

## Mid-term exam solution

### Question 1

- a - what is data mining? Mention steps involved in data Mining when viewed as a Process of Knowledge discovery
- - Data Mining refers to extracting or "Mining" Knowledge from large amounts of data "Knowledge mining of data".
  - Extraction of interesting Pattern or Knowledge from huge amount of data
  - other names :
    - Knowledge discovery (mining) in databases.
    - Knowledge extraction
    - Data / Pattern analysis
    - Information harvesting
    - Business Intelligence.
  - Knowledge discovery is an iterative Sequence of The following steps:
    - 1- Data Cleaning remove noise and inconsistent data.
    - 2- Data Integration where multiple data Sources may be Combined.
    - 3- Data Selection where data relevant to the analysis task are retrieved from the data warehouse.
    - 4- Data Transformation where data are transformed or Consolidated into forms appropriate for mining by performing Summary or aggregation operations, for instance.
    - 5- Data Mining an essential Process where intelligent methods are applied in order to extract data Patterns, Knowledge discovery.
    - 6- Pattern evaluation to identify The truly interesting Patterns representing Knowledge based on some interesting measures.

7- Knowledge Presentation where visualizing and knowledge representation techniques are used to Present the mined knowledge to the users.

b- Suppose that this data shows the marks from a math exam...

Marks	13	12	9	11	14	12	10	15	11	10	7
-------	----	----	---	----	----	----	----	----	----	----	---

i- what is the mean & median of the data.

mean :

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i, \quad \frac{124}{11} = 11.27 \approx 11$$

median : Middle value if odd number of values, or average of the middle two values otherwise.

~~7~~ ~~9~~ ~~10~~ ~~10~~ ~~11~~ 11 ~~12~~ ~~12~~ ~~13~~ ~~14~~ ~~15~~

median = 11

ii- what is the mode of the data

-value that occurs most frequently in the data

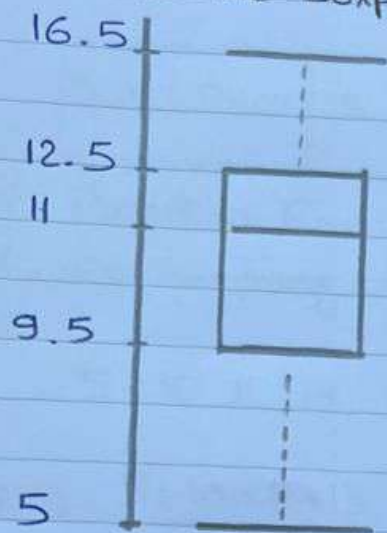
→ 10, 11, 12

iii- find roughly  $Q_1$  and  $Q_3$  of the data.

- Give the five number Summary of the data.



v. Draw The boxplot of the data.



Question - 2 - In Real world, Tuples with missing values for some attributes are a common occurrence. Describe various methods for handling this problem.

Various methods to handle missing data:

- Ignore The Tuple - usually done when class label is missing
  - not effective when the % of values per attribute varies considerably or contains several attributes with missing values
- Fill In Data/value missing manually.
  - This approach is time consuming and may not be a responsible task for large data sets with many missing values, specially when the value to be filled in is not easily determined.
- Fill In automatically with:
  - a global constant
  - the attribute mean

#	value	Count	Percentage	Quantile
1	7	1	$\frac{1}{11}$ 9.09%	
2	9	1	$\frac{2}{11}$ 18.18%	
3	10	2	$\frac{4}{11}$ 36.36%	← $Q_1$ 25%
4	11	2	$\frac{6}{11}$ 54.54%	← $Q_2$ 50%
5	12	2	$\frac{8}{11}$ 72.72%	
6	13	1	$\frac{9}{11}$ 81.81%	← $Q_3$ 75%
7	14	1	$\frac{10}{11}$ 90.90%	
8	15	1	$\frac{11}{11}$ 100%	

$$\square Q_1 = \frac{9+10}{2} = 9.5, \quad Q_2 = 11, \quad Q_3 = \frac{12+13}{2} = 12.5$$

$$\square IQR = Q_3 - Q_1 = 12.5 - 9.5 = 3$$

$$\begin{aligned} \square \text{Maximum Acceptable} &= Q_3 + 1.5 IQR \\ &= 12.5 + 1.5 (3) \\ &= 16.5 \end{aligned}$$

$$\begin{aligned} \square \text{Minimum Acceptable} &= Q_1 - 1.5 IQR \\ &= 9.5 - 1.5 (3) \\ &= 5 \end{aligned}$$

$$\square \text{outlier} = \text{Max} + \text{Min} + IQR, \quad = 16.5 + 5 + 3, \quad = 24.5$$



b - Suppose a group of 12 records has been stored as following:

Sales price	5	10	11	13	15	35	50	55	72	92	204	215
-------------	---	----	----	----	----	----	----	----	----	----	-----	-----

i - Partition Them into 3 bins by equal frequency.

ii - use Smoothing by bin means to smooth These data Partitions.

5	10	11	13	15	35	50	55	72	92	204	215
Avg $\frac{5+10+11+13}{4}$				$\frac{15+35+50+55}{4}$				$\frac{72+92+204+215}{4}$			
$= 9.75$				$= 38.75$				$= 145.75$			

iii - use min max Normalization to normalize These data onto the range  $[0.0, 1.0]$

$$\hat{v} = \frac{v - \min A}{\max A - \min A} \times (\text{new max } A - \text{new min } A) + \text{new min } A$$

$$\hat{v}_5 = \frac{5 - 5}{215 - 5} \times 1 = 0$$

$$\hat{v}_{13} = \frac{13 - 5}{215 - 5} \times 1 =$$

$$\hat{v}_{15} = \frac{15 - 5}{215 - 5} \times 1 =$$

$$\hat{v}_{35} = \frac{35 - 5}{215 - 5} \times 1 =$$

$$\hat{v}_{11} = \frac{11 - 5}{215 - 5} \times 1 =$$

$$\hat{v}_{50} = \frac{50 - 5}{215 - 5} \times 1 =$$

$$Ma \quad \hat{V}_{55} = \frac{55-5}{215-5} \times 1 =$$

$$\hat{V}_{72} = \frac{72-5}{215-5} \times 1 =$$

$$\hat{V}_{92} = \frac{92-5}{215-5} \times 1 =$$

$$\hat{V}_{204} = \frac{204-5}{215-5} \times 1 =$$

$$\hat{V}_{215} = \frac{215-5}{215-5} \times 1 = 1$$

### Question - 3 -

a - explain 3 of The Typical OLAP operations or multidimensional data.

□ Roll up (drill - up) : Summarize data

□ Drill down (Roll down) : reverse of roll - up.

□ Slice and dice : Project and select.

□ Pivot (rotate)

o other operations : drill across → involving more than one fact table  
drill Through → Through the bottom level of the cube to its backend relational tables