

CIND110  
DATA ORGANIZATION FOR DATA ANALYSTS

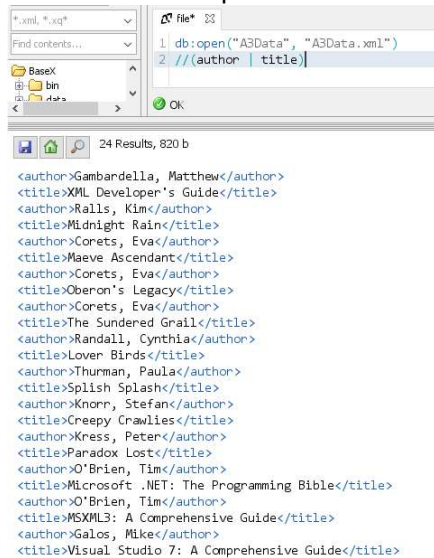
**ASSIGNMENT 2**  
**XML, Information Retrieval, Data Mining, and NoSQL**

SECTION: DK0  
SUBMITTED BY: ANN SAM  
STUDENT NUMBER: 501160843

**1 Section-A:**

**1.1 XML**

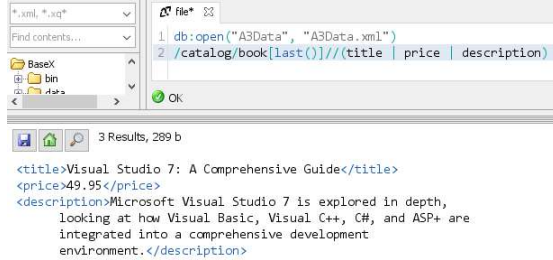
1. Write an XPath expression to find all authors along with their corresponding books.



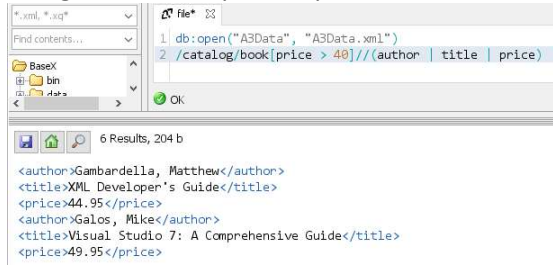
2. Write an XPath expression to find the prices of all the books and their genre.



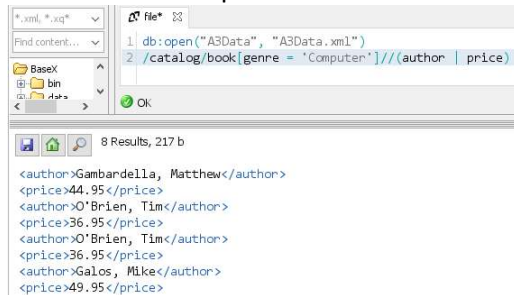
3. Write an XPath expression to find the title, price, and the description in the text of the last book in the catalog.



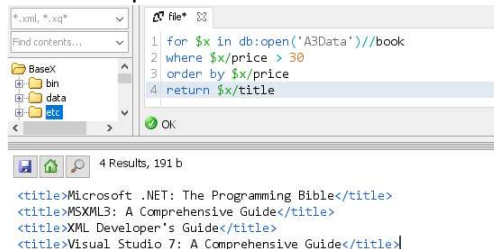
4. Write an XPath expression to find the authors and titles of the books which cost more than 40 dollars, along with the respective prices.



5. Write an XPath expression to find the authors and prices of the books belonging to Computer genre.



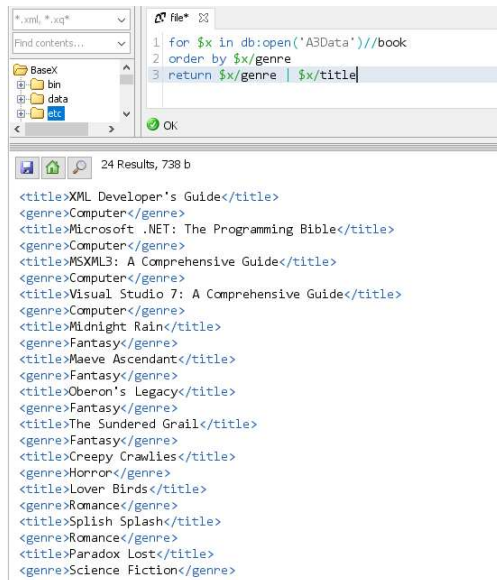
6. Write an XQuery (FLWOR) script to find the titles of the books arranged in ascending order of price, of which the price are more than 30 dollars.



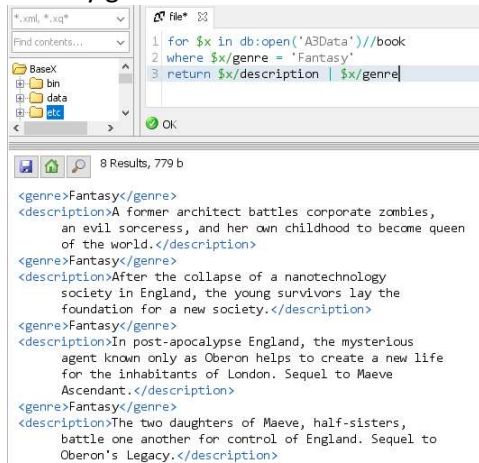
7. Write an XQuery (FLWOR) script to provide only the descriptions of the books which cost less than 5 dollars.



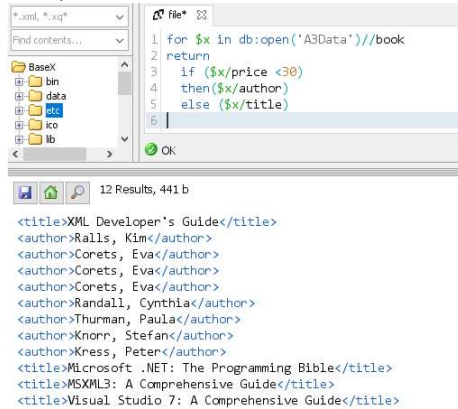
8. Write an XQuery (FLWOR) script which gives the various genre along with the text of the title of the books in them.



9. Write an XQuery (FLWOR) script which gives the description text showing that the books belongs to Fantasy genre.



10. Write an XQuery (FLWOR) script which gives the list of authors whose books cost less than 30 dollars and provides the titles of the books otherwise.



## 1.2 Information Retrieval (IR)

See *SamErb.Ann.rmd* and *SamErb.Ann.html* files.

## 2 Section-B

### 2.1 Data Mining

1. Use the K-means algorithm to cluster this dataset. You can initiate the calculations by assuming K=2 and assume that the records with RIDs 103 and 104 are used as the initial cluster centroids.

Using the Euclidean Distance formula:

$$d(x, y) = \sqrt{\sum_{i=1}^n (y_i - x_i)^2}$$

The distance/proximity was calculated between the data points:

Pairwise Euclidean distance between all pairs of datapoints						
ID	d(ID = 101,j)	d(ID = 102,j)	d(ID = 103,j)	d(ID = 104,j)	d(ID = 105,j)	d(ID = 106,j)
101		28.28	22.36	5.00	5.00	32.02
102	28.28		10.00	32.02	25.00	5.00
103	22.36	10.00		26.93	20.62	11.18
104	5.00	32.02	26.93		7.07	36.06
105	5.00	25.00	20.62	7.07		29.15
106	32.02	5.00	11.18	36.06	29.15	

Example calculations:

$$d(101, 102) = \sqrt{(30 - 50)^2 + (5 - 25)^2} = 28.28$$

$$d(101, 103) = \sqrt{(30 - 50)^2 + (5 - 15)^2} = 22.36$$

$$d(101, 104) = \sqrt{(30 - 25)^2 + (5 - 5)^2} = 5.00$$

Looking at the distance calculations at each initial centroid (103 and 104), we can find the new position of the centroid by calculating the average position of all the points in each cluster. As we have identified the first cluster around datapoints 104, 101, and 105, we can now calculate the average position as follows:

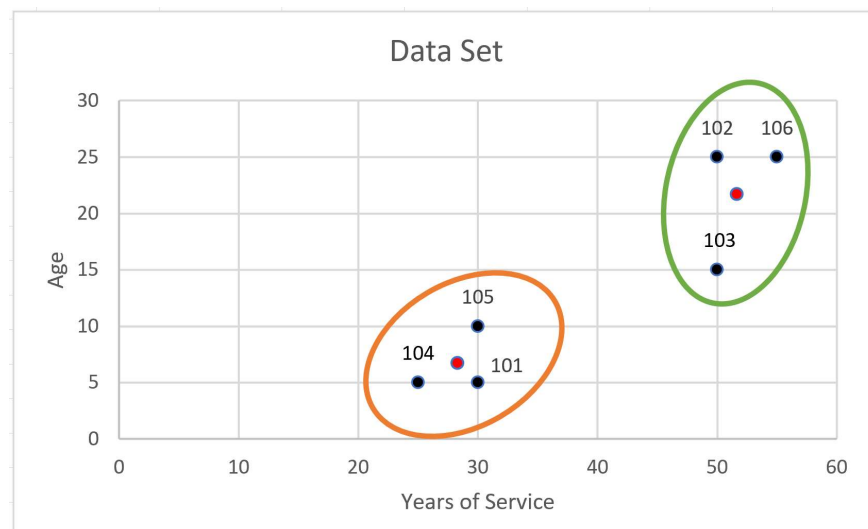
$$centroid_{c1} = \left( \frac{30+25+30}{3}, \frac{5+5+1}{3} \right) = (28.33, 6.67)$$

The same calculation can be done for the second cluster around datapoints 103, 102, and 106:

$$centroid_{c2} = \left( \frac{50+50+55}{3}, \frac{25+15+25}{3} \right) = (51.67, 21.67)$$

We can visually see the two clusters, K=2 when you plot the dataset:

Dataset for K-means Clustering			
ID	age	years	Cluster ID
101	30	5	1
102	50	25	2
103	50	15	2
104	25	5	1
105	30	10	1
106	55	25	2



- The average position of all the points of a cluster (centroid) are plotted in red

2. Provide a brief description on the difference between describing discovered knowledge using clustering and describing it using classification.

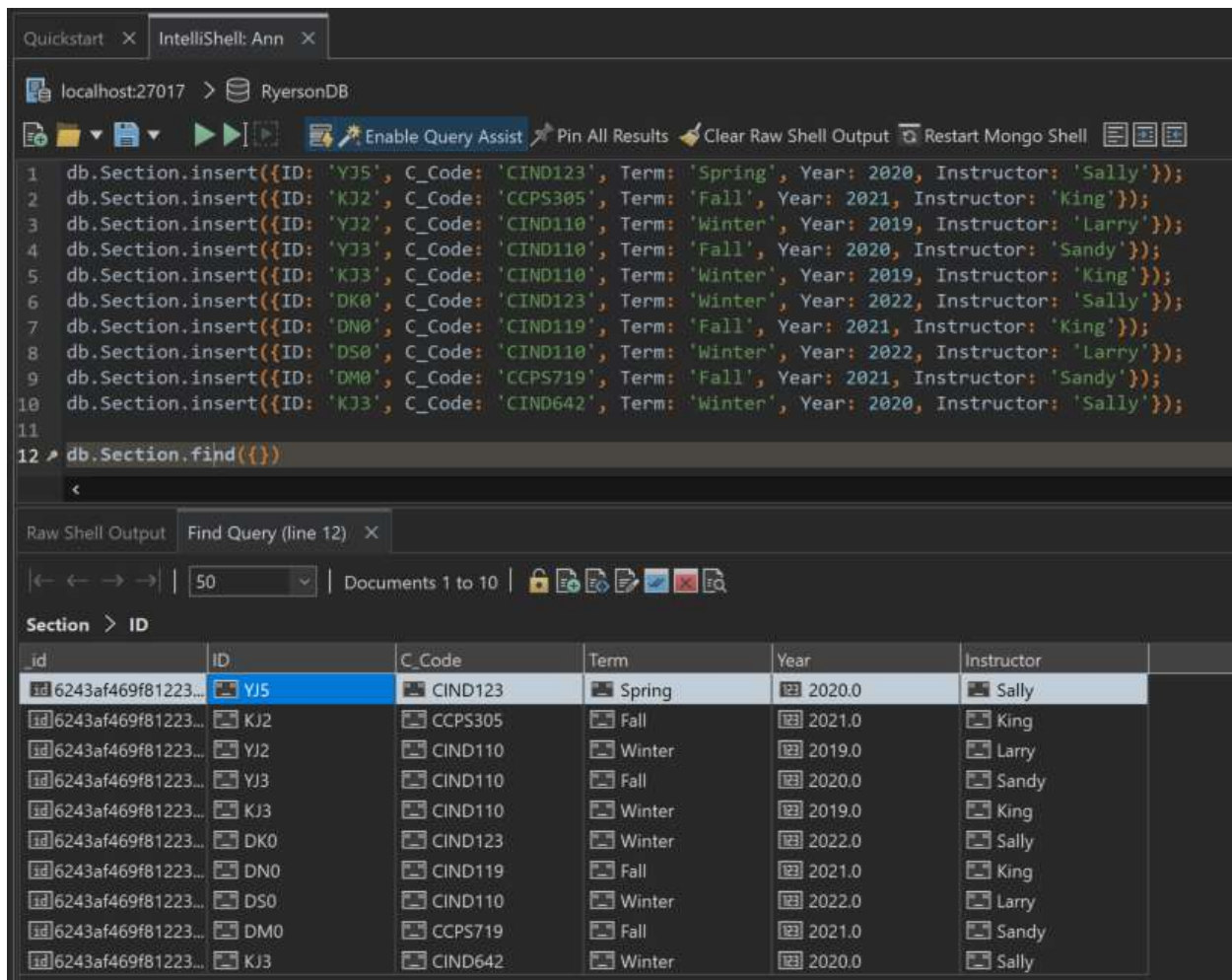
Given the small two-dimensional dataset and lack of predefined classes or labels of the data points, we were able to group the data into two clear clusters. Without any additional data we see that there is a relationship between age and years of service: the younger the age, the less years of service and the older the age, the more years of service.

In terms of classification, any new data points we receive, we can then use relationship found using K-means to determine which class/cluster the new datapoint belongs to. If we were to attempt to classify

this dataset before using the K-means algorithm, it would be much harder to determine a relationship of any between the small dataset. One data analyst could easily classify the data differently from another, while K-means clustering allows a more streamline approach.

## 2.2 NoSQL

### 1. Create the corresponding MongoDB database and Collection with the values.



The screenshot shows the MongoDB Shell interface within IntelliJ IDEA. The database 'RyersonDB' is selected. The shell contains a series of insert commands for the 'Section' collection, followed by a find query. The results of the find query are displayed in a table format.

```

1 db.Section.insert({ID: 'YJ5', C_Code: 'CIND123', Term: 'Spring', Year: 2020, Instructor: 'Sally'});
2 db.Section.insert({ID: 'KJ2', C_Code: 'CCPS305', Term: 'Fall', Year: 2021, Instructor: 'King'});
3 db.Section.insert({ID: 'YJ2', C_Code: 'CIND110', Term: 'Winter', Year: 2019, Instructor: 'Larry'});
4 db.Section.insert({ID: 'YJ3', C_Code: 'CIND110', Term: 'Fall', Year: 2020, Instructor: 'Sandy'});
5 db.Section.insert({ID: 'KJ3', C_Code: 'CIND110', Term: 'Winter', Year: 2019, Instructor: 'King'});
6 db.Section.insert({ID: 'DK0', C_Code: 'CIND123', Term: 'Winter', Year: 2022, Instructor: 'Sally'});
7 db.Section.insert({ID: 'DN0', C_Code: 'CIND119', Term: 'Fall', Year: 2021, Instructor: 'King'});
8 db.Section.insert({ID: 'DS0', C_Code: 'CIND110', Term: 'Winter', Year: 2022, Instructor: 'Larry'});
9 db.Section.insert({ID: 'DM0', C_Code: 'CCPS719', Term: 'Fall', Year: 2021, Instructor: 'Sandy'});
10 db.Section.insert({ID: 'KJ3', C_Code: 'CIND642', Term: 'Winter', Year: 2020, Instructor: 'Sally'});
11
12 db.Section.find({})
  
```

Section	ID	C_Code	Term	Year	Instructor
6243af469f81223...	YJ5	CIND123	Spring	2020.0	Sally
6243af469f81223...	KJ2	CCPS305	Fall	2021.0	King
6243af469f81223...	YJ2	CIND110	Winter	2019.0	Larry
6243af469f81223...	YJ3	CIND110	Fall	2020.0	Sandy
6243af469f81223...	KJ3	CIND110	Winter	2019.0	King
6243af469f81223...	DK0	CIND123	Winter	2022.0	Sally
6243af469f81223...	DN0	CIND119	Fall	2021.0	King
6243af469f81223...	DS0	CIND110	Winter	2022.0	Larry
6243af469f81223...	DM0	CCPS719	Fall	2021.0	Sandy
6243af469f81223...	KJ3	CIND642	Winter	2020.0	Sally



2. In the following, there are five questions to retrieve the output by NoSQL Queries from the Collection created. Write down the corresponding MongoDB codes for each question and provide the output as well.

- a. Find the course codes that instructor King taught in each term:

```
db.Section.find({Instructor: 'King'}, {C_Code: 1, Term: 1, Instructor: 1});
```

13  
14 // Find the course codes that instructor King taught in each term //  
15 db.Section.find({Instructor: 'King'}, {C\_Code: 1, Term: 1, Instructor: 1});

Raw Shell Output Find Query (line 15) X

50 Documents 1 to 3

Section > C\_Code

_id	C_Code	Term	Instructor
6243af469f81223...	CCPS305	Fall	King
6243af469f81223...	CIND110	Winter	King
6243af469f81223...	CIND119	Fall	King

- b. Find out the courses and their instructors after year 2020 onward:

```
db.Section.find({Year: {$gt:2020}}, {C_Code:1, Instructor:1});
```

16  
17 // Find out the courses and their instructors after year 2020 onward //  
18 db.Section.find({Year: {\$gt:2020}}, {C\_Code:1, Instructor:1});

Raw Shell Output Find Query (line 18) X

50 Documents 1 to 5

Section > C\_Code

_id	C_Code	Instructor
6243af469f81223...	CCPS305	King
6243af469f81223...	CIND123	Sally
6243af469f81223...	CIND119	King
6243af469f81223...	CIND110	Larry
6243af469f81223...	CCPS719	Sandy

## c. Find out the course taught in 2021 Fall Term:

```
db.Section.find({Term: 'Fall', Year: 2021}, {C_Code: 1, Term: 1, Year:1});
```

```
19
20 //Find out the courses taught in 2021 Fall Term //
21 db.Section.find({Term: 'Fall', Year: 2021}, {C_Code: 1, Term: 1, Year:1});
```

Raw Shell Output Find Query (line 21) X

Documents 1 to 3

Section > C\_Code

_id	C_Code	Term	Year
6243af469f81223...	CCPS305	Fall	2021.0
6243af469f81223...	CIND119	Fall	2021.0
6243af469f81223...	CCPS719	Fall	2021.0

## d. Find out distinct instructors of the course:

```
db.Section.distinct("Instructor");
```

```
22
23 // Find out distinct instructors of the course: //
24 db.Section.distinct("Instructor");
```

Raw Shell Output Shell Output (Array) X

Documents 1 to 1

Array

[Index]	
0	King
1	Larry
2	Sally
3	Sandy

## e. Find distinct courses taught in the program by grouping them as per their codes:

```
db.Section.distinct("C_Code").sort();
```

```
25
26 //Find distinct courses taught in the program by grouping them as per their codes //
27 db.Section.distinct("C_Code").sort();
```

Raw Shell Output Shell Output (Array) X

Documents 1 to 1

Pin Result Table V

Array

[Index]	
0	CCPS305
1	CCPS719
2	CIND110
3	CIND119
4	CIND123
5	CIND642