# Pricilla Nakyazze Research Discussion 2

## 2025-06-12

Spotify uses collaborative filtering (Implicit Data) to recommend music.

Unlike explicit ratings (e.g., 5-star scales), Spotify uses implicit data: whether a user listened to a song, how often, and how recently.

This is often treated as binary (listened = 1, not listened = 0), and used to infer preferences indirectly.

The core algorithm used is Alternating Least Squares (ALS):

It factorizes the user-item interaction matrix (e.g., user-song) into two lower-dimensional matrices: user factors and item factors.

These factors are optimized to approximate the original interaction matrix.

Formalized as predicting:

Spotify Modifies ALS to work with confidence levels from implicit feedback.

Introduces a confidence matrix to weigh observed vs. unobserved interactions.

Regression Interpretation Matrix factorization with ALS can be interpreted as regularized regression problems solved iteratively:

Fix item factors and solve for user factors.

Fix user factors and solve for item factors.

Repeat until convergence.

Industrial Data Management Challenges spotify faces. Scalability & Spark Apache Spark is used to scale ALS to massive user-item datasets.

Spark parallelizes matrix operations using distributed computing, but challenges include:

Large memory requirements

Complex data partitioning

Why It Doesn't Scale Perfectly Each iteration requires shuffling data (e.g., new user or item blocks).

Loading new data during training slows down performance.

Partitioning the data cleverly (e.g., all of a user's ratings in one block) improves performance, reduces communication overhead.

Kryo Serialization Spark uses Kryo for efficient object serialization.

Kryo is faster and more compact than Java serialization, improving performance during ALS model training and data shuffling.

Key Takeaways Most Interesting Mathematical Points: Implicit matrix factorization allows learning from non-rating data (e.g., listening history).

ALS transforms a complex optimization problem into alternating regression subproblems.

Matrix factorization predicts preferences from a latent space representing hidden tastes.

Most Important System Challenges: Handling scale: Spotify must manage millions of users and items.

Performance bottlenecks: Iterative model training, data loading, and shuffling are major challenges.

Efficient serialization (e.g., using Kryo) and smart partitioning help mitigate these issues.