# Distributed CUR Decomposition for Bi-Clustering

Stephen Kline, Kevin Shaw
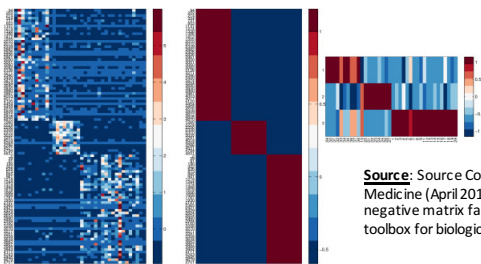
June 1, 2016

{sakline, keshaw}@stanford.edu

Stanford University, CME 323 Final Project
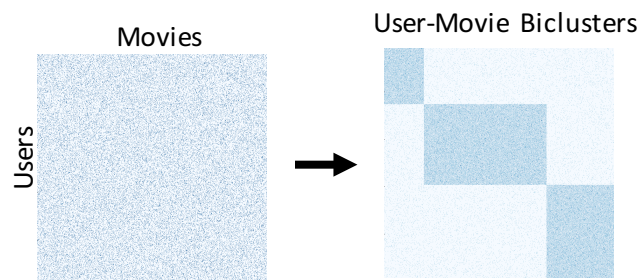
# Review of SVD: A = U$\Sigma$V$^T$

**U**

dense
big

$\Sigma$

**sparse / small**

**V$^T$**

dense / big

- **PRO - High accuracy**
  - k singular values/vectors produce the best k-rank approximation to A

- **CON - High computation / space requirements**
  - In our biclustering application with MovieLens data, the distributed SVD is "roughly square" - ARPACK (vs. "tall and skinny" $- A^TA$ trick)

**A**

sparse / huge

# Design Decisions for Distributed CUR



C
sparse
big

U
**dense / small**

Compute U locally

R
sparse / big

Only necessary to store C and R as set of indices into A

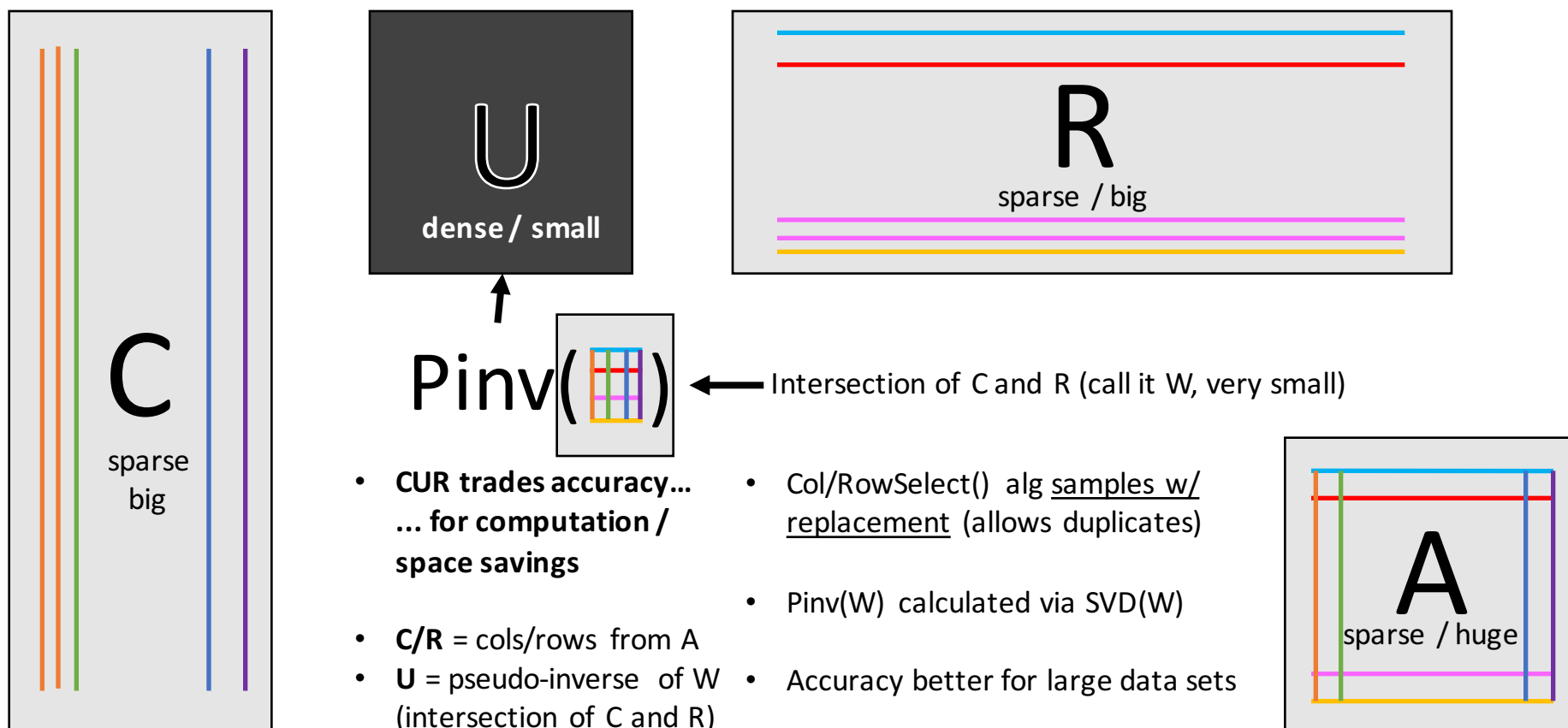There are multiple variations of CUR. We selected the algorithm as presented in: Drineas, et. al., 2006. "Fast Monte Carlo Algorithms for Matrices III" which (for example) does not remove duplicate cols/rows as some others do.

**Key Design Decision:**
Distribute <u>two instances</u> of A avoiding future all-to-all communications

A
sparse / huge

# Serial vs. Distributed CUR - Asymptotics

Serial

- Build C and R:
  - **Generate probabilities – O(mn)**
  - Create C matrix – O(mk)
  - Create R matrix – O(nk)

- Construct U
  - Compute $C^T C$ – $O(mk^2)$
  - SVD of $C^T C$ – $O(k^3)$
  - Compute A and B – $O(k^3)$
  - $U = AB^T$ – $O(k^3)$

Distributed (communication <u>cost</u> and computation <u>time</u>)

- Build C and R:
- Generate probabilities – O(mn + p) cost, O(max dense) time
  - **Create 2 RDDs by Row/Col partition – O(mn) cost, AtoA**
  - **Both instances: reduce to Row/Col sums —**
    **O(max dense) time, no communication**
  - One instance: reduce Row sum to total – O(p) cost, O(log p) time
  - Broadcast total to calculate probs – O(p) cost, O(log p) time
- Create C / R matrices
  - Locally sample k rows/cols – O(k)
  - Broadcast sample to RDDs – O(pk) cost, O(k log p) time

- Construct U
  - Same as Serial (less opportunity to distribute)

# Biclustering: Distributed CUR vs SVD - Empirics