# Large Scale Recommender Systems in Spark
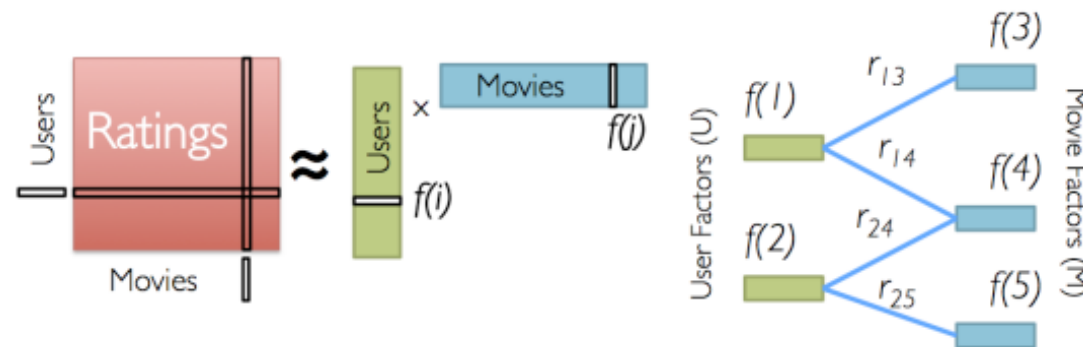
Kirk Hunter, Mrunmayee H. Bhagwat, Swetha Reddy

# Data & Tools

- Yahoo! Music Ratings (user id, song id, rating)

- Train: 700 million ratings, 1.8 million users, 136K songs

- Test: 18 million ratings, 1.8 million users, 136K songs
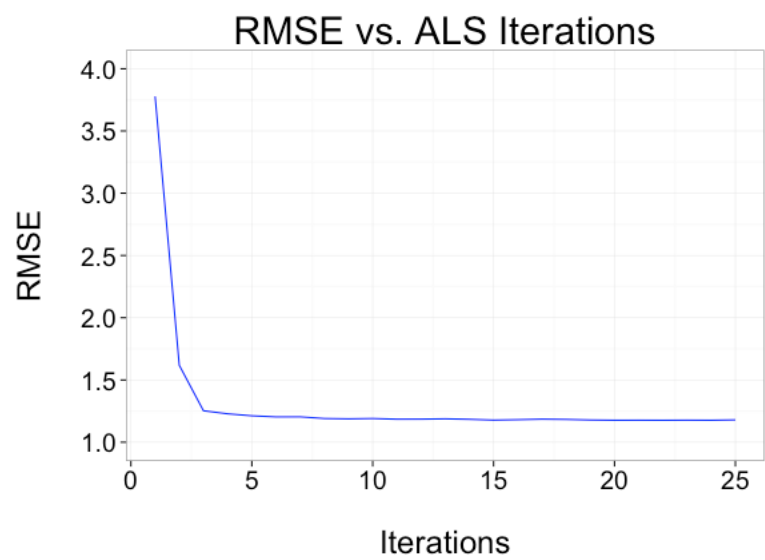
- Stored data in S3

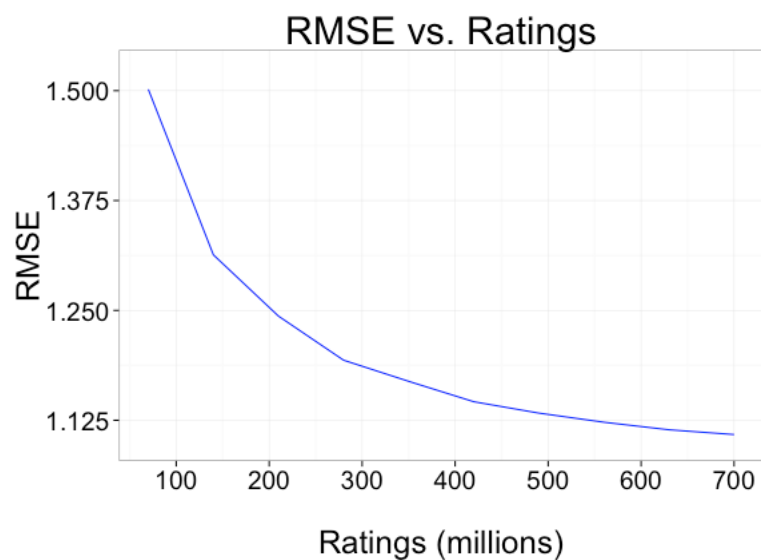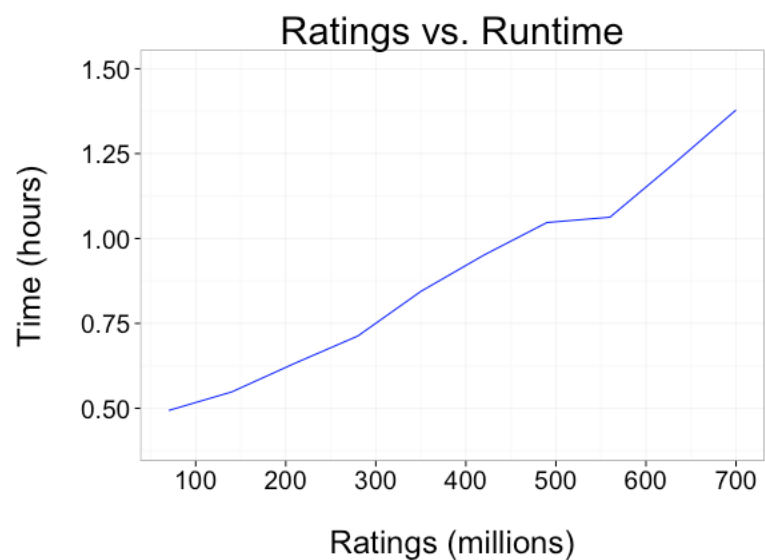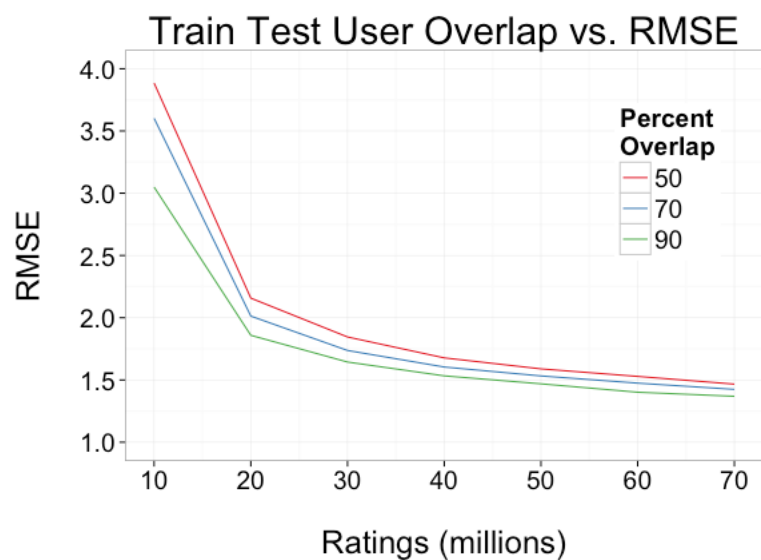- 5 node cluster running Spark on EM
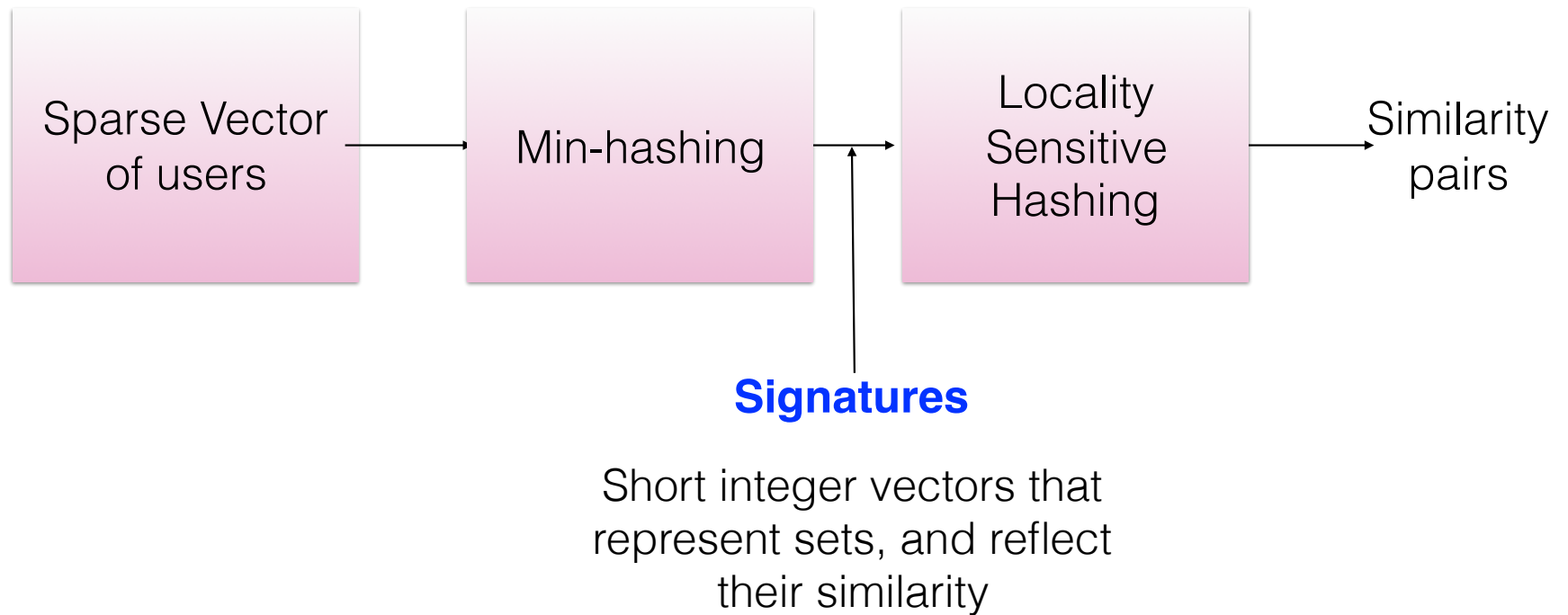
# Alternating Least Squares (ALS)

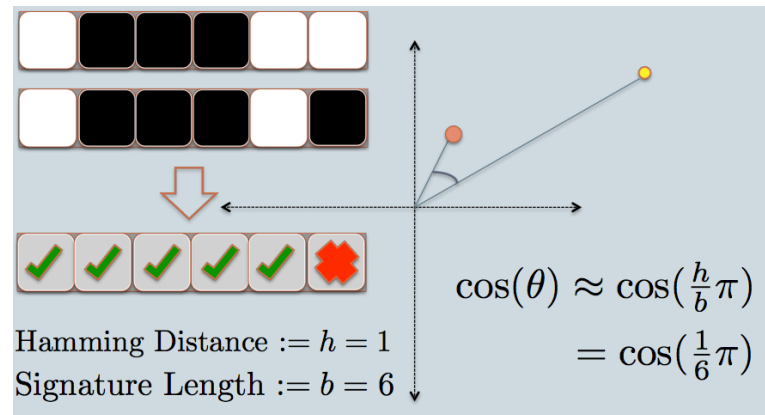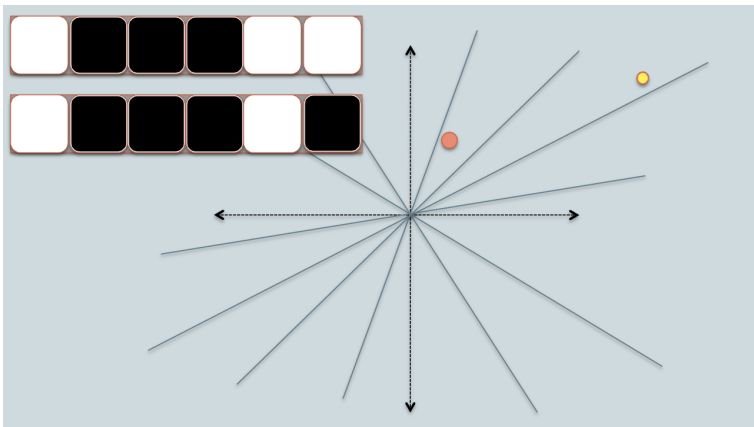Low-Rank Matrix Factorization:



Iterate:

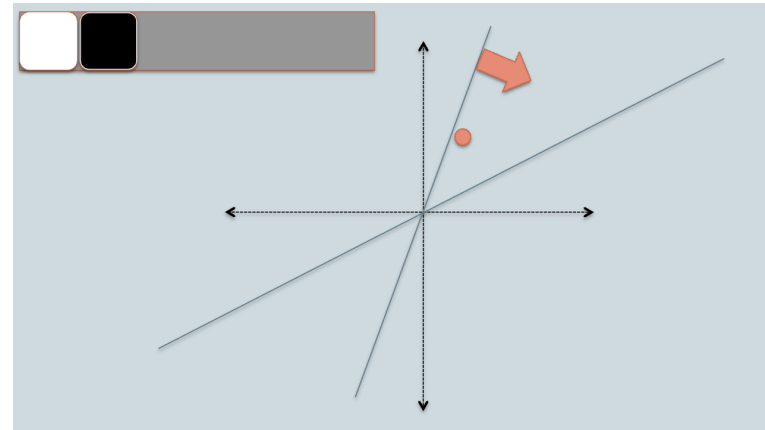$$f[i] = \arg\min_{w \in \mathbb{R}^d} \sum_{j \in \text{Nbrs}(i)} \left(r_{ij} - w^T f[j]\right)^2 + \lambda ||w||_2^2$$

## Train Test User Overlap vs. RMSE

Percent Overlap
— 50
— 70
— 90

## Ratings vs. Runtime

## RMSE vs. Ratings

## RMSE vs. ALS Iterations

# Locality Sensitive Hashing (LSH)

Sparse Vector of users → Min-hashing → Locality Sensitive Hashing → Similarity pairs

**Signatures**

Short integer vectors that represent sets, and reflect their similarity

# LSH - Cosine Similarity



Hamming Distance := $h = 1$
Signature Length := $b = 6$

$$\cos(\theta) \approx \cos(\frac{h}{b}\pi)$$
$$= \cos(\frac{1}{6}\pi)$$
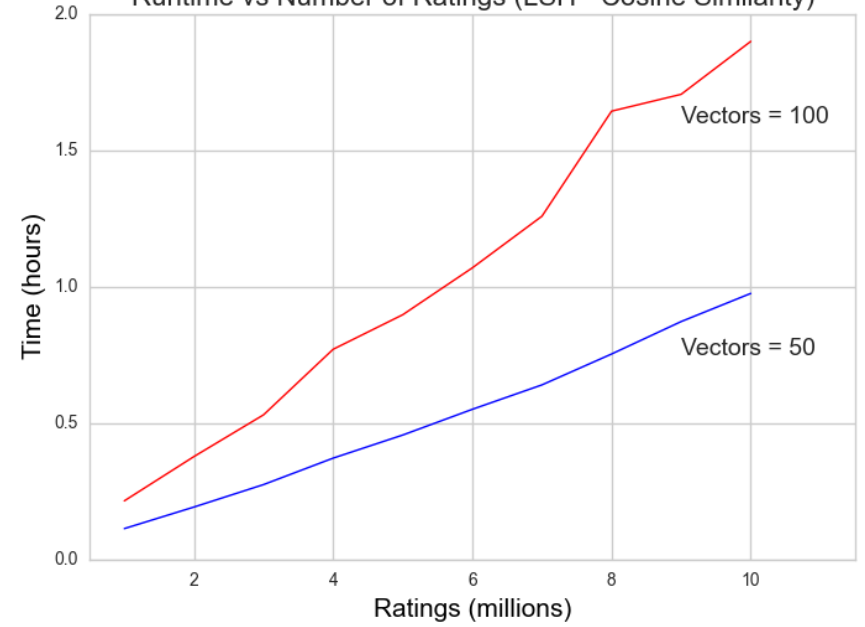
Runtime vs. Ratings (LSH - Jaccard Similarity)

Runtime vs Number of Ratings (LSH - Cosine Similarity)

# Thank You!