\text{Col total} \times \text{Row total}}{\text{Grand total}}$$

$$\text{DOF} = (\text{No. of col} - 1) \times (\text{No. of rows} - 1)$$

Eg Suppose that a group of 1500 people was surveyed. Each person was polled as to whether their preferred type of reading was fiction or nonfiction. Some have two attributes, gender, preferred reading

	Male	Female	Total	Preferred reading
fiction	250	200	450	
non-fiction	50	1000	1050	

Are gender & preferred reading are correlated?

$$\chi^2(P) = 0.071$$

Fiction → real face
 Non-fiction → imagined story
 Based story

Ans

	male	female	Total
-fection	250	200	= 450
non-fection	50	1000	= 1050
Total =	300	1200	= 1500
	1500		

(2)

for fication

	<u>observed value</u>	<u>Expected value</u>	χ^2
Male	→ 250 →	$\frac{300 \times 450}{1500}$	$\frac{(250 - 90)^2}{90}$
		= 90	

Female	→ 200 →	$\frac{1200 \times 450}{1500}$	$\frac{(200 - 360)^2}{360}$
		= 360	

for non-fication

Male	→ 50 →	$\frac{300 \times 1050}{1500}$	$\frac{(50 - 210)^2}{210}$
		= 210	

Female	→ 1000 →	$\frac{1200 \times 1050}{1500}$	$\frac{(1000 - 840)^2}{840}$
		= 840	

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

$$= [507.93] = 508$$

$$D.O.F = (2-1) \times (2-1) = 1 \quad (\because 2 \times 2 \text{ matrix})$$

$$\chi(1) = 0.021$$

$r_{A,B} < 0$, A & B are negatively correlated
 $r_{A,B} = 0$, A & B are independent
 $r_{A,B} > 0$, A & B are positively correlated

C No correlation between them

$r_{A,B} = 0$, A & B are independent

(i.e. if A \downarrow , B also increases)

$r_{A,B} > 0$, A & B are positively correlated

$$-1 \leq r_{A,B} \leq 1$$

Beta, β
i.e. values of A

(Number)
A, B → *Alpha* values

N Δ G_B
N → no. of subjects

For Number, it is correlation coefficient

i.e. $r_{A,B} = \frac{\sum_{i=1}^N (A_{i,B_i}) - N \bar{A}\bar{B}}{\sqrt{N}}$

Statistically correlated

$\chi^2 >$ value of χ^2 (df, α)

Statistically co-related

\leftarrow There are two attributes are

Q → *Gender* & *Age* not significant
Q → *Prefers reading* & *Gender* not significant

→ The hypotheses can be rejected

As 508 is above this value,

Since χ^2 value which is computed is 10.828

1	2	3
1	2	3
1	2	3

At $\alpha = 0.01$ from χ^2 distribution table, χ^2 value is 10.828

3