

# Implement Distributed Alternating Least Squares Algorithm for Matrix Completion

Varun Gandhi (vg292)

# Netflix Problem

- $V$ :  $m \times n$  matrix
- complete the matrix

	<i>Avatar</i>	<i>The Matrix</i>	<i>Up</i>
<i>Alice</i>	?	4	2
<i>Bob</i>	3	2	?
<i>Charlie</i>	5	?	3

- $W$ :  $m \times r$  (row-factor matrix)
- $H$ :  $r \times n$  (column-factor matrix)
- $W \cdot H \approx V$
- Loss function  $(V_{ij} - WH_{ij})^2$

# Motivation

Large applications involve matrices with

- millions of rows x columns;
- billions of entries

To achieve high-performance

- parallel & distributed factorisation
- keep the loss to minimum

# Algorithm

## Sequential Computation

- Initial point  $W_0$  and  $H_0$
- ALS solved for every row & column

$$\text{Compute } W_{n+1}: (\forall i) \underline{W}_{i*} \mathbf{H}_n^{(i)} = \mathbf{V}_{i*},$$

$$\text{Compute } \mathbf{H}_{n+1}: (\forall j) \mathbf{W}_{n+1}^{(j)} \underline{\mathbf{H}}_{*j} = \mathbf{V}_{*j},$$

## Parallel Computation

- Parallelise computation for rows and columns respectively

# Algorithm

## Distributed Computation

- Partition (block) the matrix with  $m_b * n_b$  matrices
- every node updates a matrix block

## Why Spark?

- In-memory algorithm
- Matrix versions cached in memory

# Progress

- Revising all linear algebra concepts
- Getting familiar with Scala and Spark
- Trying examples in Python