# Building a Netflix recommender system leveraging Kmeans Minibatch

Problem Statement

Data Description

Artifact Submission

Project Evaluation Criteria

Recommendation

Vibelines & Bounties

# 1. Understanding Content-Based Filtering - The Basics

**Content-Based Filtering** is a type of recommendation system that suggests items to users based on the characteristics (or "content") of items the user has previously liked or interacted with. The core idea is to build a profile of the user's preferences by analyzing the attributes of items they have consumed or rated highly in the past.
Here's how it generally works:

**1. Item Representation:** Each item (e.g., a movie, a TV show, an article) is described by a set of attributes or features (e.g., for a movie: its genre, director, cast, description, keywords).

**2. User Profile Creation:** A user's profile is built based on the features of items they have expressed interest in (e.g., movies they've watched, rated highly, or added to their watchlist). This profile often represents the user's "taste" or "preferences" in terms of item characteristics.

**3. Recommendation Generation:** The system then compares the user's profile to the features of unrated or unconsumed items. Items that are most similar to the user's profile are recommended.
The key principle is: "If you liked this movie, you'll like other movies that are similar to it in terms of their content attributes."

# 2. Associated Concepts in Mini-Batch K-Means

K-Means is a popular clustering algorithm that aims to partition 'n' observations into 'k' clusters, where each observation belongs to the cluster with the nearest mean (centroid). The standard K-Means algorithm processes the entire dataset in each iteration to update the cluster centroids. This can be computationally expensive and slow for very large datasets.

Mini-Batch K-Means is a variation of the K-Means algorithm designed to address this scalability issue. Instead of using the full dataset, Mini-Batch K-Means performs updates based on small, randomly sampled subsets of the data called "mini-batches".
Here's how it works:

- Initialize Centroids: Just like standard K-Means, it starts by randomly initializing 'k' cluster centroids.
- Iterative Updates with Mini-Batches:
- In each iteration, a small, random subset of data (a mini-batch) is drawn from the dataset.
- For each point in the mini-batch, its closest centroid is identified. The chosen centroids are then updated incrementally using a weighted average of the new data points assigned to them. This update is often a learning rate-like adjustment, rather than a full re-calculation based on all points assigned to that centroid in the entire dataset.
- Convergence: The process continues for a fixed number of iterations or until the centroids stabilize.

Key Advantages of Mini-Batch K-Means:
- Speed for Large Datasets: Significantly faster than standard K-Means for very large datasets because it doesn't need to load and process the entire dataset in memory for every update.
- Memory Efficiency: Requires less memory as it only operates on small mini-batches at a time.
- Approximation: While it's an approximation of K-Means, its results are often very close to those of the full K-Batch algorithm, especially with a sufficient number of iterations and well-chosen mini-batch sizes.
- Using Mini-Batch K-Means in Content-Based Filtering for Netflix Content-based filtering recommends items to users based on the attributes (content) of items they have liked in the past. For Netflix movies and shows, these attributes can include genres, directors, cast, keywords from descriptions, and even plot summaries.

# 3. Why is Content-Based Filtering Important?

- **Interpretability:** Recommendations are easily explainable because they are based on explicit item attributes (e.g., "We recommend this movie because it's a sci-fi thriller with a strong female lead, just like others you've enjoyed").

- **No Cold Start for New Items:** New items can be recommended as soon as their attributes are known, even if no one has interacted with them yet. This is crucial for platforms constantly adding new content.

- **User Independence:** Recommendations for one user are not affected by the preferences of other users, which can be useful for niche tastes.

- **Handles Niche Tastes:** Can recommend items that appeal to very specific user preferences, even if those preferences are not shared by many other users.

- **Directly Leverages Item Data:** Makes full use of the rich descriptive information available for items, which can be very detailed for digital content.

# 4. Industries where Content-Based Filtering is particularly useful:

- **Media & Entertainment (Core Application):** Recommending movies/TV shows based on genre, actors, director, plot keywords, and *semantic understanding of descriptions*; music based on artist, genre, mood, instruments, and audio features.

- **E-commerce (especially for products with rich textual descriptions):** Recommending clothing based on style/material descriptions, electronics based on specifications, or books based on plot summaries.

- **News & Content Platforms:** Suggesting articles or blog posts based on topics, keywords, authors, and the *semantic content* of articles a user has read before.

- **Job Boards:** Recommending job postings based on skills, industry, experience, and the *semantic meaning* of job descriptions and user resumes.

- **Research & Academia:** Recommending scientific papers based on keywords, authors, citations, and the *semantic content* of abstracts and full papers.

- **Online Learning Platforms:** Suggesting courses or learning modules based on subjects a student has excelled in or expressed interest in, using *semantic understanding* of course descriptions.

# Data Description

This project focuses on building a **Netflix Recommendation Engine** using the **Content-Based Filtering** approach, specifically leveraging the **Bag-of-Words (BOW)** model for item representation. The objective is to recommend movies and TV shows to users based on the textual content (genres, director, cast, description) of titles they have previously watched or liked.

**About the Dataset:**

The dataset provided is a collection of Netflix titles, containing various metadata crucial for content-based recommendations.

| Column | Description |
|---|---|
| show_id | Unique identifier for each show. |
| type | Type of content (Movie or TV Show). |
| title | Title of the show. |
| director | Director(s) of the show. |
| cast | Main actors/actresses in the show. |
| country | Country of production. |
| date_added | Date the show was added to Netflix. |
| release_year | Original release year of the show. |
| rating | TV rating (e.g., TV-MA, PG-13). |
| duration | Duration of the movie or number of seasons for a TV show. |
| listed_in | Genres/categories the show is listed under. |
| description | A brief synopsis of the show. |

# Artifact Submission

Your submission must include the following five artifacts, all packaged within a single GitHub repository.

**1. Jupyter Notebook (.ipynb)** This is the core of your submission. Your Jupyter Notebook should be a complete, well-documented narrative of your data analysis journey. It must include:
- **Detailed Explanations:** Use Markdown cells to explain your thought process, the questions you are trying to answer, and the insights you've uncovered.
- **Clean Code:** The code should be well-structured, easy to read, and free of unnecessary clutter.
- **Comprehensive Comments:** Use comments to explain complex logic and the purpose of different code blocks.
- **Key Visualizations:** All visualizations should be clear, properly labeled, and directly support your findings.

**2. Presentation (.pptx or .pdf)**
Create a compelling presentation that summarizes your team's analysis and key findings. This presentation should serve as your final pitch. It must include:
- **Executive Summary:** A concise overview of your findings.
- **Key Insights:** The most important takeaways from your analysis.
- **Data-Driven Recommendations:** Actionable steps that can be taken based on your insights.
- **Supporting Visualizations:** A selection of your best visualizations to illustrate your points.

**3. README File (.md)**
The README file is the first thing we'll look at. It should serve as a quick guide to your project and provide essential details. It must include:
- **Project Title :**
- **Brief Problem Statement:** A summary of the project and your approach.
- **Summary of Findings:** A bullet-point summary of your most significant insights.

**4. Attached Dataset**
Please include the original dataset (.csv or other format) within your repository. This ensures the judges can reproduce your analysis without any issues.

**5. GitHub Repository**
Your final submission will be your GitHub repository. The repository name **must follow this exact format:**
**Content_Based_Filtering_ProjectName_TMP**

# Challenge Evaluation Criteria

| Criteria Name | Criteria weight |
|---|---|
| Data Understanding and Exploratory Data Analysis | 20% |
| Data preprocessing and feature engineering | 25% |
| Model building and evaluation | 30% |
| Business Recommendation | 15% |
| Coding guidelines and standards | 10% |

# Recommendation for Content-Based Filtering project with BOW

Here's how Mini-Batch K-Means can be leveraged in such a system:

**1. Feature Extraction and Representation:**

- Movie/Show Content Data: Start with your Netflix dataset containing features like title, genres, director, cast, description, etc..
- Text Preprocessing: Clean and preprocess textual features (like description, genres, director, cast) using techniques like lowercasing, removing stop words, lemmatization, and potentially replacing spaces with underscores for multi-word entities (e.g., "science fiction" -> "science_fiction").
- Vectorization (Embeddings): Convert these textual features into numerical vectors.
- TF-IDF: If your features are simple bags of words (e.g., combined all_features string), you can use TF-IDF to get a numerical representation for each movie/show.
- Word Embeddings (e.g., Word2Vec, FastText): A more sophisticated approach is to use pre-trained (or custom-trained) word embeddings. You can average the embeddings of all keywords/genres/directors in a movie's all_features list to get a single "movie embedding" vector. This captures semantic meaning better.
- Node2Vec (as discussed previously): If you build a graph connecting movies to genres, directors, actors, etc., Node2Vec can generate embeddings for each movie node and each feature node. You'd then use these movie embeddings. Result: Each movie/show in your Netflix catalog is now represented as a dense numerical vector. This is your "dataset" for clustering.

**2. Clustering Movies/Shows with Mini-Batch K-Means:**

Why Cluster? The idea is that movies/shows within the same cluster are "content-similar" to each other. They share common genres, directors, cast members, or thematic elements.

**3. Applying Mini-Batch K-Means:**

You would feed the vectorized movie/show data (e.g., 20-dimensional embeddings from Node2Vec, or 300-dimensional averaged Word2Vec vectors) to the Mini-Batch K-Means algorithm.
You'll need to choose the number of clusters, 'k'. This can be determined using methods like the elbow method or silhouette score.
Mini-Batch K-Means will efficiently group your vast Netflix catalog into 'k' content-based clusters.

**4. Recommendation Generation:**

User Profile: When a user likes a particular movie or show, you identify which cluster that liked item belongs to.

**5. Recommendation Strategy:**

- Intra-Cluster Recommendation: Recommend other movies/shows from the same cluster that the user has not yet seen or rated. This is the most straightforward content-based recommendation.
- Inter-Cluster Exploration: You could also recommend items from neighboring clusters (clusters whose centroids are numerically close in the embedding space) to offer a slightly diversified but still relevant set of recommendations.
- User-Centric Clustering (Optional): If you have user-specific content preference vectors (e.g., average of embeddings of all movies a user has watched), you could also cluster users based on these vectors using Mini-Batch K-Means. Then, for a given user, recommend items from clusters preferred by users in their cluster.

# Benefits of using Mini-Batch K-Means here:

- **Scalability for Netflix Catalog:** Netflix has a massive catalog. Standard K-Means would be too slow and memory-intensive. Mini-Batch K-Means allows you to cluster millions of items efficiently.

- **Fast Model Updates:** If new movies are added frequently, you could potentially update clusters or re-assign new items to existing clusters much faster than with full K-Means.

- **Personalization:** By grouping content, you can easily identify content niches and recommend items that align with a user's demonstrated content preferences, even for new users (cold start) or new items.

- **Complementary to Collaborative Filtering:** Content-based recommendations (from clustering) can be combined with collaborative filtering techniques (which look at user-item interactions) to build a robust hybrid recommender system.

# Lets Go