

# Information Storage and Retrieval

Tuesday, December 09, 2025 3:42 PM

## Unit 06: Advanced Information Retrieval

Nov – Dec 2022

Q7) a) Define Recommender system? Explain in brief Collaborative Filtering.  
Ans.

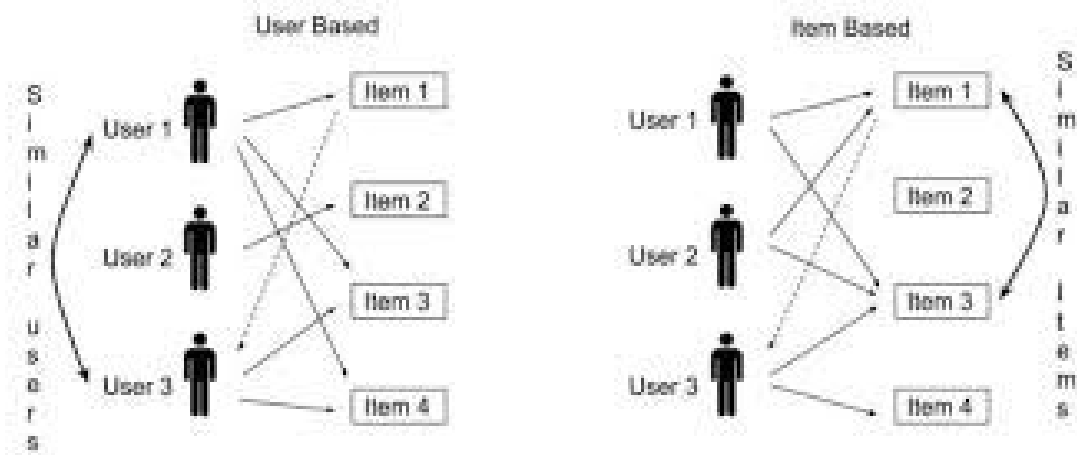
### What is a Recommender System?

A **Recommender System** is an information filtering tool that suggests items to users based on their preferences, behavior, or similarities with other users/items.

It is widely used in platforms like **Netflix (movie recommendations)**, **Amazon (product suggestions)**, **Spotify (music recommendations)**, etc.

### Purpose of a Recommender System

- To help users discover relevant items from a large collection.
- To improve user satisfaction and engagement.
- To personalize user experience.



### What is Collaborative Filtering?

**Collaborative Filtering (CF)** is one of the most common approaches used in recommender systems.

It makes predictions about a user's interests by **collecting preferences from many other users**.

The key idea is:

*"Users who agreed in the past tend to agree in the future."*

### Types of Collaborative Filtering

#### 1. User-Based Collaborative Filtering

- Finds users similar to the target user.

- Recommends items those similar users liked.
- Example: If User A and User B like the same movies, and User B likes another movie X, then recommend X to User A.

## 2. Item-Based Collaborative Filtering

- Finds similarity between items.
- Recommends items similar to the ones the user already likes.
- Example: If many users who watched Movie A also watched Movie B, then Movie B is recommended to a user who watched A.

## How Collaborative Filtering Works (Steps)

1. **Collect user-item interaction data**
  - Ratings, clicks, purchases, watch history.
2. **Calculate similarity**
  - Between users (user-based) or items (item-based).
  - Using cosine similarity, Pearson correlation, etc.
3. **Predict rating or preference**
  - Based on similar users/items.
4. **Recommend top items**
  - Highest predicted ratings or relevance.

## Advantages of Collaborative Filtering

- Does not require item information/content.
- Learns complex user preferences automatically.
- Provides personalized recommendations.

## Limitations

- **Cold Start Problem**  
New users/items have little data, making recommendations hard.
- **Sparsity**  
Most users rate only a few items, causing sparse matrices.
- **Scalability**  
Large datasets require heavy computation.

Q7) b) Explain semantic web in details.

Ans.

## What is the Semantic Web?

The **Semantic Web** is an extension of the current World Wide Web that aims to make data **understandable by machines**.

Instead of just displaying information for humans to read, the Semantic Web allows computers to **interpret**, **share**, and **connect** data automatically.

It was proposed by **Tim Berners-Lee**, the inventor of the Web.

## Key Idea

“Give meaning (semantics) to web data so machines can understand relationships between things.”

This transforms the web from a collection of documents into a **web of interconnected data**.

## Why do we need the Semantic Web?

- People understand meaning, but machines do not.
- Websites use different formats and structures.
- With semantics, machines can **combine information** from many sources.

Example:

A machine could understand that:

- “*Paris is the capital of France*”
  - “*France is in Europe*”
- So it can infer: “*Paris is a city in Europe.*”

## How the Semantic Web Works

It uses technologies that add **metadata** (data about data) to web pages so machines can understand context.

### Main Components

#### 5. RDF (Resource Description Framework)

- A model to represent information in triples:

**Subject – Predicate – Object**

Example:

Paris – isCapitalOf – France

#### 6. OWL (Web Ontology Language)

- Defines concepts, relationships, and rules (like a dictionary for machines).

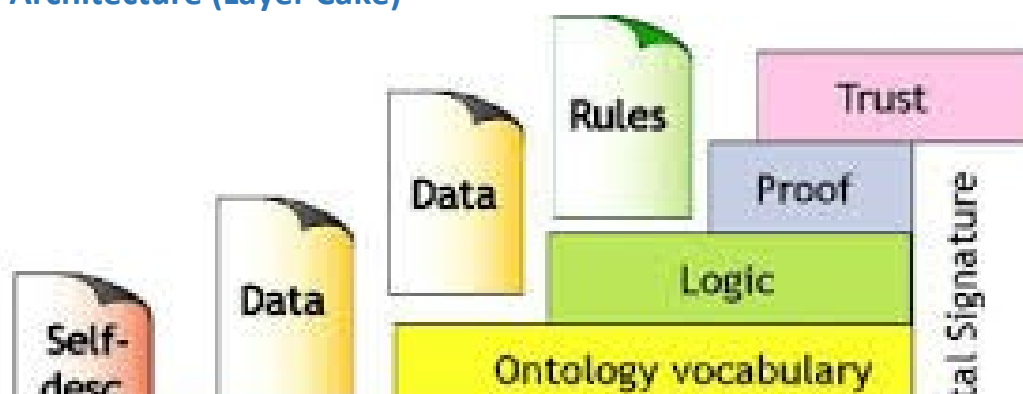
#### 7. SPARQL

- A query language used to retrieve and manipulate semantic data.

#### 8. URIs (Unique Identifiers)

- Identify resources uniquely on the web.

## Semantic Web Architecture (Layer Cake)





This diagram shows the stack of technologies used in the Semantic Web, from basic data encoding to reasoning and trust layers.

### Benefits of the Semantic Web

- Better search results (understanding meaning, not just keywords)
- Data integration across multiple platforms
- Intelligent applications (AI systems)
- Automated reasoning and inference
- More meaningful and connected information

### Example

#### Without Semantic Web:

The word “*Apple*” might mean the fruit or the company—machines cannot tell.

#### With Semantic Web:

Metadata clarifies:

- Apple (Company) → Technology organization
- Apple (Fruit) → Edible fruit

Machines can differentiate and link related knowledge.

### In Short

The Semantic Web turns the web into a **machine-understandable knowledge network**, enabling smarter applications, better search, and interconnected data.

Q8) a) Explain difference between Text-centric and Data-centric XML retrieval.  
Ans.

### Comparison Table: Text-Centric vs Data-Centric XML Retrieval

Feature	Text-Centric XML Retrieval	Data-Centric XML Retrieval
<b>Content Type</b>	Long, narrative, unstructured or semi-structured text	Highly structured, well-defined data
<b>Purpose</b>	Document organization and content markup	Data storage, exchange, and processing
<b>Structure Strictness</b>	Loose structure, flexible	Rigid, schema-driven

<b>Typical Data</b>	Articles, reports, books, blogs	Customer records, product catalogs, transactions
<b>Element Size</b>	Large elements (paragraphs, sections)	Small, atomic data fields (name, price, ID)
<b>Primary Retrieval Goal</b>	Find relevant text sections	Retrieve specific data values
<b>Query Type</b>	Keyword search, full-text search, ranking	Structured queries (XPath, XQuery)
<b>Processing Style</b>	Similar to Information Retrieval (IR)	Similar to Database Management Systems (DBMS)
<b>Tag Function</b>	Helps structure narrative content	Defines data model and relationships
<b>Examples of Use</b>	Digital libraries, news search, document databases	E-commerce databases, banking data, inventory systems

Q8) b) Explain in detail Content Based Recommendation of Documents.

Ans.

### What is Content-Based Recommendation of Documents?

**Content-Based Recommendation** is a method of suggesting documents to a user **based on the content (features)** of documents the user has shown interest in.

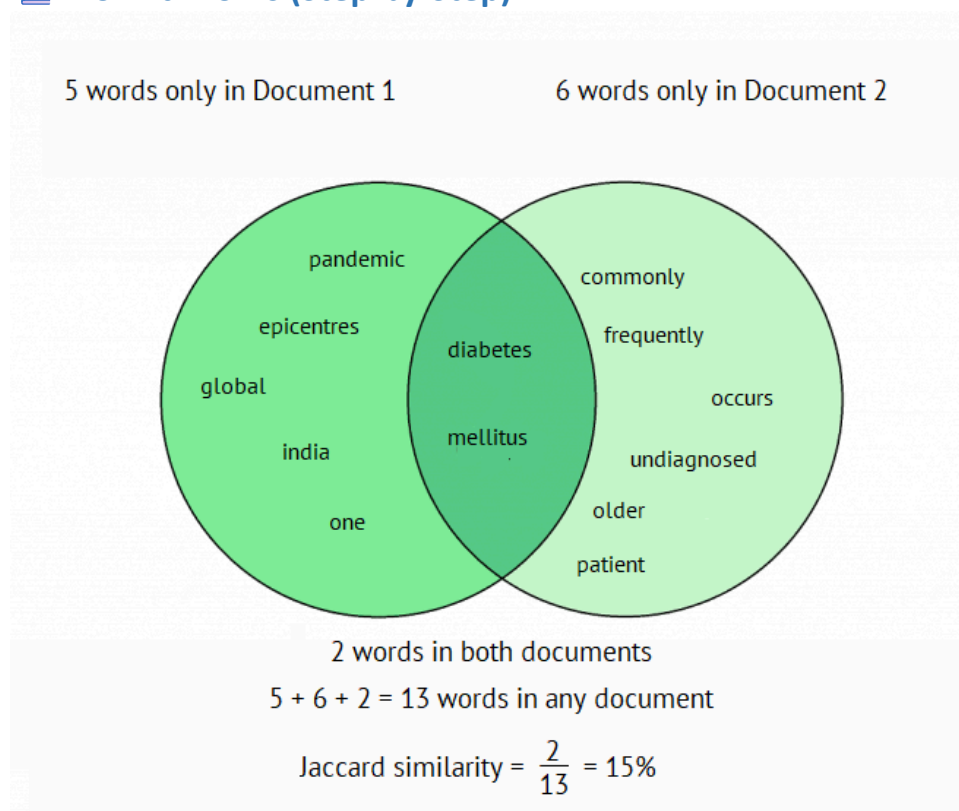
The main idea:

*“Recommend items similar to those the user already likes, based on document content.”*

This approach does **not** depend on other users’ behavior (unlike Collaborative Filtering).

It analyzes the *features* of documents such as keywords, topics, metadata, and relevance.

### How It Works (Step-by-Step)



## 1. Represent Document Content

Documents are converted into a structured form using features such as:

- Keywords
- Term Frequency – Inverse Document Frequency (TF-IDF)
- Bag-of-Words (BoW)
- LDA topic modeling
- Word embeddings (Word2Vec, BERT)

## 2. Build User Profile

The system creates a profile of user preferences based on:

- Documents the user clicked, read, rated, or liked
- The frequent terms/topics in these documents

A user profile might include:

- Preferred topics
- Top keywords
- Semantic vectors

## 3. Compute Similarity

The system measures similarity between:

- The user profile
- and**
- Other documents in the collection

Common similarity measures:

- Cosine similarity
- Jaccard similarity
- Euclidean distance

## 4. Recommend Most Similar Documents

Documents with the highest similarity score are recommended to the user.

## Example

Suppose a user frequently reads documents related to:

- “Machine Learning”
- “Neural Networks”
- “AI Models”

The system identifies these topics in their profile and recommends documents containing similar terms or topics, such as:

- “Deep Learning Overview”
- “Training Neural Networks”
- “Applications of AI in Healthcare”

## Advantages of Content-Based Recommendation

- **No cold-start for items** (new documents can be recommended immediately)
- **Personalized recommendations** based on user's own behavior
- **No need for other user data**
- **Explainable** ("Recommended because it contains similar topics as your previous reads")

## Limitations

- **Limited novelty/diversity**  
Recommends items too similar to what the user already knows.
- **Requires good feature extraction**  
Poorly represented documents → poor recommendations.
- **Cold-start for new users**  
Needs some user history to build a profile.

## In Summary

**Content-Based Document Recommendation** uses document features and user preferences to find and recommend similar documents. It works by analyzing content, building a user profile, and computing similarity between documents.

## May – Jun 2023

Q7) a) Differentiate Collaborative filtering and Content Filtering.

Ans.

### Comparison Table

Feature	Collaborative Filtering (CF)	Content-Based Filtering (CBF)
<b>Basis of Recommendation</b>	Uses preferences/behaviors of <i>similar users</i>	Uses <i>content/features</i> of items the user liked
<b>Requires Other Users' Data?</b>	Yes	No
<b>Key Idea</b>	"People similar to you liked this."	"You liked this item; here are similar ones."
<b>Data Used</b>	Ratings, clicks, interactions from many users	Item attributes (keywords, genres, topics, metadata)
<b>Cold Start Issues</b>	Struggles with new users/items (lack of data)	New items can be recommended easily; new users still a problem
<b>Novelty / Diversity</b>	Can recommend unexpected items based on community trends	Tends to recommend similar items → limited novelty
<b>Scalability</b>	Harder with huge user–item matrices (sparse data)	More scalable because it uses item features

<b>Explainability</b>	Harder to explain (“similar users liked it”)	Easy to explain (“shares features with items you liked”)
<b>Example Scenario</b>	Netflix recommends a movie because many users with similar taste watched it	Spotify recommends songs similar to the one you played
<b>Dependency on Item Features</b>	Not required	Essential (needs good feature extraction)

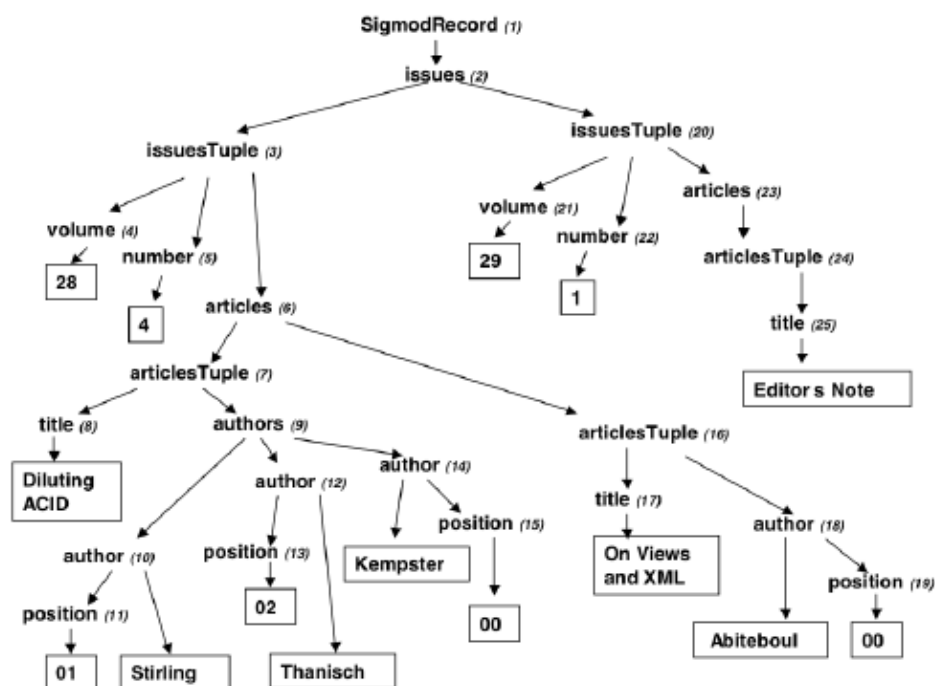
### Simple One-Line Difference

- **Collaborative Filtering:** *Recommends items based on similar users’ preferences.*
- **Content-Based Filtering:** *Recommends items similar to those the user already likes.*

Q8) a) Explain Text-Centric and Data-Centric XML retrieval.

Ans.

## 1. Text-Centric XML Retrieval



### Definition

Text-centric XML retrieval deals with XML documents where the **primary content is text**—often long, narrative, or semi-structured.

XML tags are used mainly for **marking sections**, not for enforcing strict data structures.

### Characteristics

- Contains **large blocks of human-readable text** (paragraphs, articles, reports).
- Structure is **loose and flexible**.
- Retrieval resembles **Information Retrieval (IR)**.
- Uses ranking techniques like keyword matching, TF-IDF, etc.

### Examples

- News articles
- Research papers

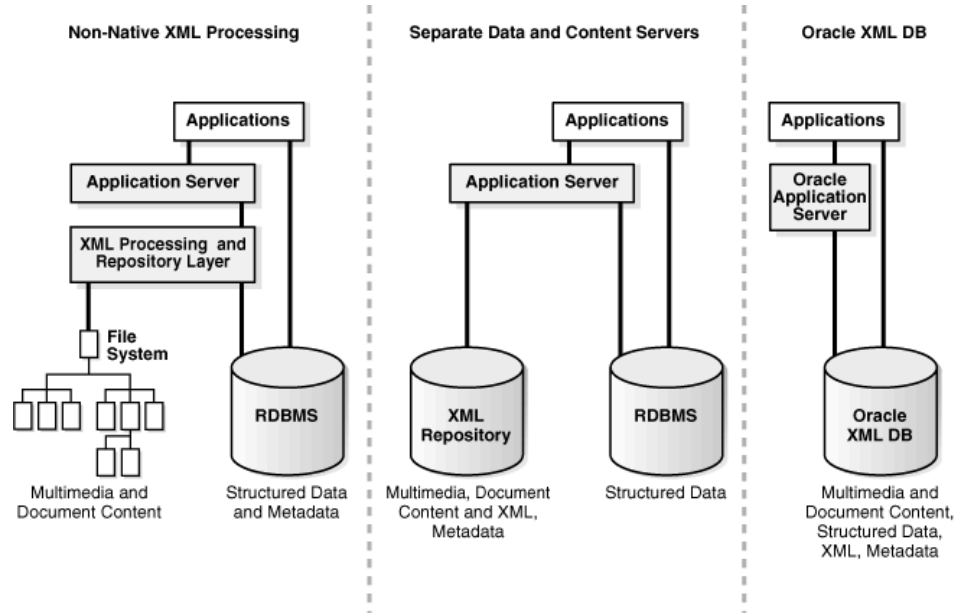


- Blogs
- Digital library content

### Retrieval Goal

- Find **relevant text fragments** or documents based on keywords or topics.

## 2. Data-Centric XML Retrieval



### Definition

Data-centric XML retrieval involves XML documents where information is **highly structured**, similar to database records.

XML is used mainly to store and exchange structured data.

### Characteristics

- Contains **small, well-defined fields** (numbers, names, dates).
- Structure is **strict and follows a schema**.
- Retrieval similar to **database querying**.
- Uses XPath, XQuery, and exact matching.

### Examples

- Banking transactions
- Product catalogs
- Customer databases
- Inventory records

### Retrieval Goal

- Retrieve **specific data values** efficiently and accurately.

## 3. Key Differences (Summary Table)

Feature	Text-Centric XML Retrieval	Data-Centric XML Retrieval
Content Type	Narrative text	Structured data
Structure	Loose, semi-structured	Rigid, schema-driven
Typical Elements	Paragraphs, sections	Small atomic data fields

Query Type	Keyword-based, ranked retrieval	Structured queries (XPath/XQuery)
Similar To	Search engines	Databases
Retrieval Output	Text fragments, whole documents	Exact records/values
Application	Digital libraries	Data storage/exchange systems

## In Simple Words

- **Text-Centric XML** → Used for documents with **lots of text**; retrieved like search engines.
- **Data-Centric XML** → Used for **structured data**; retrieved like database queries.

**Nov– Dec 2023**

Q7) b) Explain Vector Space Model for XML Retrieval.

Ans. Below is a **clear explanation** of the **Vector Space Model (VSM)** for **XML Retrieval**, along with a helpful diagram-style image.

## Vector Space Model (VSM) for XML Retrieval

The **Vector Space Model** is a mathematical model used in Information Retrieval (IR) to represent documents and queries as **vectors in a multi-dimensional space**.

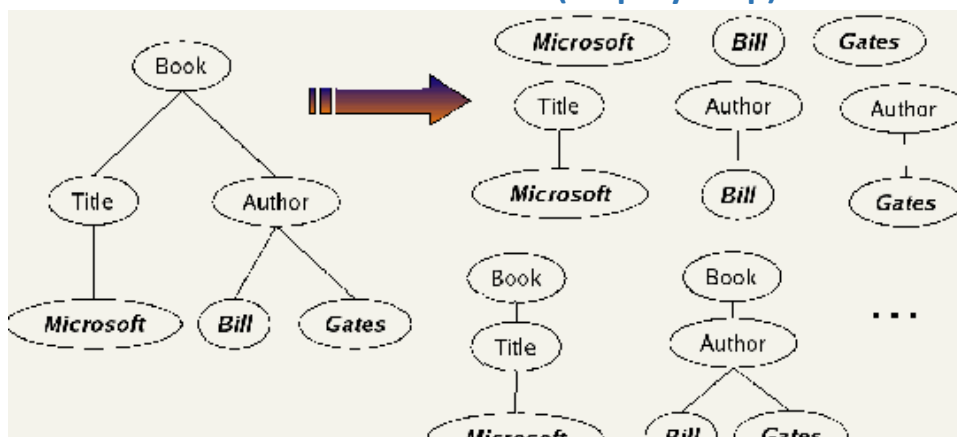
For **XML retrieval**, VSM is adapted to work not just with whole documents but also with **XML elements**, since users may want to retrieve specific sections of an XML document.

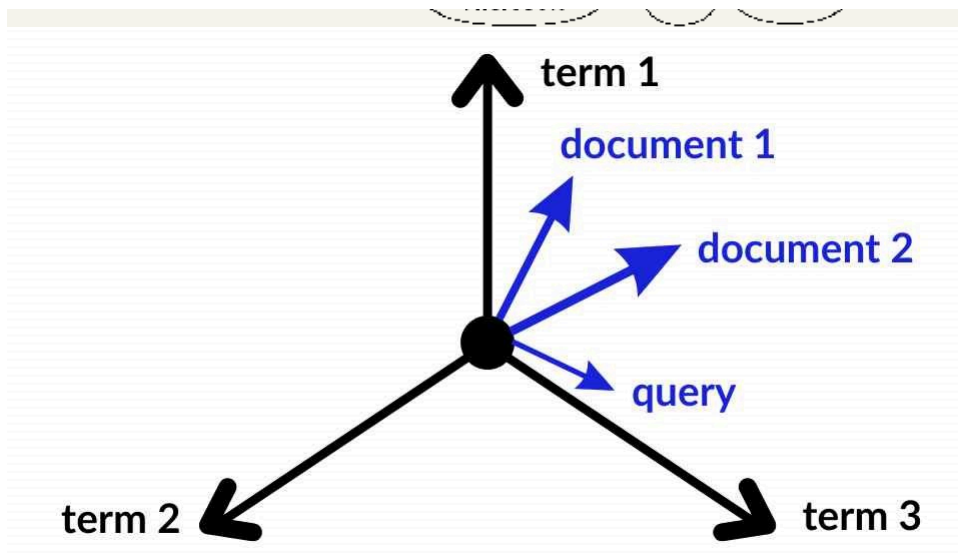
### ✓ Key Idea

*Represent XML elements and queries as vectors, compute similarity, and rank the most relevant XML elements.*

Each dimension represents a **term**, and each XML element is represented as a vector of term weights (often TF-IDF).

## How VSM Works for XML Retrieval (Step-by-Step)





### 1. Indexing XML Documents

XML documents are broken down into:

- Whole documents
- Elements (e.g., <section>, <paragraph>, <article>, etc.)

Each element is treated as a separate retrievable unit.

### 2. Representing XML Elements as Vectors

For each XML element:

- Extract terms (after stemming, stopword removal)
- Compute term weights (TF-IDF)
- Create a vector such as:  
→  $d = (w_1, w_2, w_3, \dots, w_n)$

where (  $w_i$  ) is the weight of term (  $i$  ).

### 3. Representing the Query as a Vector

A query like:

**“climate change impact”**

is represented as another vector with term weights.

### 4. Calculating Similarity (Cosine Similarity)

To rank XML elements, compute:

$$\left[ \begin{array}{l} \text{Similarity}(d,q) = \frac{\vec{d} \cdot \vec{q}}{\|\vec{d}\| \|\vec{q}\|} \end{array} \right]$$

Higher cosine value → more relevant element.

### 5. Ranking and Returning Results

The system returns:

- The most relevant **XML fragments**, not necessarily whole documents.

This makes VSM ideal for **text-centric XML**.

## Why Use VSM in XML Retrieval?

### Advantages

- Supports retrieval of **specific XML parts** (granularity).
- Produces **ranked results** based on similarity.
- Works well for **text-centric XML** (articles, reports).
- Easy to implement with TF-IDF.

### Limitations

- Struggles with **data-centric XML** (numeric/structured data).
- Ignores XML structure unless structural weighting is added.
- Result quality depends on term weighting.

## Example

Suppose an XML article has sections on:

- Climate change
- Weather patterns
- Environmental impacts

A query “climate impact” will return the **specific XML element(s)** most similar to the query vector, not just the whole document.

## In Simple Terms

The **Vector Space Model for XML retrieval**:

- Converts XML elements into vectors
- Converts the user query into a vector
- Measures similarity
- Returns the most relevant XML fragment

It is a powerful method for retrieving **textual XML content**.

**Nov– Dec 2024**

Q7) b) Write short notes on Challenges in XML Retrieval.

Ans.

## Challenges in XML Retrieval

XML Retrieval is more complex than traditional document retrieval because XML contains **hierarchical structure**, **nested elements**, and **mixed content**. Below are the major challenges:

### 1. Structural Complexity

XML documents are organized hierarchically (elements inside elements).

Different documents may use different structures even when representing similar content.

**Challenge:**

- Understanding and matching the structure during retrieval is difficult.
- Similar content may appear at different levels of the hierarchy.

## 2. Granularity of Retrieval

XML allows retrieving:

- Whole documents
- Sections
- Paragraphs
- Small elements

**Challenge:**

- Determining the correct retrieval unit (how big or small the answer should be).
- Ranking results with different sizes is difficult.

Example:

Which should rank higher: a whole article or a single relevant paragraph?

## 3. Mixed Content (Text + Structure)

XML documents contain both **textual content** and **tags**.

**Challenge:**

- How much importance to give to text vs. structure?
- Structural tags may or may not be relevant to the user's query.

## 4. Heterogeneous Schemas

Different publishers or systems may use different tag names or structures for similar data.

**Challenge:**

- Hard to match queries across multiple XML formats.
- Requires schema mapping or ontology knowledge.

Example:

One document uses <author>, another uses <writer> for the same concept.

## 5. Ranking XML Elements

Traditional ranking (like VSM or TF-IDF) works on full documents, but XML returns *fragments*.

**Challenge:**

- How to rank elements of different size and depth?
- Relevance propagation between parent and child elements is complex.

## 6. Overlapping Retrieval Results

Since XML is nested:

- A paragraph may be relevant
- Its section may also be relevant
- The entire article may also be relevant

### Challenge:

- Avoiding duplicate or overlapping results.
- Deciding which level should be shown to the user.

## 7. Query Ambiguity (Content vs Structure)

XML queries may specify:

- Only text terms
- Only structure
- Both (content + structure)

### Challenge:

Interpreting complex queries like:

Find <section> elements containing "climate change"

requires understanding both textual and structural constraints.

## 8. Efficient Indexing

XML has many small elements → millions of retrievable units.

### Challenge:

- Indexing every element is expensive.
- Need special indexes for structure-aware searching (path indexes, element indexes).

## 9. Performance and Scalability

XML repositories can be very large, with deeply nested structures.

### Challenge:

- Complex structural matching slows search.
- Memory and storage overhead is higher compared to flat documents.

## 10. Evaluation Difficulty

Traditional IR metrics (precision, recall) must be adapted.

### Challenge:

- How to evaluate partial or nested results?

- Multiple correct answers may exist at different granularities.

## Summary

XML Retrieval is challenging because XML adds **structure**, **granularity**, and **hierarchy** to traditional text retrieval.

Systems must deal with:

- Variable structures
- Element-level retrieval
- Ranking and overlap reduction
- Schema differences
- Efficient indexing

## Basic XML Concepts

XML (**Extensible Markup Language**) is a markup language designed to store and transport data in a **structured**, **self-descriptive** format. It is widely used for data exchange in web applications, databases, and documents.

### 1. XML Document Structure

An XML document consists of **elements**, **attributes**, **tags**, and **data**.

#### Example XML

```
<book> <title>XML Basics</title> <author>John Doe</author> <price  
currency="USD">29.99</price> </book>
```

### 2. Elements

Elements are the building blocks of XML.

They contain:

- Opening tag
- Content
- Closing tag

#### Example:

```
<title>XML Basics</title>
```

Elements can also contain other elements (nested structure).

### 3. Attributes

Attributes provide **additional information** about elements.

#### Example:

```
<price currency="USD">29.99</price>
```

Here, `currency="USD"` is an attribute of the `price` element.

**Note:**

Attributes store metadata; elements store actual data.

## 4. Tags

Tags define the start and end of an XML element.

- Opening tag: `<author>`
- Closing tag: `</author>`
- Empty tag: `<linebreak/>`

Tags must be:

- Properly nested
- Case-sensitive

## 5. Hierarchical (Tree) Structure

XML organizes data in a **tree structure** with:

- A **root element** (only one)
- Child elements
- Parent-child relationships
- Sibling elements

Example:

```
<library> ← Root <book> ← Child </book> </library>
```

## 6. XML Declaration

Optional line at the beginning of XML:

```
<?xml version="1.0" encoding="UTF-8"?>
```

Specifies XML version and character encoding.

## 7. Well-Formed XML

An XML document is “well-formed” if it follows syntax rules:

- One root element
- Properly nested tags
- Closing tags present
- Case-sensitive tags
- Attribute values in quotes

Example (valid):

```
<note> <to>Alice</to> <from>Bob</from> </note>
```

## 8. Valid XML

An XML document is “valid” if it follows rules defined in:



- **DTD** (Document Type Definition), or
- **XML Schema (XSD)**

Validation ensures structure and data types are correct.

## 9. XML Namespaces

Namespaces avoid naming conflicts when combining XML from different sources.

Example:

```
<bookstore xmlns:fiction="http://example.com/fiction"> <fiction:book>The  
Alchemist</fiction:book> </bookstore>
```

## 10. XML Comments

Comments are written as:

```
<!-- This is a comment -->
```

## 11. Character Data (CDATA)

CDATA sections store text that should not be parsed by XML processors.

Example:

```
<![CDATA[ <h1>This is not XML code</h1> ]]>
```

## 12. XML Advantages

- Human-readable and machine-readable
- Platform-independent
- Supports hierarchical data
- Good for data exchange
- Extensible (user-defined tags)

## In Summary

XML is a flexible, self-descriptive language used to represent structured data.

Understanding elements, attributes, tags, structure, and validation concepts is essential for working with XML in databases, web services, and information retrieval.