# studocu

# BDA-Module 4,5,6 (Z魔王)

Bachelor of Engineering (University of Mumbai)

# BDA

## MODULE 4 Mining Data Streams

## Q1. List the different issues and challenges in data stream query processing. (5)

**Data stream query processing is the process of processing and analyzing continuous data streams in real time. It is a challenging task due to the following issues and challenges:**

- High volume: Data streams can be very large, and the volume can vary over time. This can make it difficult to process the data in real time.

- **Velocity:** Data streams are continuous, and the data arrives at a high rate. This can make it difficult to keep up with the data and process it in real time.

- **Unboundedness:** Data streams are potentially unbounded, meaning that they can continue to grow indefinitely. This can make it difficult to store the data and process it in real time.

- **Inconsistency:** Data streams can be inconsistent, meaning that the data may be incomplete, inaccurate, or out of order. This can make it difficult to process the data and get accurate results.

- **Noise:** Data streams may contain noise, which is data that is irrelevant or inaccurate. This can make it difficult to identify patterns and insights in the data.

- **High data velocity and volume:** To address the high velocity and volume of data streams, data stream processing systems typically use distributed processing architectures and scalable algorithms. Distributed processing allows the system to process data in parallel on multiple nodes, which can improve performance and scalability. Scalable algorithms are designed to handle large volumes of data without sacrificing performance or accuracy.

- **Data quality**: Data cleaning techniques can be used to remove noise and inconsistencies from data streams. However, data cleaning can be challenging and computationally expensive, especially for high-velocity data streams. One approach to address this challenge is to use online

data cleaning techniques, which can clean data on the fly as it is processed.

- **Latency:** To achieve low-latency processing, data stream processing systems typically use in-memory processing techniques. In-memory processing allows the system to access data quickly without having to read it from disk. Another approach to reducing latency is to use incremental computation techniques, which can update the query results as new data arrives.

- **Scalability:** To scale to handle growing data streams, data stream processing systems typically use a variety of techniques, such as load balancing and distributed processing. Load balancing ensures that the processing workload is evenly distributed across multiple nodes. Distributed processing allows the system to process data in parallel on multiple nodes, which can improve performance and scalability.

- **Complex event processing:** Complex event processing (CEP) is a challenging task, especially for high-velocity data streams. One approach to address this challenge is to use rule-based CEP engines, which can be used to define rules for identifying complex events in data streams. Another approach is to use machine learning algorithms to identify complex events in data streams.

- **Out-of-order data:** To handle out-of-order data, data stream processing systems typically use buffering and windowing techniques. Buffering allows the system to store data temporarily until it is in order. Windowing allows the system to process data in chunks, or windows, which can simplify the processing logic.

- **Resource constraints:** To efficiently use resources on resource-constrained devices, data stream processing systems typically use specialized algorithms and data structures. One approach is to use approximate query processing techniques, which can produce approximate results with less resource usage. Another approach is to use incremental computation techniques, which can update the query results as new data arrives.

- **Security and privacy:** To protect sensitive data in data streams, data stream processing systems typically use encryption and access control techniques. Encryption can be used to protect data from unauthorized

access. Access control techniques can be used to restrict who can access and process the data.

- **Dynamic schema:** To handle dynamic schemas, data stream processing systems typically use schema-agnostic algorithms and data structures. Schema-agnostic algorithms can process data without requiring a pre-defined schema. Schema-agnostic data structures can store data in a flexible and extensible manner.

- **Fault tolerance:** To achieve fault tolerance, data stream processing systems typically use replication and checkpointing techniques. Replication can be used to replicate data across multiple nodes so that the system can continue to operate even if one node fails. Checkpointing can be used to periodically save the state of the system so that it can be restored from a checkpoint in the event of a failure.

# Q2. Clearly explain the concept of a Bloom Filter for stream data mining with the help of an example. 10

https://youtu.be/g9-NFvd01Hg?si=TVYtIYsvqAXA4EuO

A Bloom filter is a probabilistic data structure that can be used to test whether an element is a member of a set. It is a space-efficient way to store information about a large set, but it is also prone to false positives.

Bloom filters are particularly well-suited for stream data mining because they can be updated efficiently as new data arrives. This makes them ideal for applications such as anomaly detection and fraud prevention.

**Here is an example of how a Bloom filter can be used for stream data mining:**

- Suppose we have a stream of data that represents the IP addresses of visitors to a website. We want to use a Bloom filter to detect anomalous IP addresses.

- We start by initializing the Bloom filter with a fixed number of bits. Each bit represents a potential IP address.

- As each IP address arrives in the stream, we hash it to a bit in the Bloom filter and set that bit to 1.

- To test whether a particular IP address is anomalous, we hash it to the same bit in the Bloom filter and check if the bit is set to 1. If it is, the IP address is likely to be anomalous.

- However, it is important to note that a Bloom filter is prone to false positives. This means that the Bloom filter may indicate that an IP address is anomalous even if it is not. This is because it is possible for multiple IP addresses to hash to the same bit in the Bloom filter.

To reduce the number of false positives, we can increase the size of the Bloom filter. This will increase the accuracy of the Bloom filter, but it will also increase the amount of memory required.

**Bloom filters can be used for a variety of stream data mining applications, including:**

- **Anomaly detection:** Bloom filters can be used to detect anomalous data points in a stream of data.

- **Fraud prevention:** Bloom filters can be used to detect fraudulent transactions in a stream of financial transactions.

- **Network monitoring:** Bloom filters can be used to detect malicious traffic in a network.

- **Content filtering:** Bloom filters can be used to filter out unwanted content from a stream of data.

Bloom filters are a powerful tool for stream data mining, but they are important to understand the limitations of Bloom filters before using them in production applications.

# Q3. Using an example bit stream explains the working of the DGIM algorithms to count the number of 1's(ones) in a data stream. Give the updating buckets approach of DGIM algorithm. 10M

https://youtu.be/bipfcg2U5XI?si=n7mYVMW4Jr5ffWsU

https://youtu.be/8YFvrkhlzyA?si=NHHvY4E2mc52LSJA

The DGIM algorithm works by maintaining a set of buckets. Each bucket represents a range of bit positions in the data stream. The number of 1's in a bucket is estimated by counting the number of times the bucket has been updated.

To update a bucket, the DGIM algorithm first checks to see if the bit at the current position in the data stream is 1. If it is, the algorithm increments the count for the bucket. Otherwise, the algorithm does nothing.

**The following is an example of how the DGIM algorithm works to count the number of 1's in a data stream:**

(**Instead of below step refer video links to better understand**)

**Step 1: Initialize Buckets**

- Start with three buckets: B0, B1, and B2.

- B0 represents a time frame of 1 (the last bit).

- B1 represents a time frame of 2 (bits 2 and 1 before the last bit).

- B2 represents a time frame of 4 (bits 4, 3, 2, and 1 before the last bit).

**Step 2: Read the Stream**

- Read the bits of the data stream one at a time.

- If a bit is 1, update B0 (the 1-bit bucket) and increment its timer by 1. If B0's timer is more than the window size (let's say 16), then reduce it back to 1.

- If a bit is 0, there's no update for B1 or B2.

**Step 3: Move the Window**

- As you continue reading bits, some bits will move out of the window (e.g., older bits that are outside the last 16 bits).

- Ignore the bits outside the window; only focus on the most recent 16 bits.

**Step 4: Add New Bits**

- When new bits come into the window, update the buckets accordingly, just like in Step 2.

- If a new bit is 1, update B0 and potentially introduce a new bucket B3 representing a time frame of 8.

**Step 5: Counting Ones**

- To estimate the number of 1s within the 16-bit window, add the counts in buckets B0, B1, and B2.

- If B0's timer has exceeded the window size, reduce it to 1 for accuracy.

- For even more accuracy, you can also consider the timing information of each bucket and estimate the count based on that.

# Q4. Distinguish the following: DBMS and DSMS.

| Feature | DBMS | DSMS | Examples |
|---|---|---|---|
| Data type | Structured | Unstructured, structured | MySQL, PostgreSQL, Oracle Database |
| Data access | Random | Sequential | Indexes allow for efficient random access to specific data points in a DBMS. |
| Query type | One-time | Continuous | DBMSs are typically used to run one-time queries against stored data. |
| Storage capacity | Unlimited | Limited | DBMSs can typically store unlimited amounts of data on disk. |
| Rate of change | Low | High | Data in a DBMS is typically changed less frequently than data in a data stream. |
| Data processing | Batch | Real-time | DBMSs typically process data in batches, which can be inefficient for real-time applications. |
| Fault tolerance | High | Low | DBMSs typically have high fault tolerance, meaning that they can continue to operate even if there is a hardware failure. |

# BDA

## MODULE 5 Finding Similar Items and Clustering

# Q1. Write a short note on CURE algorithm. (05)

https://youtu.be/FEAznKQFzSY?si=COK0O-Mtph3GyH_b

https://youtu.be/QDBeUo3k0Nc?si=GNTCAVwnetC6taqv

CURE (Clustering Using REpresentatives) is a hierarchical clustering algorithm that is designed to be efficient and scalable for large datasets. It is based on the idea of using a set of representative points to represent each cluster, rather than using a single centroid point. This makes CURE more robust to outliers and able to identify clusters with non-spherical shapes and size variances.

**The CURE algorithm works in the following steps:**

1. **Subsampling:** The dataset is subsampled to create a smaller, more manageable dataset. This helps to improve the efficiency of the algorithm, especially for very large datasets.

2. **Clustering:** The subsampled dataset is clustered using a standard clustering algorithm, such as k-means clustering or hierarchical clustering.

3. **Shrinking:** The representative points for each cluster are shrunk towards the centroid of the cluster by a fraction. This helps to make the clusters more robust to outliers and noise.

4. **Merging:** The representative points for the different clusters are merged together to form a smaller set of representative points. This step is repeated until the desired number of clusters is reached.

**The CURE algorithm has a number of advantages over other clustering algorithms, including:**

- It is efficient and scalable for large datasets.

- It is robust to outliers and noise.

- It is able to identify clusters with non-spherical shapes and size variances.

CURE is a popular clustering algorithm in big data analytics. It is used in a variety of applications, including customer segmentation, image analysis, and fraud detection.

**Example of CURE algorithm in BDA**

Suppose we have a large dataset of customer data, and we want to cluster the customers into two segments: high-value customers and low-value customers. We can use the CURE algorithm to cluster the customers based on their purchase history.

First, we would need to subsample the dataset to create a smaller, more manageable dataset. Then, we would use a standard clustering algorithm, such as k-means clustering, to cluster the subsampled dataset into two groups. Next, we would shrink the representative points for each cluster towards the centroid of the cluster. Finally, we would merge the representative points for the two clusters together to form a smaller set of representative points.

Once the clustering is complete, we can analyse the two clusters to identify the high-value customers and the low-value customers. For example, we might find that the high-value customers are the customers who have made the most purchases in the past year. This information can then be used to develop targeted marketing campaigns for each customer segment.

**Q2. Find the Jaccard distance and cosine distance between the following pair of set X= (0,1,2,4,5,3) and Y= (5,6,7,9,10,8).**

3) Find the Jaccard distance & Cosine distance between the following pairs of set :

$\quad X = (0,1,2,4,5,3)$ and $Y = (5,6,7,9,10,8)$

$X = (0,1,2,4,5,3)$
$Y = (5,6,7,9,10,8)$

① Jaccard Distance :-

$$J(X,Y) = \frac{|A \cap B|}{|A \cup B|} = \frac{|A \cap B|}{|A| + |B| - |A \cap B|}$$

$$= \frac{1}{6+6-1}$$

$$= \frac{1}{12-1}$$

$$= \frac{1}{11} = 0.09$$

∴ Jaccard distance $= 1 - STM (X,Y)$

$\qquad = 1 - 0.09$

$\qquad = 0.91$

② Cosine Distance :-

Dot Product $= (0 \times 5) + (1 \times 6) + (2 \times 7) + (4 \times 9) +$
$\qquad\qquad (8 \times 10) + (3 \times 8)$

$\qquad = 130$

$L^2$ Norm for $x = \sqrt{(0)^2 + (1)^2 + (2)^2 + (4)^2 + (5)^2 + (3)^2}$

$\qquad = \sqrt{55}$

$L^2$ Norms for $y = \sqrt{(5)^2 + (6)^2 + (7)^2 + (9)^2 + (10)^2}$

$$= \sqrt{355}$$

$$\cos \theta = \frac{130}{\sqrt{55} \times \sqrt{355}} = 0.93$$

$$\cos \theta = 0.93$$
$$\theta = \cos^{-1}(0.43)$$
$$\theta = 21.56°$$

https://youtu.be/gwtTM9M6ugU?si=XEnzieScgjtfUS0O

https://youtu.be/PdiHBII-5KQ?si=yaEVzPX2aBc3CUZE

# Q3. write a note on different distance measures that can be used to find similarity / dissimilarity between data points in a big data set. 10M

1. **Euclidean distance:** The Euclidean distance between two points is the length of the line segment connecting them. It is the most common distance measure used in machine learning.
   https://youtu.be/yM2JbjM2-kU?si=w14Nbd8vKkyThX4Z

$$\text{Distance} = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

2. **Manhattan distance:** The Manhattan distance between two points is the sum of the absolute differences between their coordinates. It is also known as the L1 distance.

$$\text{Distance} = |x_2 - x_1| + |y_2 - y_1|$$

3. **Minkowski distance:** The Minkowski distance is a generalization of the Euclidean distance and Manhattan distance. It is defined as the p-th root of the sum of the p-th powers of the absolute differences between the coordinates of the two points.
   https://youtu.be/PdiHBIl-5KQ?si=yaEVzPX2aBc3CUZE

$$\text{Distance} = \left( |x_2 - x_1|^p + |y_2 - y_1|^p \right)^{\frac{1}{p}}$$

4. **Cosine similarity:** The cosine similarity between two vectors is the cosine of the angle between them. It is a measure of how similar the direction of the two vectors is.

$$\text{Similarity} = \frac{A \cdot B}{\|A\| \|B\|}$$

5. **Jaccard similarity:** The Jaccard similarity between two sets is the number of elements in the intersection of the two sets divided by the number of elements in the union of the two sets. It is a measure of how similar the two sets are in terms of their elements.

$$\text{Similarity} = \frac{|A \cap B|}{|A \cup B|}$$

6. **Hamming distance:** The Hamming distance between two binary strings is the number of bits that differ between the two strings. It is a measure of how similar the two strings are in terms of their bit patterns.

$$H(A, B) = \sum_{i=1}^{n}(A_i \neq B_i)$$

7. **Chebyshev distance:** The Chebyshev distance between two points is the maximum absolute difference between their coordinates. It is also known as the L∞ distance.

$$D(P_1, P_2) = \max(|x_2 - x_1|, |y_2 - y_1|)$$

8. **Mahalanobis distance:** The Mahalanobis distance between two points is a measure of the distance between the two points in terms of the standard deviation of the data. It is useful for data that is not normally distributed.

$$D(\mathbf{x}, \boldsymbol{\mu}) = \sqrt{(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu})}$$

9. **Spearman's rank correlation coefficient:** Spearman's rank correlation coefficient is a measure of the monotonic relationship between two variables. It is useful for ordinal data, which is data that is ranked in order but does not have equal intervals between the ranks.

$$\rho = 1 - \frac{6 \sum d^2}{n(n^2 - 1)}$$

Where:

- $\rho$ is Spearman's rank correlation coefficient.
- $d$ represents the difference between the ranks of corresponding pairs of data points.
- $n$ is the number of data points.

10. **Kendall's tau:** Kendall's tau is another measure of the monotonic relationship between two variables. It is similar to Spearman's rank correlation coefficient, but it is more robust to outliers

$$\tau = \frac{2C}{n(n-1)}$$

Where:

- $\tau$ represents Kendall's tau.
- $C$ is the number of concordant pairs.
- $n$ is the total number of data points.

# Q4. Define Edit distance. Explain with examples.

https://youtu.be/9t5RGyFht5U?si=GCrJeQ9rvkCq2Yei

Edit distance is a string metric that measures the difference between two strings by counting the minimum number of edit operations required to transform one string into the other. The edit operations can be insertions, deletions, or substitutions of characters.

**Edit distance is a useful metric for many big data applications, such as:**

- Natural language processing: Edit distance can be used to correct spelling errors, identify synonyms, and compare text documents.

- Information retrieval: Edit distance can be used to find similar search results, identify duplicate content, and detect plagiarism.

- Machine learning: Edit distance can be used to train machine learning models to perform tasks such as text classification, sequence prediction, and anomaly detection.

**Here are some examples of edit distance:**

- The edit distance between the strings "cat" and "dog" is 2. This is because two edits are required to transform one string into the other: a deletion of the letter "c" and an insertion of the letter "d".

- The edit distance between the strings "algorithm" and "altruism" is 1. This is because only one edit is required to transform one string into the other: a substitution of the letter "g" for the letter "r".

- The edit distance between the strings "big data" and "big data analytics" is 2. This is because two edits are required to transform one string into the other: an insertion of the word "analytics" and a space.

**In big data analytics, edit distance can be used to solve a variety of problems. For example, edit distance can be used to:**

- Cluster text documents based on their similarity.

- Detect fraudulent transactions by comparing them to known fraudulent transactions.

- Identify duplicate records in a database.

- Correct spelling errors in large datasets.

- Find similar products in a catalog for recommendation systems.

Edit distance is a powerful and versatile metric that can be used to solve a variety of big data problems.

**Here is an example of how edit distance can be used in big data analytics:**

A company has a large dataset of customer transactions. The company wants to identify fraudulent transactions. One way to do this is to use edit distance to compare each transaction to a database of known fraudulent transactions. If the edit distance between a transaction and a known fraudulent transaction is below a certain threshold, then the transaction is flagged as potentially fraudulent.
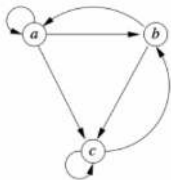
**Another example is:**

A company has a large dataset of text documents. The company wants to cluster the documents based on their similarity. One way to do this is to use edit distance to calculate the similarity between each pair of documents. The documents are then clustered based on their similarity scores.

# BDA

## MODULE 6 Real Time Big Data Models

# Q1. Compute simplified page rank using damping factor d = 0.8 for web for 3 iterations. (10)

**(See video to understand the problem)**

# PageRank Algorithm

## Example

Formula, $r_{i+1} = M \times r_i$

Iteration 1

$$
\begin{array}{c}
M \\
\begin{array}{ccc} A & B & C \end{array} \\
\begin{array}{c} A \\ B \\ C \end{array}
\begin{pmatrix}
1/3 & 1/2 & 0 \\
1/3 & 0 & 1/2 \\
1/3 & 1/2 & 1/2
\end{pmatrix}
\end{array}
\times
\begin{array}{c}
r_0 \\
\begin{pmatrix} 1/3 \\ 1/3 \\ 1/3 \end{pmatrix}
\end{array}
=
\begin{array}{c}
r_1 \\
\begin{pmatrix} 5/18 \\ 5/18 \\ 8/18 \end{pmatrix}
\end{array}
$$

---

# PageRank Algorithm

## Example

Formula, $r_{i+1} = M \times r_i$

Iteration 2

$$
\begin{array}{c}
M \\
\begin{array}{ccc} A & B & C \end{array} \\
\begin{array}{c} A \\ B \\ C \end{array}
\begin{pmatrix}
1/3 & 1/2 & 0 \\
1/3 & 0 & 1/2 \\
1/3 & 1/2 & 1/2
\end{pmatrix}
\end{array}
\times
\begin{array}{c}
r_1 \\
\begin{pmatrix} 5/18 \\ 5/18 \\ 8/18 \end{pmatrix}
\end{array}
=
\begin{array}{c}
r_2 \\
\begin{pmatrix} 25/108 \\ 34/108 \\ 49/108 \end{pmatrix}
\end{array}
$$

---

# PageRank Algorithm

## Example

Formula, $r_{i+1} = M \times r_i$

Iteration 3

$$
\begin{array}{c}
M \\
\begin{array}{ccc} A & B & C \end{array} \\
\begin{array}{c} A \\ B \\ C \end{array}
\begin{pmatrix}
1/3 & 1/2 & 0 \\
1/3 & 0 & 1/2 \\
1/3 & 1/2 & 1/2
\end{pmatrix}
\end{array}
\times
\begin{array}{c}
r_2 \\
\begin{pmatrix} 25/108 \\ 34/108 \\ 49/108 \end{pmatrix}
\end{array}
=
\begin{array}{c}
r_3 \\
\begin{pmatrix} 152/648 \\ 197/648 \\ 299/648 \end{pmatrix}
\end{array}
$$

## Q2. Explain collaborative filtering systems. How is it different from a content-based system? 10M

https://youtu.be/W30s-EoP3d0?si=fg7hQDEdVXDazWyM

Collaborative filtering (CF) systems are a type of recommender system that uses the preferences of other users to recommend items to a particular user. CF systems work by finding users who have similar tastes to the user in question, and then recommending items that those similar users have liked.

CF systems are based on the idea that people who like the same things are likely to like other similar things. For example, if two users both have rated the movie "The Shawshank Redemption" highly, then the CF system might recommend the movie "The Godfather" to both users, even if neither user has rated that movie yet.

CF systems can be implemented using a variety of different algorithms, but they all work by finding similarity between users. Some common algorithms include:

- **Neighbourhood-based methods**: These methods find similar users to the target user by computing a similarity score between users. The similarity score can be based on a variety of factors, such as the items that the users have rated or purchased in the past. Once the similar users have been found, the system recommends items that those similar users have liked.

- **Matrix factorization methods:** These methods learn a latent representation of users and items. The latent representation of a user captures their preferences for items, and the latent representation of an item captures the characteristics that make it similar to other items. Once the latent representations have been learned, the system can recommend items to users by finding items that are similar to items that the user has liked in the past.

**Here are some examples of collaborative filtering systems:**

- Product recommendation systems on e-commerce websites

- Movie recommendation systems on streaming services

- Music recommendation systems on streaming services

- Book recommendation systems on online retailers

- Social media platforms that recommend new friends to users

**Here are some of the advantages of collaborative filtering systems:**

- They can make more serendipitous recommendations than content-based systems.

- They are able to learn over time and become more accurate at making recommendations.

- They are relatively easy to implement.

**Here are some of the disadvantages of collaborative filtering systems:**

- They require a large amount of user data to be effective.

- They can be susceptible to cold start problems, which is when the system does not have enough data about a new user or item to make accurate recommendations.

- They can be susceptible to spam and fake reviews.

# Collaborative vs content based (see table foe better understanding)

- Content based systems examine the properties of the item examined or recommended.

- These systems take input from the user profile and the contextual parameters along with product features to make the recommendation list.

- Similarity of items is determined by measuring the similarity in their properties.

- These systems need some information related to the content of available items, but no complete information.

- It also requires user profiling describing what user likes.

- In movie recommendation system- content based filtering systems will recommend the movie based on users profile information like age, gender, etc as well as properties of movie like genre, actors, etc.

| Feature | Collaborative filtering (CF) | Content-based |
|---|---|---|
| Recommendation strategy | Recommends items based on the preferences of other users. | Recommends items based on the features of items that the user has liked in the past. |
| Ability to make serendipitous recommendations | Yes, CF systems can recommend items to users that are outside of their usual interests. | No, content-based systems are more likely to recommend items that are similar to items that the user has liked in the past. |
| Ability to learn over time | Yes, CF systems can learn over time to better understand user preferences and make more accurate recommendations. | Yes, content-based systems can also learn over time, but to a lesser extent than CF systems. |
| Data requirements | Requires a large amount of user data to be effective. | Requires less user data to be effective than CF systems. |
| Cold start problem | More susceptible to the cold start problem than content-based systems. | Less susceptible to the cold start problem than CF systems. |

# Q3. Explain model for recommendation system in detail.10

A recommendation system is a software application that provides suggestions for items to users based on their past behaviours and preferences. Recommendation systems are used in a wide variety of applications, including product recommendation systems on e-commerce websites, movie recommendation systems on streaming services, music recommendation systems on streaming services, and book recommendation systems on online retailers.

Recommendation systems work by collecting data about users and items. This data can include ratings, purchases, views, and other interactions. The recommendation system then uses this data to learn about the user's preferences and to generate recommendations for items that the user is likely to like.

**There are many different models that can be used for recommendation systems. Some of the most popular models include:**

- **Collaborative filtering (CF) models:** CF models use the preferences of other users to recommend items to a particular user. CF models work by finding users who have similar tastes to the user in question, and then recommending items that those similar users have liked. CF models are often divided into two categories: neighborhood-based methods and matrix factorization methods.

- **Content-based models:** Content-based models recommend items to a user based on the features of items that the user has liked in the past. For example, a content-based model that recommends movies might recommend movies with the same genre, actors, or director as movies that the user has rated highly in the past.

### 1.Collaborative filtering (CF) models

CF models are the most popular type of model for recommendation systems. CF models work by finding users who have similar tastes to the user in question, and then recommending items that those similar users have liked.

**There are two main types of CF models:**

- **Neighborhood-based methods:** Neighborhood-based methods find similar users to the target user by computing a similarity score between users. The similarity score can be based on a variety of factors, such as the items that the users have rated or purchased in the past. Once the similar users have been found, the system recommends items that those similar users have liked.

- **Matrix factorization methods**: Matrix factorization methods learn a latent representation of users and items. The latent representation of a user captures their preferences for items, and the latent representation of an item captures the characteristics that make it similar to other items. Once the latent representations have been learned, the system can recommend items to users by finding items that are similar to items that the user has liked in the past.

CF models are used in a wide variety of applications, including product recommendation systems on e-commerce websites, movie recommendation systems on streaming services, and music recommendation systems on streaming services.

## 2.Content-based models

Content-based models recommend items to a user based on the features of items that the user has liked in the past. For example, a content-based model that recommends movies might recommend movies with the same genre, actors, or director as movies that the user has rated highly in the past.

Content-based models are simpler to implement than CF models, but they are not as effective at making serendipitous recommendations. Serendipitous recommendations are recommendations for items that are outside of the user's usual interests, but that the user might still like.

Content-based models are often used in conjunction with CF models to create hybrid recommender systems. Hybrid recommender systems can combine the advantages of both CF and content-based models to provide more accurate and relevant recommendations.

**Evaluation of recommendation system models**

**There are a number of different metrics that can be used to evaluate recommendation system models. Some of the most common metrics include:**

- **Accuracy:** Accuracy measures how well the model predicts the user's ratings for items.

- **Recall:** Recall measures how well the model recommends items that the user likes.

- **Precision:** Precision measures how well the model recommends items that the user actually clicks on or purchases.

- **Diversity:** Diversity measures how well the model recommends items that are different from each other.

- **Novelty:** Novelty measures how well the model recommends items that the user has not seen before.

**Examples of recommendation systems:**

**Netflix:** Recommends movies and TV shows based on viewing history and ratings.

**Amazon:** Recommends products based on purchase history and browsing behavior.

**Spotify:** Recommends songs and albums based on listening history and preferences.

**YouTube:** Recommends videos based on viewing history and watchlist.