# BIG DATA ANALYTICS (CS-431)

**Dr. Sriparna Saha**
Associate Professor

**Website**: https://www.iitp.ac.in/~sriparna/
**Google Scholar:** https://scholar.google.co.in/citations?user=Fj7jA_AAAAAJ&hl=en
**Research Lab:** SS_Lab
**Core Research AREA:** NLP, GenAI, LLMs, VLMs, Multimodality, Meta-Learning, Health Care, FinTech, Conversational Agents

**TAs**: Sarmistha Das, Nitish Kumar, Divyanshu Singh, Aditya Bhagat, Harsh Raj

# Apache Mahout: Scalable Machine Learning

**What is Apache Mahout?**

- An open-source project by the Apache Software Foundation.
- Its primary goal is to provide **scalable, distributed machine learning algorithms**.
- The name "Mahout" is derived from the Hindi word for an elephant driver, reflecting its original close ties to Apache Hadoop (the elephant).
- It is not an end-to-end machine learning platform, but rather a **library of implemented algorithms** that can be integrated into applications.
- Focuses on the core ML domains:
  - **Collaborative Filtering** (Recommender Engines)
  - **Clustering** (Unsupervised Learning)
  - **Classification** (Supervised Learning)

# Why Do We Need Mahout? The Challenge of "Big ML"

Traditional machine learning algorithms are often designed to run on a single machine. This approach fails when dealing with "Big Data" due to:

1. **Memory Constraints:** The entire dataset cannot fit into the RAM of a single server.
2. **Processing Power:** The computational complexity of training a model on billions of data points is too much for a single CPU.

Mahout was created to solve this problem by implementing ML algorithms on top of distributed computing paradigms like **Apache Hadoop MapReduce** and, more recently, **Apache Spark**. This allows for horizontal scaling by distributing the computational load across a cluster of machines.

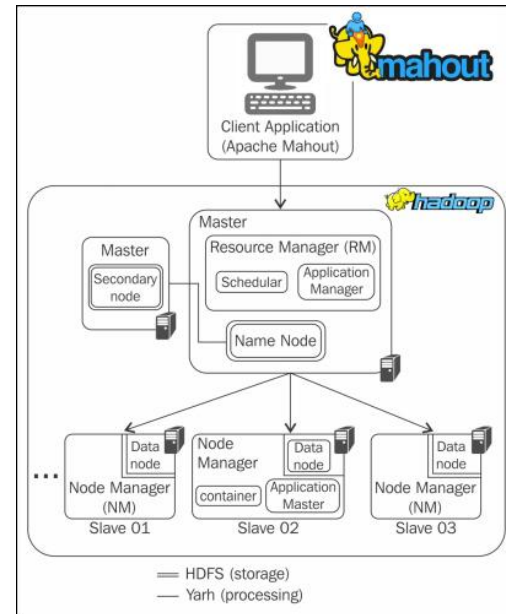# The Evolution of Mahout: From MapReduce to Samsara

Mahout has undergone a significant architectural shift.

- **Phase 1: Hadoop MapReduce (The Original)**
  - Algorithms were implemented as a series of MapReduce jobs.
  - Pros: Extremely scalable and fault-tolerant.
  - Cons: High latency due to repeated disk I/O between MapReduce stages, making it slow for iterative algorithms. Difficult to program and debug.
- **Phase 2: "Samsara" (The Modern Mahout)**
  - A complete rewrite focused on a modern, backend-agnostic architecture.
  - Introduced a Scala-based Domain Specific Language (DSL) that looks similar to the R programming language for statistics.
  - This DSL allows developers to express complex linear algebra operations concisely.

# Mahout's Early Architecture: Deep Integration with Hadoop (Pre-Samsara)

- ○ In its early days, Apache Mahout was fundamentally built on top of **Apache Hadoop**. This architecture leveraged Hadoop's capabilities for both distributed storage and distributed processing.
- ○ **HDFS (Hadoop Distributed File System):** All large datasets used by Mahout algorithms were stored on HDFS. This provided fault-tolerant, high-throughput access to application data.
- ○ **MapReduce:** Mahout algorithms were implemented as a series of MapReduce jobs. Each step of an algorithm (e.g., in K-Means, the assignment and update steps) would often translate into one or more MapReduce jobs.
- ○ **YARN (Yet Another Resource Negotiator):** Hadoop YARN managed the resources (CPU, memory) across the cluster, allocating containers for Mahout's MapReduce tasks to run on the DataNodes.
- ○ **Client Application:** Users would write their Mahout applications (typically in Java) and submit them to the Hadoop cluster, which would then coordinate the execution.
- ○ **Limitations:** While highly scalable and robust, this tight coupling with MapReduce led to performance bottlenecks, especially for iterative machine learning algorithms due to the high I/O overhead of writing intermediate results to HDFS between job stages.



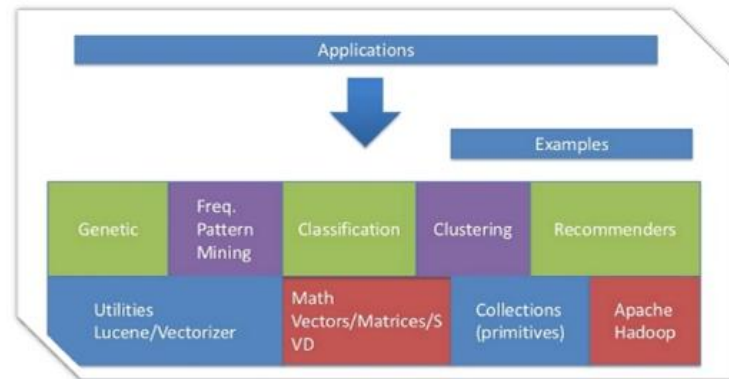Hadoop MapReduce (The Original)

# Mahout's Early Architecture: Deep Integration with Hadoop (Pre-Samsara)

This diagram illustrates Mahout's conceptual, multi-tiered architecture, focusing on the logical separation of concerns.

- **Applications Layer:** At the top, this represents user-developed applications that utilize Mahout's functionalities. It also includes examples provided within the Mahout project.
- **Algorithms Layer:** This is the core engine containing the implemented machine learning algorithms. This includes categories like:
  - **Classification:** For supervised learning tasks.
  - **Clustering:** For unsupervised grouping.
  - **Recommenders:** For building recommendation systems.
  - Other specialized algorithms such as Frequent Pattern Mining or Genetic algorithms.

## General Architecture



**Three-tiers architecture**
(Application, Algorithms and Shared Libraries)

# Mahout's Early Architecture: ==Deep Integration with Hadoop== (Pre-Samsara) Continuation :
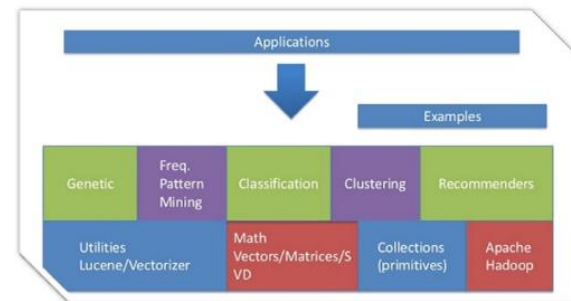
**Core Math Libraries Layer:** This is a crucial foundational layer. It provides highly optimized, distributed implementations of fundamental linear algebra operations and data structures:

- **Vectors and Matrices:** Efficient representation and manipulation of numerical vectors and matrices, which are the building blocks of most ML algorithms.
- **SVD (Singular Value Decomposition):** A powerful matrix factorization technique used in many ML applications, including dimensionality reduction and latent semantic analysis.
- **Collections (Primitives):** Specialized distributed collection types.

**Utilities and Infrastructure Layer:** At the base, this layer provides supporting tools and the link to the underlying distributed computing environment:

- **Utilities (e.g., Lucene, Vectorizer):** Components for data preparation, such as text tokenization (historically Lucene) and converting raw data into numerical vectors.
- **Apache Hadoop:** The foundational distributed processing and storage (MapReduce, HDFS) that historically underpinned Mahout's distributed execution.



General Architecture

**Three-tiers architecture**
(Application, Algorithms and Shared Libraries)

# The Emergence of Samsara: Addressing MapReduce's Limitations

- As the Big Data landscape evolved, the limitations of the MapReduce paradigm for complex, iterative machine learning algorithms became evident.
  - **High Latency:** Each iteration in an algorithm (e.g., refining cluster centroids in K-Means) required writing intermediate results to HDFS and reading them back, incurring significant disk I/O overhead.
  - **Low-Level Abstraction:** Implementing sophisticated ML algorithms directly as sequences of MapReduce jobs was complex, verbose, and difficult to maintain.
  - **Need for In-Memory Processing:** Modern algorithms and larger datasets demanded faster, in-memory processing capabilities.
  - **Mahout's Response: Samsara:** This led to a complete re-architecture of Mahout, known as "Samsara." The goal was to provide a more flexible, efficient, and developer-friendly platform for scalable machine learning.
  - **Key Vision:** Abstract away the execution engine details and focus on providing a powerful, high-level language for expressing machine learning algorithms.
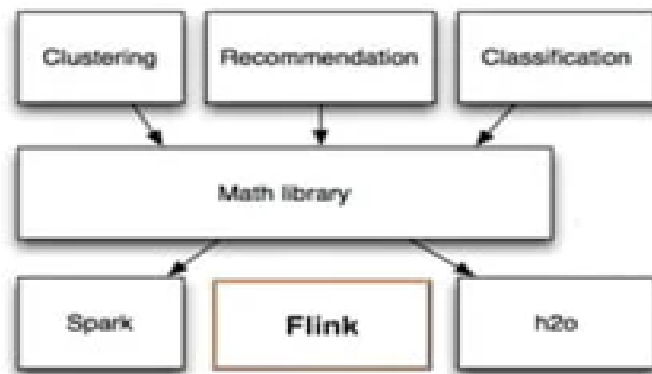
# Samsara's Core Architectural Philosophy: Backend Agnostic

**Algorithms as Algebraic Expressions:** Instead of tightly coupling to a specific execution engine, Mahout Samsara allows users to express machine learning algorithms as high-level **algebraic operations** on distributed matrices and vectors. This is achieved through a Scala-based Domain Specific Language (DSL).

**Central "Math Library" (Algebraic Optimizer):** This is the heart of Samsara. It takes these high-level algebraic expressions and translates them into optimized execution plans for various distributed computing backends. It acts as an abstraction layer between the algorithm definition and its execution.

**Pluggable Execution Backends:** This is the most significant change. Mahout Samsara is no longer solely tied to Hadoop MapReduce. It can execute its optimized plans on different distributed processing frameworks:

- **Apache Spark:** The primary and most mature backend for Mahout Samsara, leveraging Spark's in-memory computation for significant performance gains, especially for iterative algorithms.
- **Apache Flink:** Another powerful distributed processing engine, particularly strong in stream processing, which Mahout can also target.
- **Other Backends:** The architecture is designed to potentially support other engines like H2O or even a single-threaded local backend for development and testing.



The diagram above represents the fundamental shift in Mahout's architecture with Samsara. The core philosophy is **backend-agnosticism**.

# Deep Dive into Samsara: Key Components

- **Scala-based DSL:** Provides a powerful, R-like syntax for expressing linear algebra operations on distributed data. This significantly simplifies algorithm implementation and makes it more accessible to data scientists.
- **Distributed Row Matrix (DRM):** The fundamental data abstraction in Samsara. It's a distributed matrix where rows are partitioned across the cluster. This enables parallel processing of large matrices that cannot fit into a single machine's memory.
- **Algebraic Optimizer:** A sophisticated component that takes the high-level DSL expressions and converts them into optimized execution graphs (e.g., Spark DAGs) tailored for the chosen backend. This ensures efficient execution.
- **Abstracted Backends:** Provides interfaces to connect to different distributed processing engines like Spark, Flink, etc. This makes Mahout flexible and future-proof.

# The Distributed Row Matrix (DRM): Mahout's Foundation

- **Core Data Structure:** The DRM is the cornerstone for representing and processing large matrices in a distributed fashion within Mahout Samsara.
- **Concept:** It's a logical matrix whose rows are physically distributed and stored across the nodes of a cluster.
- **Scalability:** By distributing the rows, Mahout can handle matrices of virtually any size, limited only by the cluster's total storage and memory.
- **Efficient Operations:** The Mahout DSL and backend optimizers are designed to perform linear algebra operations (addition, multiplication, transpose, etc.) directly on DRMs in a highly parallel manner.
- **Integration with Backends:** On Spark, a DRM is typically backed by an RDD (Resilient Distributed Dataset) of `Vector` objects, where each vector represents a row of the matrix.

# Mahout's Scala DSL: An R-like Syntax for Big Data ML

- **Motivation:** Traditional distributed programming (like raw MapReduce) is often verbose and complex, especially for mathematical algorithms.
- **The Solution:** Mahout's DSL offers a high-level, expressive syntax inspired by statistical languages like R.
- **Example (Conceptual):**
  - Instead of writing low-level Spark transformations to compute A×B, you might simply write `A %*% B` where A and B are DRMs.
  - For element-wise operations: `A + 5` or `A * B`.
- **Benefits for Developers:**
  - **Conciseness:** Reduces boilerplate code, allowing developers to focus on the mathematical logic of the algorithm.
  - **Readability:** Code looks more like mathematical notation, making it easier for data scientists and mathematicians to understand.
  - **Productivity:** Speeds up the development and prototyping of new distributed machine learning algorithms.

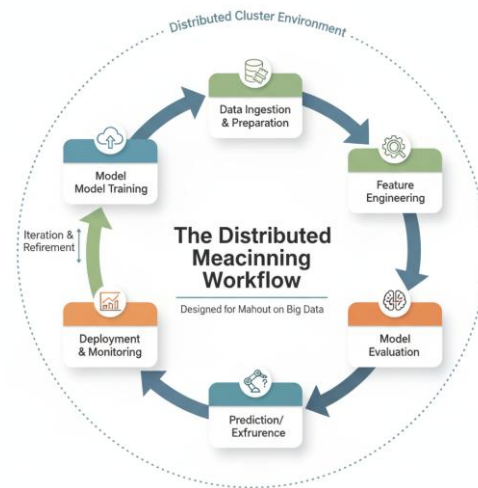# Mahout's Algebraic Optimizer: Bridging DSL and Execution

- **The Translator:** The Algebraic Optimizer is a key innovation in Samsara. It acts as a sophisticated translator between the high-level DSL and the low-level execution engine.
- **Input:** Takes the Mahout DSL expressions (which represent linear algebra operations on DRMs).
- **Optimization:** It analyzes the entire computation graph (DAG - Directed Acyclic Graph) of these operations. It applies various optimizations to reduce data shuffling, minimize I/O, and improve computational efficiency.
- **Output:** Generates an optimized physical execution plan specific to the chosen backend (e.g., a highly optimized Spark DAG).
- **Why it's crucial:** This optimizer ensures that even though developers write simple, high-level code, the underlying execution is highly efficient and performs optimally on the distributed cluster. It effectively hides the complexity of distributed computing from the algorithm developer.

# The Distributed Machine Learning Workflow

Building a machine learning model with a distributed tool like Mahout follows a standard, cyclical workflow. The key difference is that each step is designed to operate on data distributed across a cluster.

1. **Data Ingestion & Preparation:** Loading raw data and transforming it into a usable numerical format (vectors).
2. **Feature Engineering:** Creating meaningful features from the raw data to improve model performance.
3. **Model Training:** Running a distributed algorithm on the prepared data to produce a model. This is the most computationally intensive step.
4. **Model Evaluation:** Assessing the model's performance on unseen data using various metrics.
5. **Prediction/Inference:** Using the trained model to make predictions on new data.
6. **Iteration:** Repeating the process, tuning parameters, and refining features to improve the model.
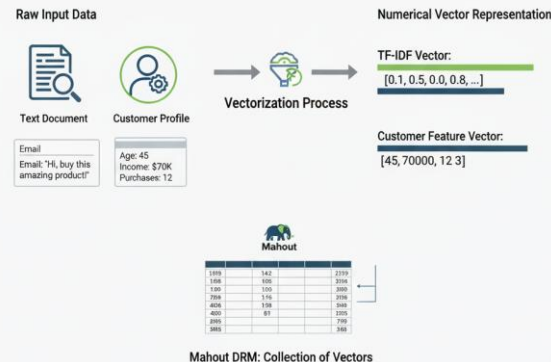
# Core Concept: Data Vectorization

Machine learning algorithms are mathematical; they don't understand text, images, or raw log files. **Vectorization** is the critical process of converting this raw data into numerical vectors that algorithms can process.

- **What is a Vector?** An ordered list of numbers. Each number in the vector represents a "feature" or characteristic of a single data point (e.g., a customer, a document).
- **Mahout's Role:** Mahout operates on DRMs (Distributed Row Matrices), which are essentially large collections of these vectors distributed across a cluster.
- **Example:**
  - A customer might be represented as a vector: `[age, income, number_of_purchases, days_since_last_visit]`.
  - A text document might be converted into a TF-IDF vector, where each number represents the importance of a specific word.
- This transformation is the first and most fundamental step in any Mahout job.



**Data Vectorization**

Raw Input Data — Text Document, Customer Profile → Vectorization Process → Numerical Vector Representation

TF-IDF Vector: [0.1, 0.5, 0.0, 0.8, ...]

Customer Feature Vector: [45, 70000, 12 3]
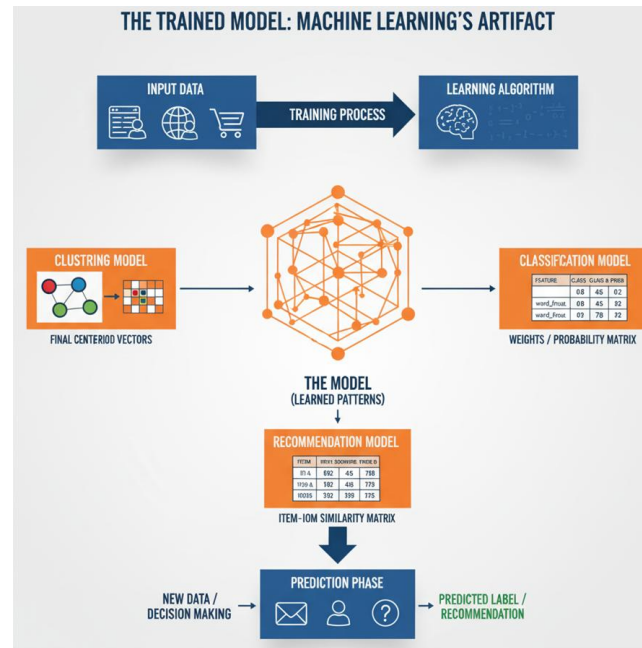
Mahout DRM: Collection of Vectors

# What is a "Model" in Distributed ML?

In the context of Mahout and distributed machine learning, a "model" is the **output of the training process**. It's the learned artifact that encapsulates the patterns discovered in the training data.

- A model is not abstract; it's a concrete data structure, which is often just another matrix or a set of vectors.
- **Examples of Models:**
    - **Clustering:** The model consists of the final centroid vectors (e.g., k vectors for K-Means).
    - **Classification:** The model could be a set of weights or probabilities for each feature and class (e.g., in Naive Bayes).
    - **Recommendations:** The model is often an item-item similarity matrix, where each row represents an item and the columns contain its similarity score to other items.

This trained model is then used during the prediction phase to make decisions on new data.



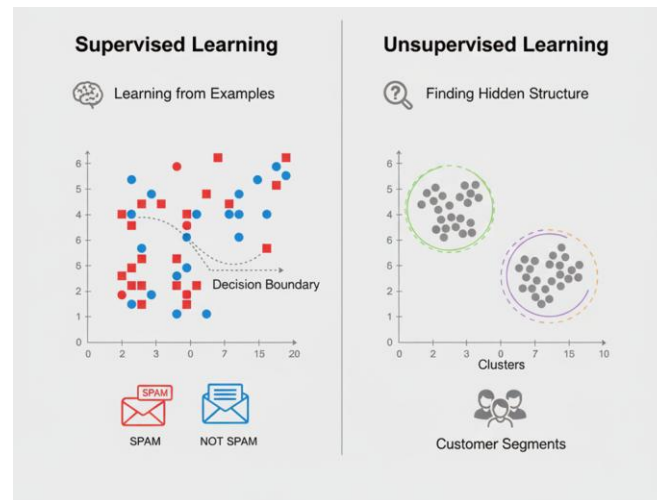THE TRAINED MODEL: MACHINE LEARNING'S ARTIFACT

# Fundamental Theory: Supervised vs. Unsupervised Learning

All machine learning algorithms fall into a few broad categories. Mahout provides tools for the two most common ones.

- **Supervised Learning (Learning from Examples)**
  - **Goal:** To predict a known outcome or label.
  - **Data Requirement:** Requires "labeled" training data, where the correct answer is already known (e.g., emails marked as "spam" or "not spam").
  - **Mahout Algorithms:** Classification algorithms like Naive Bayes and Random Forests.
- **Unsupervised Learning (Finding Hidden Structure)**
  - **Goal:** To discover inherent patterns, groupings, or structures in the data without any pre-existing labels.
  - **Data Requirement:** Works on unlabeled data.
  - **Mahout Algorithms:** Clustering algorithms like K-Means and Canopy Clustering.

# Feature Engineering at Scale

"Garbage in, garbage out." The quality of the features fed into a model determines the quality of its output. Feature engineering is the art and science of creating predictive variables from raw data.

- **Importance:** This is often the most time-consuming but most impactful part of the ML workflow.
- **Common Techniques:**
  - **TF-IDF (Term Frequency-Inverse Document Frequency):** As discussed, for converting text into meaningful numerical vectors.
  - **One-Hot Encoding:** Converting categorical variables (e.g., "country" = "USA", "Canada") into binary vectors ([1,0], [0,1]).
  - **Normalization/Standardization:** Scaling numerical features to a common range (e.g., 0 to 1) to prevent features with large values from dominating the algorithm.
- **Distributed Challenge:** These transformations must be performed efficiently on the entire distributed dataset before training.



Feature Engineering

TF-IDF

text data → [w=2, x÷1] → [0.1, 0.5, 0.2]

One-Hot Encoding

| | Country | USA | Canada |
|---|---|---|---|
| USA | USA | 1.0 | [1 |
| Canada | USA | [1.1 | [0 ] |

Normization

Before: 10, 100, 500 → Scale (0-1) After: 02, 22, 1.0

Distributed Processing → Mahout Algorithms

# Evaluating Model Performance

A model is only useful if it performs well. Model evaluation involves using statistical metrics to quantify a model's predictive accuracy on data it has not seen before.

- **For Classification:**
  - **Confusion Matrix:** A table showing the relationship between predicted labels and actual labels (True Positives, False Positives, True Negatives, False Negatives).
  - **Accuracy:** The percentage of correct predictions overall.
  - **Precision & Recall:** Measures of correctness for a specific class, crucial for imbalanced datasets.
- **For Clustering:**
  - **Silhouette Score:** Measures how similar a data point is to its own cluster compared to other clusters.
  - **Intra-cluster and Inter-cluster Distances:** Good clustering has low distance within clusters and high distance between clusters.
- **Method:** Typically involves splitting the data into a **training set** and a **testing set**. The model is built on the training set and evaluated on the testing set.



| | Predicted Yes | Predicted No |
|---|---|---|
| **Actual Yes** | ⊘ TRUE POSITIVE (TP) | FALSE POSITIVE (FP ✗ |
| | ✗ NEGATIVE ✗ | ⊘ TRUE NEGATIVE (TN) |

Accuracy = (TP + TN) / Total
Precision = TP / (TP + FP)
Recall = TP / (TP + FN)
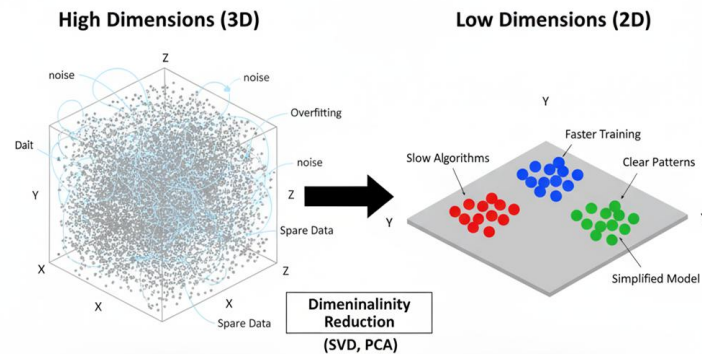
# The Challenge of Dimensionality

Big data often means "wide" data—datasets with thousands or even millions of features (dimensions). This leads to the **"Curse of Dimensionality."**

- **Problems with High Dimensions:**
  - **Computational Cost:** Algorithms become much slower.
  - **Overfitting:** Models may learn noise from the training data instead of the true signal, leading to poor performance on new data.
  - **Sparsity:** Data points become very spread out, making it difficult to find patterns.
- **Solution: Dimensionality Reduction**
  - Techniques to reduce the number of features while retaining as much of the important information as possible.
  - **Mahout's tools:** Mahout's core math library is well-suited for techniques like **Singular Value Decomposition (SVD)** and **Principal Component Analysis (PCA)**, which are powerful methods for dimensionality reduction.



High Dimensions (3D) — noise, Overfitting, Slow Algorithms, Dait, Spare Data → Dimeninalinity Reduction (SVD, PCA) → Low Dimensions (2D) — Faster Training, Clear Patterns, Simplified Model
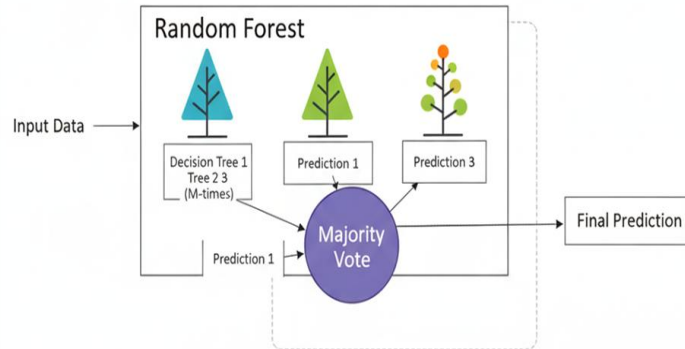
# Ensemble Methods: The Power of Many

Ensemble methods are techniques that combine the predictions of several base models to improve overall performance. The idea is that a "committee" of models will make better decisions than any single individual model.

- **Key Principle:** The individual models should be diverse. They can be trained on different subsets of the data or use different algorithms.
- **Example:** **Random Forests**
  - This is a powerful ensemble method provided by Mahout.
  - It builds a large number of individual **decision trees**.
  - Each tree is trained on a random sample of the data and a random subset of the features.
  - The final prediction is made by taking a majority vote of all the trees.
- **Benefit:** Ensembles are highly accurate and very robust against overfitting.
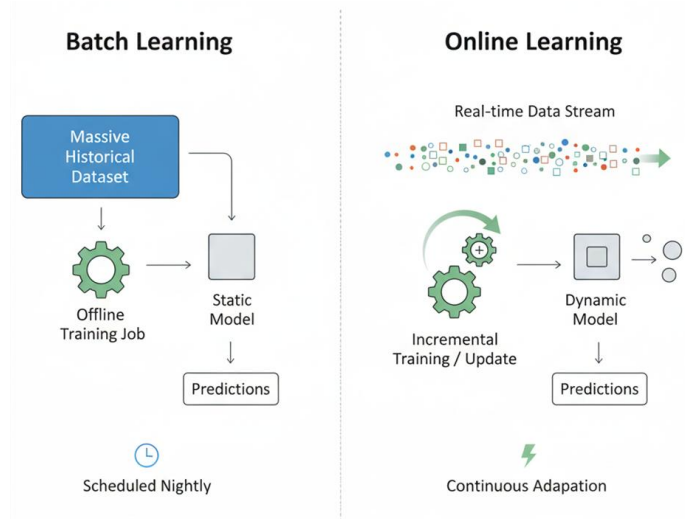
# Batch vs. Online Learning Paradigms

This is a key theoretical distinction in how machine learning models are trained.

- **Batch Learning (Mahout's primary focus):**
  - The model is trained on the **entire available dataset** at once.
  - This is done offline in a scheduled "batch" job (e.g., nightly).
  - The trained model is then deployed and used for predictions until the next batch run.
  - **Pros:** Simple, can leverage massive datasets.
  - **Cons:** Cannot adapt to new data in real-time; retraining can be slow and resource-intensive.
- **Online (or Stream) Learning:**
  - The model is updated incrementally as new data points arrive one by one or in mini-batches.
  - **Pros:** Adapts quickly to new trends; computationally less demanding at any given moment.
  - **Cons:** More complex to implement; sensitive to the order of data.

While Mahout is mainly a batch tool, its architecture can be used within larger systems that might incorporate stream processing (e.g., using Spark Streaming to prepare data for a Mahout batch job).

# Mahout's Environment and Development Workflow

**Primary Language:** Scala. The Mahout DSL is written in Scala, making it the most natural language for developing Mahout applications.

**Build Tools:** Projects using Mahout typically rely on:

- **Apache Maven:** A popular build automation tool for Java projects, widely used for Scala as well.
- **SBT (Scala Build Tool):** The de facto standard build tool for Scala projects, offering more native Scala integration.

**Execution Engines:** Apache Spark is the recommended and best-supported execution engine for modern Mahout (Samsara).

**Integration:** Mahout is consumed as a library. You add the Mahout Spark dependency to your `pom.xml` (Maven) or `build.sbt` (SBT) file, and then write your Scala code using the Mahout DSL.

**Command Line Interface (CLI):** Mahout also provides a CLI for running common algorithms with predefined parameters, useful for quick experiments or scripting batch jobs.

# Mahout vs. Spark MLlib: When to Choose Which?

| Feature | Apache Mahout (Samsara) | Apache Spark MLlib |
|---|---|---|
| **Primary Abstraction** | Distributed Row Matrix (DRM) and Scala DSL for linear algebra. | DataFrames for ML Pipelines (ML), RDDs for low-level API (MLlib). |
| **API Style** | R-like, algebraic, functional Scala DSL for expressing complex mathematical operations. | Pipeline API with `Transformer` and `Estimator` concepts, designed for end-to-end ML workflows. |
| **Core Focus** | A high-performance **linear algebra engine** and a library for *building and customizing* sophisticated distributed ML algorithms. | A comprehensive, end-to-end library with a **wider range of pre-built, easy-to-use algorithms** and ML pipeline tools. |
| **Use When...** | You need to implement **custom, mathematically complex algorithms**, require fine-grained control over algebraic operations, or want to explore novel ML approaches. | You need to quickly build standard ML pipelines, leverage a broad set of common algorithms, and prioritize ease of use and integration with Spark. |
| **Flexibility** | Backend-agnostic (Spark, Flink). | Tightly integrated with Spark; benefits from Spark's ecosystem. |
| **Flexibility** | Data scientists/ML engineers comfortable with linear algebra and Scala, building specialized solutions. | Data scientists/ML engineers looking for off-the-shelf, robust ML tools within the Spark ecosystem. |

# Future of Mahout and its Place in Big Data ML

- **Continued Evolution**: Mahout continues to evolve, adapting to new distributed computing paradigms and mathematical techniques.
- **Specialization:** While Spark MLlib aims for broad coverage, Mahout can serve as a powerful tool for more specialized, research-oriented, or highly optimized mathematical machine learning implementations.
- **Open Source Contribution**: As an Apache project, Mahout benefits from community contributions, ensuring its continued development and relevance.
- **Complementary, Not Competing**: In many real-world scenarios, Mahout and Spark MLlib can be seen as complementary tools. One might use MLlib for standard tasks and Mahout for custom components or for algorithms not yet optimized in MLlib.
- **Focus on Primitives:** Mahout's strength lies in providing robust, distributed linear algebra primitives, allowing users to build almost any algorithm from scratch with confidence in its scalability.

# Summary: Key Mahout Concepts

Apache Mahout is an open-source library for **scalable machine learning on Big Data**.

It has transitioned from a Hadoop MapReduce-centric framework to **Samsara**, leveraging **Apache Spark** (and other engines) for improved performance and flexibility.

**Distributed Row Matrix (DRM)** is its core data abstraction for distributed linear algebra.

Mahout's **Scala-based DSL** and **Algebraic Optimizer** simplify the development and ensure efficient execution of complex ML algorithms.

It provides powerful implementations for key ML tasks:

- **Recommender Systems** (e.g., Correlated Co-occurrence)
- **Clustering** (e.g., K-Means, Canopy)
- **Classification** (e.g., Naive Bayes, Random Forests)

Mahout is a valuable tool for developers and data scientists seeking to build **custom, high-performance, distributed ML solutions**.