

- 2.2 Suppose that the data for analysis includes the attribute *age*. The *age* values for the data tuples are (in increasing order) 13, 15, 16, 16, 19, 20, 20, 21, 22, 22, 25, 25, 25, 25, 30, 33, 33, 35, 35, 35, 35, 36, 40, 45, 46, 52, 70.
- What is the *mean* of the data? What is the *median*?
 - What is the *mode* of the data? Comment on the data's modality (i.e., bimodal, trimodal, etc.).
 - What is the *midrange* of the data?
 - Can you find (roughly) the first quartile (Q_1) and the third quartile (Q_3) of the data?
 - Give the *five-number summary* of the data.
 - Show a *boxplot* of the data.
 - How is a *quantile-quantile plot* different from a *quantile plot*?

- 2.3 Suppose that the values for a given set of data are grouped into intervals. The intervals and corresponding frequencies are as follows:

<i>age</i>	<i>frequency</i>
1–5	200
6–15	450
16–20	300
21–50	1500
51–80	700
81–110	44

Compute an *approximate median* value for the data.

- 2.4 Suppose that a hospital tested the age and body fat data for 18 randomly selected adults with the following results:

<i>age</i>	23	23	27	27	39	41	47	49	50
%fat	9.5	26.5	7.8	17.8	31.4	25.9	27.4	27.2	31.2
<i>age</i>	52	54	54	56	57	58	58	60	61
%fat	34.6	42.5	28.8	33.4	30.2	34.1	32.9	41.2	35.7

- Calculate the mean, median, and standard deviation of *age* and %fat.
 - Draw the boxplots for *age* and %fat.
 - Draw a *scatter plot* and a *q-q plot* based on these two variables.
- 2.5 Briefly outline how to compute the dissimilarity between objects described by the following:
- Nominal attributes
 - Asymmetric binary attributes

- (c) Numeric attributes
 - (d) Term-frequency vectors
- 2.6 Given two objects represented by the tuples (22, 1, 42, 10) and (20, 0, 36, 8):
- (a) Compute the *Euclidean distance* between the two objects.
 - (b) Compute the *Manhattan distance* between the two objects.
 - (c) Compute the *Minkowski distance* between the two objects, using $q = 3$.
 - (d) Compute the *supremum distance* between the two objects.
- 2.7 The *median* is one of the most important holistic measures in data analysis. Propose several methods for median approximation. Analyze their respective complexity under different parameter settings and decide to what extent the real value can be approximated. Moreover, suggest a heuristic strategy to balance between accuracy and complexity and then apply it to all methods you have given.
- 2.8 It is important to define or select similarity measures in data analysis. However, there is no commonly accepted subjective similarity measure. Results can vary depending on the similarity measures used. Nonetheless, seemingly different similarity measures may be equivalent after some transformation.

Suppose we have the following 2-D data set:

	A_1	A_2
x_1	1.5	1.7
x_2	2	1.9
x_3	1.6	1.8
x_4	1.2	1.5
x_5	1.5	1.0

- (a) Consider the data as 2-D data points. Given a new data point, $x = (1.4, 1.6)$ as a query, rank the database points based on similarity with the query using Euclidean distance, Manhattan distance, supremum distance, and cosine similarity.
- (b) Normalize the data set to make the norm of each data point equal to 1. Use Euclidean distance on the transformed data to rank the data points.

2.7 Bibliographic Notes

Methods for descriptive data summarization have been studied in the statistics literature long before the onset of computers. Good summaries of statistical descriptive data mining methods include Freedman, Pisani, and Purves [FPP07] and Devore [Dev95]. For

CPSC 8650 - Assignment 1

Parampreet Singh

Q1 Given -

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

a) Mean of the data:

$$\bar{x} = \frac{1}{n} \left(\sum_{i=1}^n x_i \right)$$

$$= \frac{(13 + 15 + 16 + \dots + 70)}{27}$$

$$= 29.062 \approx 29.96$$

Median \rightarrow 27 observation i.e odd

$$\text{median} = \left(\frac{27+1}{2} \right)^{\text{th}} \text{ term}$$

$$\text{median} = 25$$

b) Mode \rightarrow 25 - frequency 4
30 - frequency 4

\therefore mode is both 25 & 30

L (lower quartile) = 0.25(3295)

c) Mid range of data is

smallest + largest / 2

$$\begin{array}{ccccccccc} 26 & 66 & 66 & 16 & \underline{13 + 70} & 31 & 35 & 21 & 51 \\ 25 & 25 & 24 & 28 & 28 & 38 & 38 & 26 & 26 \\ & & & & 2 & & & & \end{array}$$
$$= 41.5$$

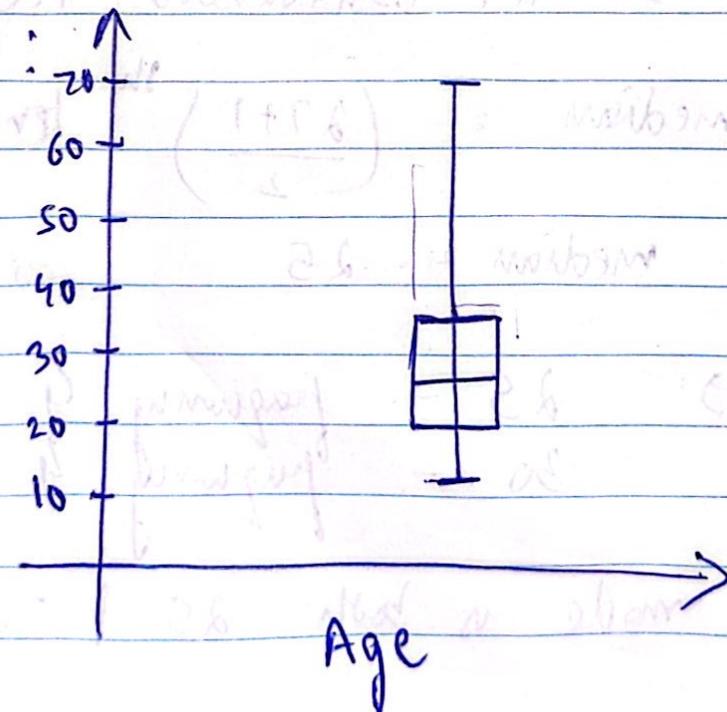
d) ϕ_1 (25th percentile) = 20

ϕ_3 (75th percentile) = 35

e) 5 number summary of data:

minimum	ϕ_1	medium	ϕ_3	maximum
13	20	25	35	70

f) Box plot:



g) Quantile - Quantile is a graphical representation graph which compares data from one set to another set, it is a univariate distribution against the corresponding quantiles of another univariate distribution, both axis display the range of values measured for their corresponding distribution

Quantile plot is a graphical representation method which displays all the data in the increased order which plots quantile information, measured for the independent variable and plotted against their corresponding quantile plot

Ans 2.B :-

Given data grouped into intervals

age	frequency	cumulative frequency
1-5	200	200
6-10	450	650
11-15	300	950
16-20	1500	2450
21-25	700	3150
26-30	44	3194

$$\text{Median} = l + \frac{\frac{N}{2} - (\sum \text{freq})_{\text{med}}}{\text{freq}_{\text{med}}} \times \text{width}$$

$N \rightarrow$ Sum of freq

\rightarrow lower boundary of median class = 21

$$(\sum \text{freq.}) \neq l = 950$$

$$\sum \text{freq}_{\text{med}} = 1500$$

$$\text{width} = 30$$

$$\text{Median} = 21 + \frac{(3194 - 950)}{1500} \times 30$$

$$= 21 + 12.94$$

$$\boxed{\text{Median} = 33.94}$$

Ans 2.4

$$(a) \text{ Mean (sample mean)} \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

for age

$$= \frac{1}{18} (23+23+\dots+61)$$

$$= \frac{836}{18} = 46.44$$

$$\begin{aligned}
 \text{Mean for } J. \text{ fat} &= \frac{1}{18} (9.5 + 26.5 + \dots + 35.7) \\
 &= \frac{518.1}{18} \\
 &= 28.783 \\
 &\approx \underline{\underline{28.78}}
 \end{aligned}$$

$$\text{Median for age} = \frac{50+52}{2} = 51$$

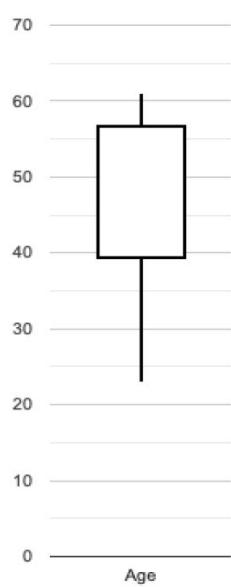
$$\text{Median for } J. \text{ fat} = \frac{30.2 + 31.2}{2} = 30.7$$

standard deviation (σ) for sample age

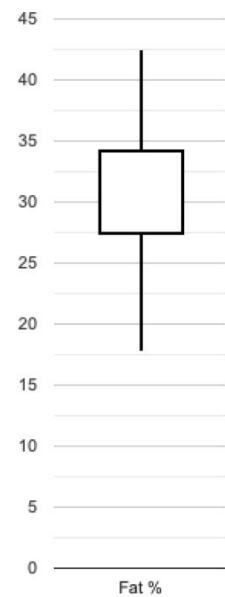
$$\begin{aligned}
 \sigma &= \sqrt{\text{Variance}} \\
 &= \sqrt{\frac{1}{n-1} \sum (x_i - \bar{x})^2} \\
 &= \sqrt{\frac{1}{17} [(23-46.44)^2 + \dots + (61-46.44)^2]} \\
 &= \sqrt{\frac{2970.4447}{17}} \approx 13.2186 \approx 13.2
 \end{aligned}$$

$$\begin{aligned}
 \sigma_{J. \text{ fat}} &= \sqrt{\frac{1}{17} [(9.5 - 28.78)^2 + \dots + (35.7 - 28.78)^2]} \\
 &= \sqrt{\frac{1485.943}{17}} \approx 9.2564 \approx 9.2
 \end{aligned}$$

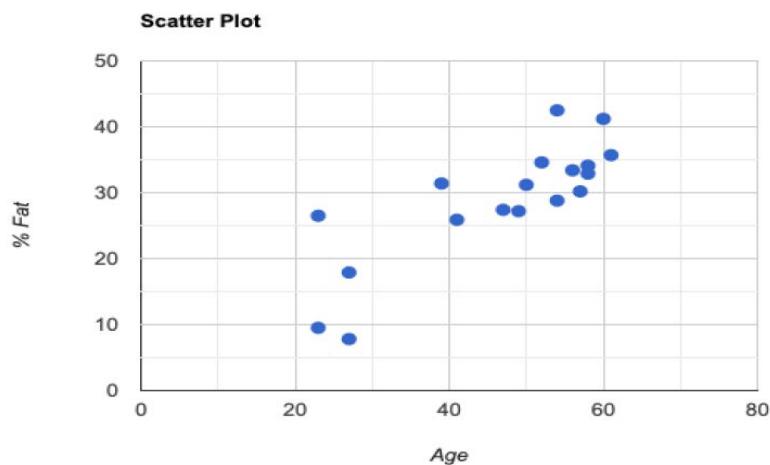
Box Plot for age



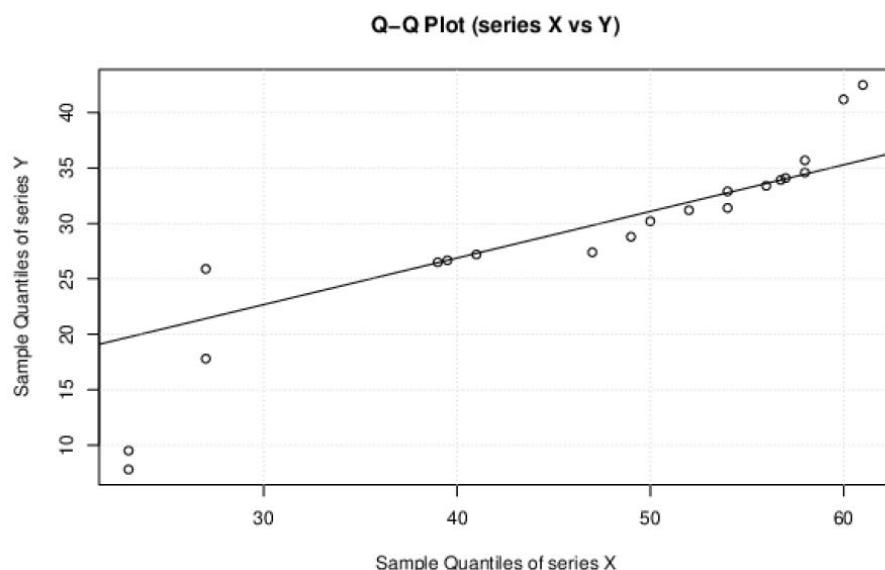
Box Plot for %fat



Scatter Plot



Q-Q Plot



2.5 (a) Nominal attribute:

We can compute dissimilarity b/w objects for the nominal attribute using simple matching method

$$d(i, j) = \frac{p-m}{p} \text{ where } m = \text{no. of matches}$$

$p = \text{total no. of variables}$

(b) Asymmetric binary attribute:

Generally binary attribute values are 0 & 1 both has some meaning and a binary attribute has 2 types to measure the distance they are asymmetric and symmetric

Contingency data table for binary data:

		object d (i)		sum
obj d (j)	1	0		
	1	q	r	$q+r$
0	s	t		$s+t$

distance measure for asymmetric binary variables are

$$d(i, j) = \frac{r+s}{r+s+t}$$

c) Numeric attribute -

we can find the dissimilarity b/w objects using the numeric attribute with minkowski distance with p dimensional data objects and further we can find using (-) norm (Manhattan), L-2 norm (Euclidean), supremum.

Minkowski distance

$$d(i, j) = \sqrt[n]{|x_{i1} - x_{j1}|^p + |x_{i2} - x_{j2}|^p + \dots + |x_{ip} - x_{jp}|^p}$$

where $i = x_{i1}, x_{i2}, \dots, x_{ip}$

$j = x_{j1}, x_{j2}, \dots, x_{jp}$

are 2-D data objects, n is the order
(The distance so defined is also called L-n norm)

L-1 norm (Manhattan Distance)

$$d(i, j) = |x_{i1} - x_{j1}| + |x_{i2} - x_{j2}| + \dots + |x_{ip} - x_{jp}|$$

L-2 norm (Euclidean Distance)

$$d(i, j) = \sqrt{|x_{i1} - x_{j1}|^2 + |x_{i2} - x_{j2}|^2 + \dots + |x_{ip} - x_{jp}|^2}$$

If when $n \rightarrow \infty$,
then it is called supremum distance
(ℓ_∞ norm or L_∞ norm)

This is the maximum difference b/w
any component attribute of the vectors

$$\begin{aligned} d(i, j) &= \lim_{n \rightarrow \infty} \left(\sum_{f=1}^p |x_{if} - x_{jf}|^n \right)^{1/n} \\ &= \max_f |x_{if} - x_{jf}| \end{aligned}$$

(d) Term frequency vectors:

it is also called cosine similarity

where it counts the frequency of the word
particularly used in document.

It will use cosine measure to calculate the similarity b/w the documents. It is also known as term frequency vector

$$\cos(d_1, d_2) = \frac{d_1 \cdot d_2}{\|d_1\| \|d_2\|}$$

Q. b Given 2 objects represented by the tuples $(22, 1, 42, 10)$ & $(20, 0, 36, 8)$

(a) Euclidean distance

$$d(i, j) = \sqrt{(x_{i1} - x_{j1})^2 + (x_{i2} - x_{j2})^2 + \dots + (x_{ip} - x_{jp})^2}$$

$$= \sqrt{(22 - 20)^2 + (1 - 0)^2 + (42 - 36)^2 + (10 - 8)^2}$$

$$= \sqrt{4 + 1 + 36 + 4}$$

$$= \sqrt{45}$$

$$= 6.708203$$

$$\approx 6.71$$

(b) Manhattan distance

$$d(i, j) = |x_{i1} - x_{j1}| + |x_{i2} - x_{j2}| + \dots + |x_{ip} - x_{jp}|$$

$$= |122 - 20| + |11 - 0| + |42 - 36| + |10 - 8| \\ = 11$$

(c) Minkowski Distance for $q=3$

$$d(i, j) = \sqrt[3]{|x_{i1} - x_{j1}|^3 + |x_{i2} - x_{j2}|^3 + |x_{i3} - x_{j3}|^3 + |x_{in} - x_{jn}|^3}$$

$$= \sqrt[3]{(122 - 20)^3 + (11 - 0)^3 + (42 - 36)^3 + (10 - 8)^3}$$

$$= \sqrt[3]{233}$$

$$= 6.1534 \approx 6.15$$

d) supremum distance

$$d(i, j) = \max_f |x_{if} - x_{jf}|$$

$$= |42 - 36|$$

$$= 6$$

Aus 2.7 So median = $L_1 + \left(\frac{N_c - (\text{freq})d}{\text{freq medians}} \right) \text{width}$

where L_1 is the lower boundary of the median interval, N_c is the cumulative frequency of the entire dataset and $(\leq \text{freq})d$ is the sum of the frequency

above the median interval, f_{median} is the frequency of the median's interval and width is the width of the interval. The most straight forward way is to divide all data into k equal length intervals and the error incurred will be decreased as k becomes larger & larger, however time used will also increase.

So the error made and time used is a good optimality measure, and there are some other approaches to approximate median, one of them would be hierarchically divide the whole dataset into intervals, firstly divide into k regions, then find the regions in which median resides then secondly divide this particular region into k sub regions find the sub region in which median resides and the same procedure goes on. we will repeatedly do this, until the width of the sub region reaches the predefined sub threshold then median approximation formula as above mentioned is applied. This would be faster than globally partitioning all data into

shorter intervals which is expensive.

Ans 2.8

(a) New data point $x = (1.4, 1.6)$

Euclidean distance

$$d(1, x) = \sqrt{(1.5 - 1.4)^2 + (1.7 - 1.6)^2} = 0.1414$$

$$d(2, x) = \sqrt{(2 - 1.4)^2 + (1.9 - 1.6)^2} = 0.6708$$

$$d(3, x) = \sqrt{(1.6 - 1.4)^2 + (1.8 - 1.6)^2} = 0.2828$$

$$d(4, x) = \sqrt{(1.2 - 1.4)^2 + (1.5 - 1.6)^2} = 0.2336$$

$$d(5, x) = \sqrt{(1.5 - 1.4)^2 + (1.0 - 1.6)^2} = 0.6083$$

Manhattan Distance

$$d(1, x) = |1.5 - 1.4| + |1.7 - 1.6| = 0.2$$

$$d(2, x) = |2 - 1.4| + |1.9 - 1.6| = 0.9$$

$$d(3, x) = |1.6 - 1.4| + |1.8 - 1.6| = 0.4$$

$$d(4, x) = |1.2 - 1.4| + |1.5 - 1.6| = 0.3$$

$$d(5, x) = |1.5 - 1.4| + |1.0 - 1.6| = 0.7$$

Supremum dist

$$d(1, x) = 0.1$$

$$d(2, x) = 0.6$$

$$d(3, x) = 0.2$$

$$d(4, x) = 0.2$$

$$d(S, n) = 0.6$$

Cosine Similarity

$$\cos(d_1, n) = \frac{(1.5 \times 1.4) + (1.7 \times 1.6)}{\sqrt{1.5^2 + 1.7^2} \times \sqrt{1.4^2 + 1.6^2}} = 0.9999$$

$$\cos(d_2, n) = \frac{(2 \times 1.4) + (1.9 \times 1.6)}{\sqrt{2^2 + 1.8^2} \times \sqrt{1.4^2 + 1.6^2}} = 0.99575$$

$$\cos(d_3, n) = \frac{(1.6 \times 1.4) + (1.8 \times 1.6)}{\sqrt{1.6^2 + 1.8^2} \times \sqrt{1.4^2 + 1.6^2}} = 0.99997$$

$$\cos(d_4, n) = \frac{(1.2 \times 1.4) + (1 \times 1.6)}{\sqrt{1.2^2 + 1^2} \times \sqrt{1.4^2 + 1.6^2}} = 0.96536$$

when comparing values of the above data points, the ranking is below

Euclidean Distance : n_1, n_4, n_3, n_5, n_2

Manhattan Distance : n_1, n_4, n_3, n_5, n_2

Supremum Distance : n_1, n_4, n_3, n_5, n_2

Cosine Similarity : n_1, n_3, n_4, n_2, n_5

b) $x = (1.4, 1.6)$

Normalized value of x is $\frac{x}{\text{norm}(x)}$

Here $\text{norm}(x) = \sqrt{1.4^2 + 1.6^2}$

Normalized x is $\left(\frac{1.4}{\sqrt{1.4^2 + 1.6^2}}, \frac{1.6}{\sqrt{1.4^2 + 1.6^2}} \right) = (0.6573, 0.7512)$

x_1 value $(1.5, 1.7)$ normalized $x_1 = \sqrt{1.5^2 + 1.7^2}$
normalized $x_1 = \left(\frac{1.5}{\sqrt{1.5^2 + 1.7^2}}, \frac{1.7}{\sqrt{1.5^2 + 1.7^2}} \right) = 0.6608$ $x_1 = 2.27$

$x_2 (2, 1.9)$ normalized $(x_2) = \sqrt{2^2 + 1.9^2} = 2.76$

normalized $x_2 = \left(\frac{2}{2.76}, \frac{1.9}{2.76} \right)$
 $= (0.7246, 0.6883)$

Similarly for

	A_1	A_2
x_3	0.6639	0.7464
x_4	0.6250	0.7813
x_5	0.8333	0.5555

So euclidean distance, normalizing π ,
to normalized datapoint

$$\pi_1 = 0.004$$

$$\pi_2 = 0.09$$

$$\pi_3 = 0.007$$

$$\pi_4 = 0.04$$

$$\pi_5 = 0.26$$

So, based on these values
Ranking

$$\hookrightarrow \pi_1, \pi_3, \pi_4, \pi_1, \pi_2, \pi_5$$