

## **PART A**

1. Structure data is organized and follows a predefined format, often in databases or tables, making it easy to search, analyze, and process. It typically includes numbers, dates, and categorical information. Unstructured data, on the other hand, lacks a specific format and is not easily organized, making it more challenging to analyze. It includes text, images, videos, and audio, often requiring advanced techniques like natural language processing or machine learning for meaningful insights.

2. The content of a lookup table is typically considered static data because it is meant to provide reference information that doesn't change frequently or at all during the normal operation of a system or application. Lookup tables are used for mapping values or translating one set of data into another, and their stability ensures consistency and reliability in the functioning of a system. Modifying the content of a lookup table frequently could introduce errors or inconsistencies in the data translation process, making it essential to keep them static for accurate and predictable results.

3 Keyword-based search offers several advantages, including:

1. Simplicity: It's user-friendly and easy for people to understand and use, as it mimics natural language queries.

2. Speed: Keyword searches are often fast, providing quick access to relevant information.

3. Flexibility: Users can search for a wide range of topics and adapt their queries easily.

4. Familiarity: Most internet users are already familiar with keyword-based search engines like Google.

5. Versatility: It can be applied to various types of content, including text, images, and documents.

6. Broad Coverage: It can help uncover information across diverse sources and topics.

7. Scalability: Keyword-based search can handle large datasets and is widely used in web search engines and databases.

4 . Model simulation is a computational technique used to imitate the behavior or processes of a real-world system or phenomenon using a mathematical or computational model. It involves running a computer program or algorithm to generate predictions, observations, or outcomes that mimic what would occur in the actual system under specific conditions. Model simulations are valuable for testing hypotheses, making predictions, and understanding complex systems without the need for real-world experimentation.

5 .Categorization of data involves organizing or classifying information into distinct groups or categories based on common characteristics, attributes, or criteria. This helps in simplifying data management, retrieval, and analysis by grouping similar data together, making it easier to understand and work with.

6 . Converting raw data into a matrix format is essential in data modeling because it enables efficient mathematical and statistical analysis. Matrix representation organizes data into rows and columns, with each row representing an observation or data point, and each column representing a variable or feature. This format is advantageous because:

1. Compatibility: Many machine learning and statistical algorithms require data in a matrix format for processing.

2. Simplicity: It simplifies data manipulation, transformation, and computation.

3. Standardization: Matrix format ensures uniformity in data representation, making it easier to work with.

4. **Mathematical Operations:** Matrices allow for various mathematical operations, such as matrix multiplication, which is fundamental in many data modeling techniques.

Overall, converting raw data to a matrix format streamlines the data modeling process and facilitates a wide range of analytical methods.

7 . The four "V's" of data are:

1. **\*Volume:\*** Refers to the sheer amount of data, including its size and scale.

2. **\*Velocity:\*** Denotes the speed at which data is generated, processed, and analyzed.

3. **\*Variety:\*** Encompasses the diversity of data types and sources, such as text, images, videos, and structured data.

4. **\*Veracity:\*** Refers to the reliability and accuracy of data, highlighting the need to ensure data quality and trustworthiness.

## **PART B**

### **1.DEALING WITH STREAMED DATA**

Dealing with streamed data presents unique challenges compared to working with static or batch data. Streamed data is continuously generated and processed in real-time, which requires careful consideration of various factors to ensure data quality, scalability, and efficiency. Here's a detailed explanation of the key factors to consider when working with streamed data:

1. **\*\*Data Sources and Formats\*\*:**

- **\*\*Source Type\*\*:** Identify the sources of your streamed data. It could be sensor data, logs, social media feeds, or any other continuous data source.

- **\*\*Data Formats\*\*:** Understand the format of the data, such as JSON, XML, CSV, or binary, as this impacts how you parse and process it.

2. **\*\*Data Velocity\*\*:**

- **\*\*Rate of Ingestion\*\*:** Determine the speed at which data is generated and ingested. This can vary from slow to high-velocity streams, requiring different processing strategies.

3. **\*\*Data Volume\*\*:**

- **\*\*Data Size\*\*:** Assess the amount of data generated over time. Handling large volumes of data may require distributed processing and storage solutions.

4. **\*\*Data Variety\*\*:**

- **Structured vs. Unstructured**: Determine whether the data is structured, semi-structured, or unstructured. This affects how you extract and transform the data.

#### 5. **Data Quality**:

- **Data Cleansing**: Implement data cleansing and validation processes to handle incomplete or erroneous data.

- **Duplicate Handling**: Consider how to handle duplicates that may occur due to network issues or retries.

- **Outlier Detection**: Implement mechanisms to detect and handle outliers or anomalies in real-time.

#### 6. **Latency Requirements**:

- **Processing Time**: Define the acceptable latency for processing data. Some applications require low-latency processing, while others can tolerate higher delays.

#### 7. **Scalability**:

- **Horizontal Scaling**: Ensure that your streaming architecture can scale horizontally to handle increased data loads.

- **Load Balancing**: Implement load balancing strategies to distribute the processing load evenly across multiple instances.

#### 8. **Fault Tolerance**:

- **Redundancy**: Design your system with redundancy to ensure high availability and fault tolerance.
- **Checkpointing**: Use checkpointing mechanisms to recover from failures without data loss.

#### 9. **Storage**:

- **Data Retention**: Decide how long to retain streamed data. This affects storage costs and compliance requirements.
- **Storage Format**: Choose an appropriate storage format (e.g., databases, data lakes, or distributed file systems) based on your data processing needs.

#### 10. **Processing Framework**:

- **Stream Processing Framework**: Select a stream processing framework or platform that suits your requirements. Popular choices include Apache Kafka, Apache Flink, Apache Spark Streaming, and AWS Kinesis.
- **Complex Event Processing (CEP)**: Consider using CEP systems for detecting patterns and correlations in the streaming data.

#### 11. **Schema Evolution**:

- Plan for changes in data schema over time. Ensure your system can handle schema evolution gracefully without breaking downstream processing.

#### 12. **Security and Privacy**:

- Implement encryption, access controls, and data masking to protect sensitive data.
- Comply with data privacy regulations like GDPR or HIPAA when processing personal or sensitive information.

#### 13. **Monitoring and Alerting**:

- Set up monitoring and alerting systems to detect anomalies, bottlenecks, and performance issues in real-time.
- Monitor the health of your streaming infrastructure, including message queues, processing nodes, and storage.

#### 14. **Testing and Validation**:

- Develop robust testing strategies for streaming applications, including unit tests, integration tests, and end-to-end testing.
- Perform load testing to ensure your system can handle peak data loads.

#### 15. **Documentation and Documentation**:

- Document the architecture, data flows, and processing logic comprehensively to facilitate collaboration and troubleshooting.

#### 16. **Cost Management**:

- Keep track of the costs associated with streaming data, including storage, processing, and data transfer fees, and optimize accordingly.

## 17. **Compliance**:

- Ensure that your streaming solution complies with relevant industry standards and regulations.

Dealing with streamed data is a complex task, but by considering these factors and making informed decisions, you can build a robust and efficient streaming data processing system that meets your business needs. It's essential to continually monitor and adapt your streaming infrastructure as data volumes and processing requirements evolve over time.

## **2.CONTENT BASED SEARCH**

Content-based search, also known as content-based retrieval or content-based image retrieval (CBIR), is a search technique used in information retrieval systems to find data or documents that are similar in content to a given query. This approach relies on the actual content or features of the data itself, rather than metadata or textual descriptions. Content-based search is commonly used in various domains, including image and video search, music recommendation, and document retrieval. Here's a detailed explanation of how content-based search works:

### **1. Feature Extraction:**

- In content-based search, the first step is to extract relevant features from the data being searched. The type of features extracted depends on the domain. For example:

- **Image Search**: Features can include color histograms, texture descriptors, shape characteristics, or deep learning-based feature vectors (e.g., CNN embeddings).

- **Music Recommendation**: Features might include audio spectrograms, mel-frequency cepstral coefficients (MFCCs), or music fingerprints.

### **2. Feature Representation:**

- Once features are extracted, they need to be transformed into a suitable representation for comparison. This often involves creating feature vectors, which are numerical representations of the extracted features. These vectors capture the essence of the content.

- Feature vectors are typically normalized to ensure that different features contribute equally to the similarity calculation.

### **3. Querying:**

- When a user submits a query, the system extracts features from the query itself, following the same process as for the data.

- The query features are transformed into a feature vector for comparison.

### **4. Similarity Calculation:**

- To perform the actual search, the system calculates the similarity between the query feature vector and the feature vectors of the data items in the database.

- Various distance metrics can be used for this purpose, such as Euclidean distance, cosine similarity, or Manhattan distance.

- The lower the distance or higher the similarity score, the more similar the items are considered to be.

#### **\*\*5. Ranking and Retrieval:\*\***

- The retrieved items are ranked based on their similarity scores. The most similar items are presented to the user as search results.

- The number of results and the ranking strategy can be customized to suit the application's requirements.

#### **\*\*6. User Feedback and Refinement:\*\***

- Content-based search systems often include feedback mechanisms to allow users to provide feedback on the retrieved results.

- User feedback can be used to refine the search results by adapting the query or the weighting of features.

#### **\*\*7. Scalability and Efficiency:\*\***

- Handling large datasets efficiently is a challenge in content-based search. Indexing and caching strategies are often employed to speed up retrieval.

- In distributed systems, distributed databases and parallel processing can be used to enhance scalability.

#### **\*\*8. Continuous Learning:\*\***

- Some content-based search systems incorporate machine learning techniques to improve search quality over time. These systems learn from user interactions and adapt their models accordingly.

#### **\*\*9. Evaluation:\*\***

- Content-based search systems are evaluated using metrics such as precision, recall, and F1-score to assess the accuracy of retrieved results.

#### **\*\*10. Challenges:\*\***

- Variability in data quality, lighting conditions, or noise can affect the performance of content-based search systems.

- Ensuring the system's ability to handle diverse content types and domains is a challenge.

- Balancing the trade-off between feature dimensionality (complexity) and search efficiency is crucial.

In summary, content-based search is a technique that relies on the intrinsic features of data to perform similarity-based searches. It has applications in various fields, including multimedia retrieval, recommendation systems, and information retrieval, and it involves several steps, from feature extraction to similarity calculation and result ranking. The success of content-based search systems often depends on the quality of feature extraction and representation, as well as efficient search algorithms.

### 3. DATA VISUALIZATION:

Data visualization is the graphical representation of data to help people understand the patterns, trends, and insights within the data. It's a crucial part of data analysis and communication because it makes complex data more accessible and interpretable. Various data visualization tools are available to create different types of charts, graphs, and interactive visuals. Here, I'll explain data visualization and provide examples of popular data visualization tools:

#### **\*\*Data Visualization:\*\***

Data visualization serves several purposes, including:

1. **\*\*Exploration:\*\*** Data visualization helps analysts explore datasets to identify trends, outliers, and patterns that may not be immediately apparent in raw data.
2. **\*\*Communication:\*\*** It enables effective communication of data-driven insights to non-technical stakeholders, making it easier for them to grasp complex information.
3. **\*\*Decision-Making:\*\*** Visualizations facilitate data-driven decision-making by providing a clear and concise representation of key information.
4. **\*\*Storytelling:\*\*** Visualizations can be used to tell a compelling data-driven story, making it easier for an audience to understand and remember the message.

#### **\*\*Examples of Data Visualization Tools:\*\***

##### 1. **\*\*Tableau:\*\***

- Tableau is a widely-used data visualization tool that offers a range of interactive visualization options.

- Example: A Tableau dashboard can display sales trends over time, allowing users to filter data by region, product, or time period, and interactively explore the data.

##### 2. **\*\*Microsoft Power BI:\*\***

- Power BI is a business analytics service by Microsoft that provides interactive data visualizations and business intelligence capabilities.

- Example: Power BI can be used to create a dynamic bar chart showing revenue by product category, with filters for date range and geographic region.

##### 3. **\*\*D3.js (Data-Driven Documents):\*\***

- D3.js is a JavaScript library for creating custom and interactive data visualizations on the web.

- Example: A D3.js-based visualization might create a force-directed graph to show relationships between nodes in a network.

##### 4. **\*\*Matplotlib:\*\***

- Matplotlib is a popular Python library for creating static, animated, and interactive visualizations.

- Example: A Matplotlib scatter plot can display the correlation between two variables, such as temperature and ice cream sales.

##### 5. **\*\*ggplot2:\*\***



- ggplot2 is an R package for creating elegant and customizable data visualizations.
- Example: ggplot2 can generate a density plot to visualize the distribution of exam scores for a group of students.

#### 6. **QlikView:**

- QlikView is a business discovery platform that enables users to create and interact with visual data analytics.
- Example: QlikView can generate a pivot table and associated bar chart to analyze sales by product and region.

#### 7. **Plotly:**

- Plotly is a Python graphing library that allows for interactive and web-based visualizations.
- Example: Plotly can create an interactive 3D surface plot to visualize elevation data.

#### 8. **Google Data Studio:**

- Google Data Studio is a free tool for creating interactive dashboards and reports.
- Example: Google Data Studio can create a report that displays website traffic data, including the number of users, sessions, and pageviews.

#### 9. **Excel (with Power Query and Power Pivot):**

- Excel is a familiar spreadsheet tool that includes advanced data visualization features through add-ins like Power Query and Power Pivot.
- Example: Excel can create a pivot chart to visualize sales trends, with data imported and transformed using Power Query.

#### 10. **Infogram:**

- Infogram is an online platform for creating infographics and interactive charts.
- Example: Infogram can be used to create an infographic summarizing survey results, with interactive elements for exploration.

These data visualization tools cater to a range of user preferences, skill levels, and data analysis needs. The choice of tool depends on factors such as data complexity, interactivity requirements, and the target audience for the visualizations.

## 4. WORD FLOCKING ALGORITHM

Word flocking, also known as text flocking or word cloud visualization, is a data visualization technique used to represent text data in a visually appealing and informative way. This technique transforms a collection of text documents or a large body of text into a visual display where words are clustered together, forming a "flock" or cloud-like structure. The size and position of each word in the cloud are determined by various factors related to the frequency and importance of the word in the text data. Word flocking algorithms aim to create word clouds that highlight the most relevant and significant terms within the text.

Here's an explanation of the key elements and steps involved in word flocking algorithms:

### 1. **Text Preprocessing**:

- Before creating a word flock, the input text data is preprocessed. This typically involves tokenization (breaking text into words or phrases), removing stop words (common words like "the," "and," "in" that don't carry significant meaning), and stemming or lemmatization to reduce words to their root form.

### 2. **Word Frequency Calculation**:

- Word flocking algorithms calculate the frequency of each word in the text data. Words that appear more frequently are often given more prominence in the word cloud.

### 3. **Word Importance Scoring**:

- Some word flocking algorithms take into account the importance or significance of words beyond simple frequency. For example, TF-IDF (Term Frequency-Inverse Document Frequency) can be used to assign scores to words based on their frequency in the current document relative to their frequency in the entire corpus of documents. Words with higher TF-IDF scores are considered more important.

### 4. **Word Positioning**:

- The positioning of words within the word cloud is a critical aspect of word flocking. Algorithms determine the x and y coordinates for each word to create an aesthetically pleasing and balanced layout.

- Words with similar meanings or contexts are often clustered together to improve readability and provide semantic meaning to the visualization.

### 5. **Font Size and Color**:

- The size of each word in the word cloud is often proportional to its frequency or importance. Larger words are more prominent.

- Words can also be assigned different colors based on various criteria, such as sentiment or topic, to provide additional insights.

### 6. **Layout Optimization**:

- Word flocking algorithms may employ optimization techniques to fine-tune the layout of the word cloud, ensuring that words do not overlap, and the overall arrangement looks visually appealing.

## 7. **Visualization**:

- Once the word cloud is generated, it is presented as a graphical visualization. Users can interact with the word cloud, such as clicking on words to see their context or obtaining additional information.

## 8. **Interactivity**:

- Many word flocking visualizations are interactive, allowing users to explore the text data further. Interactivity might include zooming in on specific areas of the word cloud, searching for specific words, or filtering words based on criteria like frequency or importance.

Word flocking algorithms are commonly used in various applications, including:

- **Text summarization**: Identifying the most important keywords or phrases in a large text document.
- **Topic modeling**: Visualizing the key topics or themes in a collection of documents.
- **Sentiment analysis**: Highlighting sentiment-related words to understand the sentiment expressed in text data.
- **Content analysis**: Identifying patterns and trends in social media discussions, customer reviews, or news articles.

Word flocking is a visually engaging way to gain insights from text data and make it more accessible to a wider audience. Different word flocking tools and libraries are available for creating custom word cloud visualizations in various programming languages, making it a versatile technique for text analysis.

## 5.ACO:

### Ant colony optimisation

\* It is proposed by Marco Dorigo in 1992

\* It is a Bio-Inspired algorithm which is based on the swarm intelligence.

The main motive of Ant colony optimisation algorithm is Ant should move from ant colony to the food in a shortest path with the help of Pheromones.

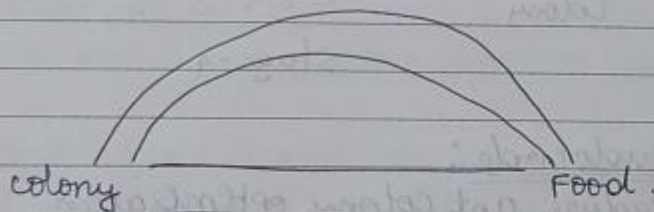
↳ here pheromone is a saliva of ant

### 4 Stages of Ant colony optimisation.

#### Stage - 1

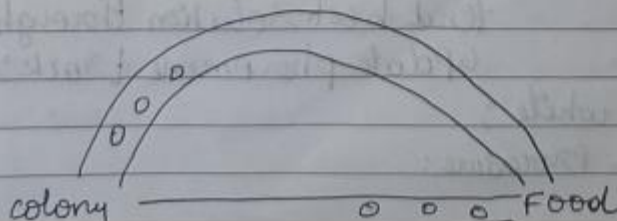
⇒ All ants are in Ant colony.

⇒ No ant is moving in stage 1, so there is no pheromone in the environment since we don't know which path have to be taken.



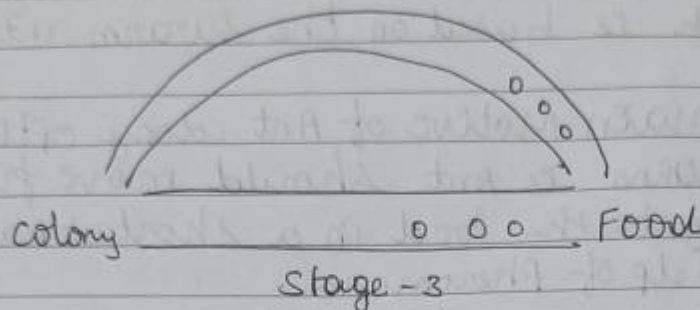
#### Stage - 2

⇒ 0.5% they are trying with different path find the shortest path.



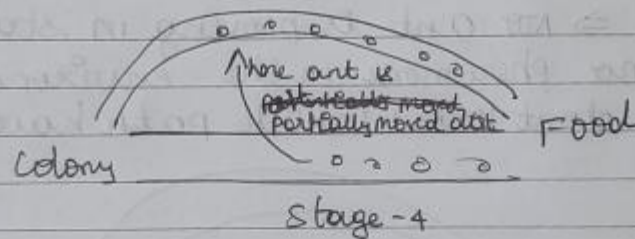
stage-3

in the stage 3 shortest path is found



stage-4

Possible ways are observed & the ant moves into the shortest path.



Pseudocode:

Procedure ant colony optimization

Initialize necessary parameter & Pheromone trails;

while not termination do:

    Generate ant population;

    calculate fitness value association with each ant;

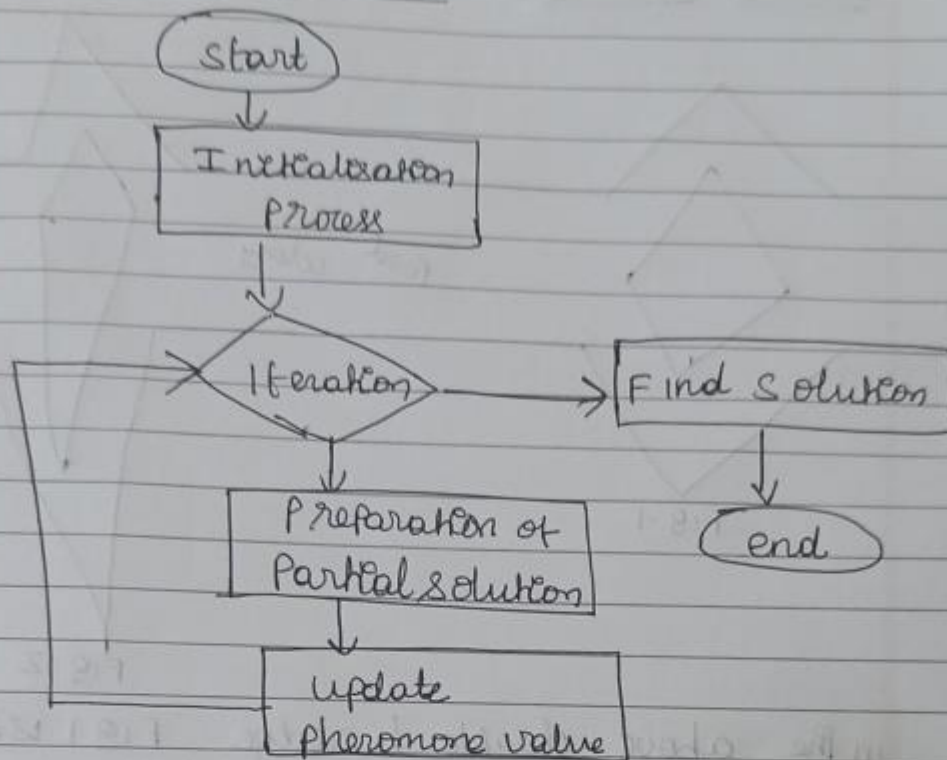
    find best solution through selection method

    update pheromone trails;

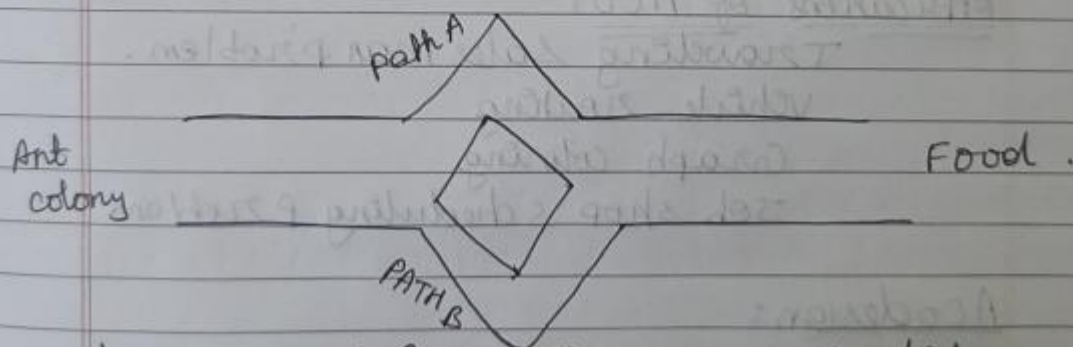
end while;

end procedure;

## Flowchart of ACO.



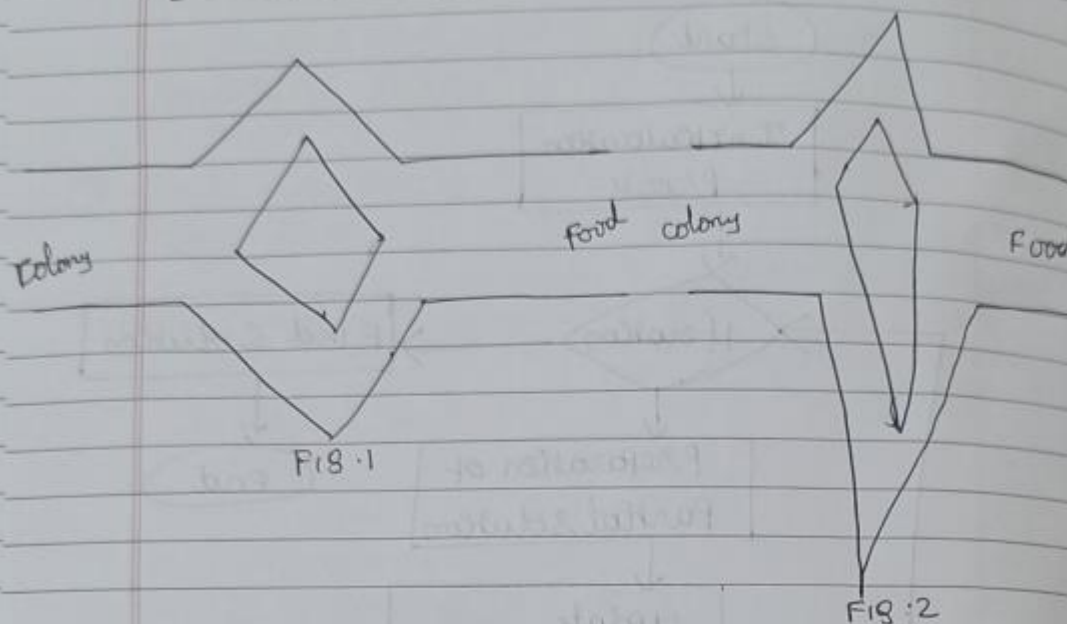
## Foraging behaviour of ACO:



here PATH A & PATH B are equal distance so what ever path we follow it takes a same time to reach a destination.



## double bridge.



in the above double bridge FIG. 1 is the shortest path to reach the destination so fig 1 is the best shortest path.

## Applications of ACO:

Traveling sales man problem.

Vehicle routing

Graph coloring

Job shop scheduling problem

## ACodesign:

Let  $G$  be a graph  $G = (V, E)$

$V \rightarrow$  vertex  $E \rightarrow$  Edges

$E_1, E_2$  - two edges,  $L_1, L_2$  - Length

$R_1, R_2$  - pheromone values.

$V_s$  - source vertex (Ant colony)

$V_d$  - Destination vertex (Food).

Probability for selecting a path  $P_i = \frac{R_i}{R_1 + R_2}$ ;  $i = 1, 2$

## **6. PREPROCESSING AND STRUCTURING OF DATA**

Pre-processing and structuring of data are critical steps in data analysis and machine learning. They involve cleaning, transforming, and organizing raw data into a suitable format for analysis and modeling. Here are common techniques used for data pre-processing and structuring:

### **\*\*1. Data Cleaning:\*\***

- **\*\*Handling Missing Values:\*\*** Deal with missing data by either removing rows or columns with missing values, imputing missing values using statistical methods, or using machine learning algorithms to predict and fill in missing values.
- **\*\*Outlier Detection and Treatment:\*\*** Identify and handle outliers by using statistical methods or machine learning models to mitigate their impact on analysis.
- **\*\*Noise Reduction:\*\*** Remove or reduce noise in the data, such as errors or inconsistencies, by using techniques like smoothing or filtering.

### **\*\*2. Data Transformation:\*\***

- **\*\*Normalization:\*\*** Scale numeric features to have a common range (e.g., between 0 and 1) to ensure that no variable dominates others due to its scale.
- **\*\*Standardization:\*\*** Standardize data by subtracting the mean and dividing by the standard deviation to give the data a mean of 0 and a standard deviation of 1.
- **\*\*Encoding Categorical Variables:\*\*** Convert categorical variables into numeric form using techniques like one-hot encoding, label encoding, or binary encoding.
- **\*\*Feature Engineering:\*\*** Create new features or modify existing ones to capture relevant information and improve model performance.
- **\*\*Aggregation:\*\*** Aggregate data at different levels (e.g., grouping sales data by month or region) to create summary statistics.

### **\*\*3. Data Reduction:\*\***

- **\*\*Dimensionality Reduction:\*\*** Use techniques like Principal Component Analysis (PCA) or feature selection to reduce the number of features while retaining the most relevant information.
- **\*\*Sampling:\*\*** If working with large datasets, use techniques like random sampling to select a representative subset of the data for analysis.

### **\*\*4. Data Integration:\*\***

- **\*\*Combine Data Sources:\*\*** Merge or join data from multiple sources or databases to create a unified dataset for analysis.



- **Data Alignment:** Ensure that data from different sources are aligned correctly, especially when dealing with time-series data or data with different time granularities.

#### **5. Data Resampling:**

- **Temporal Resampling:** Convert time-series data to different time frequencies (e.g., daily to monthly) to match the analysis requirements.

- **Upsampling and Downsampling:** Adjust the sampling rate of data to capture more or less detail, depending on the analysis needs.

#### **6. Data Imbalance Handling:**

- **Over-sampling:** Increase the number of instances in the minority class to balance class distributions in classification problems.

- **Under-sampling:** Reduce the number of instances in the majority class to balance class distributions.

#### **7. Data Formatting and Structuring:**

- **Date and Time Parsing:** Convert date and time fields into standardized formats for consistent analysis.

- **Text Preprocessing:** Tokenize, remove stop words, and perform other text processing tasks when dealing with textual data for natural language processing (NLP).

- **Data Binning and Discretization:** Convert continuous data into discrete intervals or bins for easier analysis.

#### **8. Data Validation and Quality Checks:**

- **Data Validation:** Verify that data adheres to expected constraints, such as data types, ranges, and relationships between variables.

- **Data Profiling:** Profile the data to identify data quality issues, anomalies, and patterns.

#### **9. Data Splitting:**

- **Train-Test Split:** Divide the dataset into training and testing subsets to evaluate model performance.

- **Cross-Validation:** Use k-fold cross-validation to assess model generalization and reduce the risk of overfitting.

#### **10. Data Storage and Management:**

**- Choose appropriate data storage formats and databases to ensure efficient data retrieval and management.**

**Effective data pre-processing and structuring are essential for building accurate and reliable models and gaining meaningful insights from data. These techniques help enhance data quality, reduce noise, and improve the suitability of data for analysis, leading to more robust and valuable outcomes in data-driven projects.**