

Internal Assessment Test 3 – December 2022

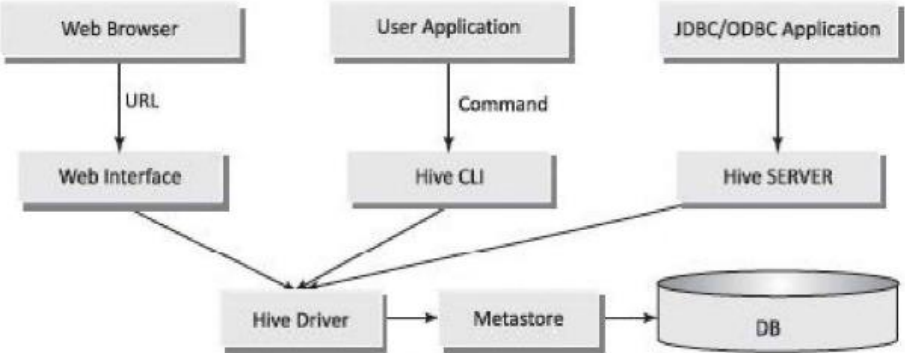
Sub:	BIG DATA AND ANALYTICS				Sub Code:	18CS72	Branch:	ISE	
Date:	27/12/2022	Duration:	90 min's	Max Marks:	50	Sem/Sec:	VII / A, B & C		OBE
<u>Answer any FIVE FULL Questions</u>							MARKS	CO	RBT
1	<p>Illustrate main features and Architecture of Hive with neat diagram. Scheme: Architecture of HIVE with diagram = 10M Solution: Components of Hive architecture are:</p> <ul style="list-style-type: none"> • Hive Server (Thrift) – An optional service that allows a remote client to submit requests to Hive and retrieve results. Requests can use a variety of programming languages. Thrift Server exposes a very simple client API to execute HiveQL statements. • Hive CLI (Command Line Interface) – Popular interface to interact with Hive. Hive runs in local mode that uses local storage when running the CLI on a Hadoop cluster instead of HDFS. • Web Interface – Hive can be accessed using a web browser as well. This requires a HWI Server running on some designated code. The URL <code>http://hadoop:<port no.> / hwi</code> command can be used to access Hive through the web. • Metastore – It is the system catalog. All other components of Hive interact with the Metastore. It stores the schema or metadata of tables, databases, columns in a table, their data types and HDFS mapping. • Hive Driver – It manages the life cycle of a HiveQL statement during compilation, optimization and execution. 					[10]	CO4	L2	
 <pre> graph TD WB[Web Browser] -- URL --> WI[Web Interface] UA[User Application] -- Command --> HCLI[Hive CLI] JOA[JDBC/ODBC Application] --> HS[Hive SERVER] WI --> HD[Hive Driver] HCLI --> HD HS --> HD HD --> MS[Metastore] MS --> DB[(DB)] </pre>									
2	<p>Using HiveQL for the following: Scheme: create table + Alter Table = 4+6 = 10M Solution: i) Create table with partition</p>					[10]	CO4	L2	

Table Partitioning

Create a table with Partition using command:

```
CREATE [EXTERNAL] TABLE <table name> (<column name 1>
<data type 1>, ..... )
PARTITIONED BY (<column name n> <data type n> [COMMENT
<column comment>], ...);
```

ii) Add, rename and drop a partition to a table

Rename a Partition in the existing Table using the following command:

```
ALTER TABLE <table name> PARTITION partition_spec
RENAME TO PARTITION partition_spec;
```

Add a Partition in the existing Table using the following command:

```
ALTER TABLE <table name> ADD [IF NOT EXISTS] PARTITION
partition_spec
[LOCATION 'location1'] partition_spec [LOCATION
'location2'] ...;
partition_spec: (p_column = p_col_value, p_column =
p_col_value, ...)
```

Drop a Partition in the existing Table using the following command:

```
ALTER TABLE <table name> DROP [IF EXISTS] PARTITION
partition_spec, PARTITION partition_spec;
```

3

Describe Pig data types and operator: Group, Filter, Limit, Split.
Scheme: Data types+Group+Filter+Limit+Split = 2+2+2+2+2 = 10M
Solution:

Pig data types

Data type	Description	Example
bag	Collection of tuples	{{(1,1), (2,4)}
tuple	Ordered set of fields	(1,1)
map (data map)	Set of key-value pairs	[Number#1]
int	Signed 32-bit integer	10
long	Signed 64-bit integer	10L or 10l
float	32-bit floating point	22.7F or 22.7f
double	64-bit floating point	3.4 or 3.4e2 or 3.4E2
chararray	Char [], Character array	data analytics
bytearray	BLOB (Byte array)	ff00

Operators: Group, Filter, Limit, Split.

[10]

CO4

L2

Group GROUP statement collects records with the same key. There is no direct connection between group and aggregate functions in Pig Latin unlike SQL.

<p>Collects all records with the same value for the provided key into a bag. Then it can pass to aggregate function, if required or do other things with that.</p>	<pre>A = load 'input' as (name: chararray, rollno:long, marks: float); grpds = group A by marks; B = foreach grpds generate name, COUNT(A);</pre>
--------------------------------------------------------------------------------------------------------------------------------------------------------------------	---------------------------------------------------------------------------------------------------------------------------------------------------

Filter FILTER gives a simple way to select tuples from a relation based on some specified conditions (predicate). It is Pig's *select* command.

<p>Loads an entire record, then selects the tuples with marks more than 75 from each record</p>	<pre>A = load 'input' as (name:chararray, rollno:long, marks:float); B = filter A by marks > 75.0;</pre>
<p>Find name (chararray) that do not match a regular expression by preceding the text without a given character string. Output is all names that do not start with P.</p>	<pre>A = load 'input' as (name:chararray, rollno:long, marks:float); B = filter A by not name matches 'P.*';</pre>

Limit LIMIT gets the limited number of results.

<p>Outputs only first five tuples from the relation.</p>	<pre>A = load 'input' as (name: chararray, city: chararray); B = Limit A 5;</pre>
----------------------------------------------------------	-----------------------------------------------------------------------------------

Split SPLIT partitions a relation into two or more relations

<p>Outputs A relation A splits into two relations P and Q</p>	<pre>A = load 'input' as (name:chararray, rollno:long, marks:float); Split A into P if marks >50.0, Q if marks ≤ 50.0;</pre>
---------------------------------------------------------------	---------------------------------------------------------------------------------------------------------------------------------

4 **Describe the web content mining and three phases for web usage mining. Scheme: web content mining + web usage mining = 5+5 =10M**

[10] CO5 L2

Solution:

Web Content Mining:

Web Content Mining is the process of information or resource discovery from the content of web documents across the World Wide Web. Web content mining can be (i) direct mining of the contents of documents or (ii) mining through search engines. They search fast compared to direct method.

- Web content mining relates to both, data mining as well as text mining. Following are the reasons:
- (i) The content from web is similar to the contents obtained from database, file system or through any other mean. Thus, available data mining techniques can be applied to the web.
 - (ii) Content mining relates to text mining because much of the web content comprises texts.
 - (iii) Web data are mainly semi-structured and/or unstructured, while data mining is structured and the text is unstructured.

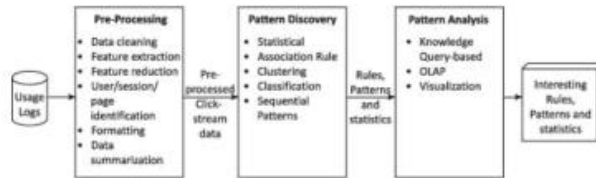
Applications

- Following are the applications of content mining from web documents:
1. Classifying the web documents into categories
 2. Identifying topics of web documents
 3. Finding similar web pages across the different web servers
 4. Applications related to relevance:

Web Usage Mining:

Web usage mining discovers and analyses the patterns in click streams. Web usage mining also includes associated data generated and collected as a consequence of user interactions with web resources.

Figure 9.7 shows three phases for web usage mining.



The phases are:

1. Pre-processing – Converts the usage information collected from the various data sources into the data abstractions necessary for pattern discovery.
2. Pattern discovery – Exploits methods and algorithms developed from fields, such as statistics, data mining, ML and pattern recognition.
3. Pattern analysis – Filter outs uninteresting rules or patterns from the set found during the pattern discovery phase.

Usage data are collected at server, client and proxy levels. The usage data collected at the different sources represent the navigation patterns of the overall web traffic. This includes single-user, multi-user, single-site access and multi-site access patterns.

5 **In Machine learning, explain linear and non linear relationship with essential graphs.**

[10]

CO5

L2

Scheme: linear and nonlinear relationship with graph = 10M

Solution:

Correlation is a statistical technique that measures and describes the ‘strength’ and ‘direction’ of the relationship between two variables.

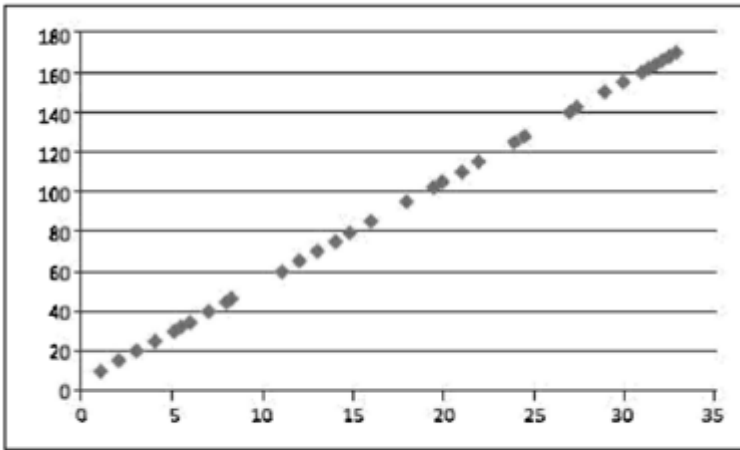
Correlation means analysis which lets us find the association or the absence of the relationship between two variables, x and y . Correlation gives the strength of the relationship between the model and the dependent variable on a convenient 0-100% scale.

R-Square R is a measure of correlation between the predicted values y and the observed values of x . R -squared (R^2) is a goodness-of-fit measure in linear-regression model. It is also known as the coefficient of determination. R^2 is the square of R , the coefficient of multiple correlations, and includes additional independent (explanatory) variables in regression equation.

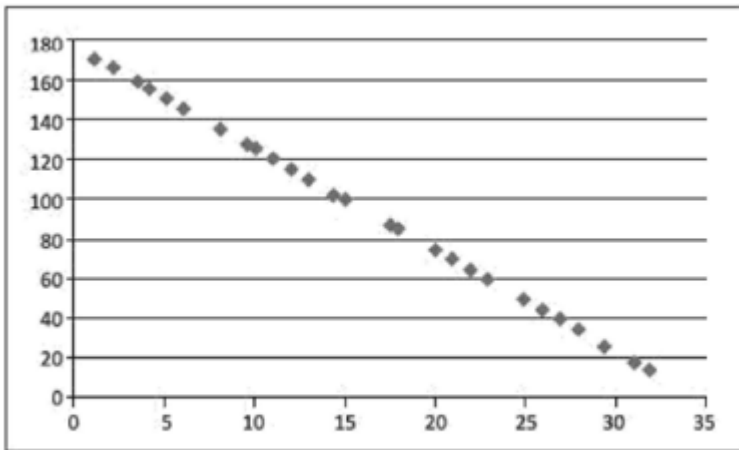
- (i) Constrained Pearson correlation – It is a variation of Pearson correlation that uses midpoint instead of mean rate.
- (ii) Spearman rank correlation – It is similar to Pearson correlation, except that the ratings are ranks.
- (iii) Kendall’s G correlation – It is similar to the Spearman rank correlation, but instead of using ranks themselves, only the relative ranks are used to calculate the correlation.

Table 6.1 The strength of the relationship as a function of r

Value of r	Strength of relationship
-1.0 to -0.5 or 1.0 to 0.5	Strong
-0.5 to -0.3 or 0.3 to 0.5	Moderate
-0.3 to -0.1 or 0.1 to 0.3	Weak
-0.1 to 0.1	None or very weak



Perfect Positive Linear Relationship ($r = 1$)



Perfect Negative Linear Relationship ($r = -1$)

6 How does the Apriori algorithm work? For the following table, find frequent item set and describe the different steps of forming Association rules using Apriori algorithm. Consider minimum support count as 2.

ID	List of Items
T1	I1, I2, I5
T2	I2, I4
T3	I2, I3
T4	I1, I2, I4
T5	I1, I3
T6	I2, I3
T7	I1, I3
T8	I1, I2, I3, I5
T9	I1, I2, I3

Scheme: Apriori algorithm + Association rule = 2+8 = 10M

Solution:

[10]

CO5

L3

The Apriori principle can reduce the number of itemsets needed to be examined. Apriori principle suggests if an itemset is frequent, then all of its subsets must also be frequent. For example, if itemset {A, B, C} is a frequent itemset, then all of its subsets {A}, {B}, {C}, {A, B}, {B, C} and {A, C} must be frequent. On the contrary, if an itemset is not frequent, then none of its supersets can be frequent.

Assume X and Y are two itemsets. Apriori principle holds due to the following property of support measure:

$$\forall X, Y: (X \subseteq Y) \rightarrow s(X) \geq s(Y) \quad (6.24)$$

Explanation: \forall means for all, and \subseteq means 'subset of' and can be 'equal to or included in'. Support of an itemset never exceeds the support of its subsets. This is known as the *anti-monotone property* of support.

Apriori algorithm evaluates candidates for association as follows:

C_k : Set of candidate-itemsets of size k

F_k : Set of frequent itemsets of size k

$F_1 = \{\text{large items}\}$

for ($k=1; F_k \neq \emptyset; k++$) do {

C_{k+1} = New candidates generated from F_k

for each transaction t in the database do

Increment the count of all candidates in C_{k+1} that are contained in t

F_{k+1} = Candidates in C_{k+1} with minimum support

}

Steps of the algorithm can be stated in the following manner:

1. Candidate itemsets are generated using only large itemsets of the previous iteration. The transactions in the database are not considered while generating candidate itemsets.
2. The large itemset of the previous iteration is joined with itself to generate all itemsets having size higher by 1.
3. Each generated itemset that does not have a large subset is discarded. The remaining itemsets are candidate itemsets.

Solution:

Support count = 2

<u>Item</u>	<u>Count</u>
I ₁	6
I ₂	6
I ₃	6
I ₄	2
I ₅	2

Min support count is 2.
So all items selected from
1-candidate itemset to
form 2-candidate itemset.

<u>Item</u>	<u>Count</u>
I ₁ I ₂	4 ✓
I ₁ I ₃	4 ✓
I ₁ I ₄	1 ✗
I ₁ I ₅	2 ✓
I ₂ I ₃	4 ✓
I ₂ I ₄	2 ✓
I ₂ I ₅	2 ✓
I ₃ I ₄	0 ✗
I ₃ I ₅	1 ✗
I ₄ I ₅	0 ✗

Min support count is 2.
4 itemset pruned. to form
3-candidate itemset.
selected itemsets are

I₁ I₂
I₁ I₃
I₁ I₅
I₂ I₃
I₂ I₄
I₂ I₅

<u>Item</u>	<u>Count</u>
I ₁ I ₂ I ₃	2 ✓
I ₁ I ₂ I ₅	2 ✓
I ₁ I ₂ I ₄	1 ✗
I ₁ I ₃ I ₅	1 ✗

2 frequent itemset
selected are

I₁ I₂ I₃
I₁ I₂ I₅

to form 4 frequent itemset

<u>Item</u>	<u>Count</u>
I ₁ I ₂ I ₃ I ₅	1 ✗

I_1, I_2, I_3, I_4 - Support count is 1 so I_1 is not selected. Consider 3 frequent itemset which satisfy min. support count are,

I_1, I_2, I_3
 I_1, I_2, I_4

Association Rule:-

I_1, I_2, I_3

$I_1 \rightarrow I_2, I_3$

$$\text{Confidence} = \frac{\text{Support}\{I_1, I_2, I_3\}}{\text{Support}\{I_1\}} = \frac{2}{6} \times 100 = 33\%$$

$I_2 \rightarrow I_1, I_3$

$$\text{Confidence} = \frac{\text{Support}\{I_1, I_2, I_3\}}{\text{Support}\{I_2\}} = \frac{2}{6} \times 100 = 33\%$$

$I_3 \rightarrow I_1, I_2$

$$\text{Confidence} = \frac{\text{Support}\{I_1, I_2, I_3\}}{\text{Support}\{I_3\}} = \frac{2}{6} \times 100 = 33\%$$

$I_1, I_2 \rightarrow I_3$

$$\text{Confidence} = \frac{\text{Support}\{I_1, I_2, I_3\}}{\text{Support}\{I_1, I_2\}} = \frac{2}{4} \times 100 = 50\%$$

$I_1, I_3 \rightarrow I_2$

$$\text{Confidence} = \frac{\text{Support}\{I_1, I_2, I_3\}}{\text{Support}\{I_1, I_3\}} = \frac{2}{4} \times 100 = 50\%$$

$I_2, I_3 \rightarrow I_1$

$$\text{Confidence} = \frac{\text{Support}\{I_1, I_2, I_3\}}{\text{Support}\{I_2, I_3\}} = \frac{2}{4} \times 100 = 50\%$$

This shows that last three transactions are strong if minimum confidence threshold is 50%.

Similarly need to derive rules for I_1, I_2, I_4 frequent itemset also.