

# Thresholding Begets Descent

Inderjit S. Dhillon  
The University of Texas at Austin

Banff International Research Station  
Jan 20, 2011

Joint work with Prateek Jain, Raghu Meka and Ambuj Tewari

# Overview

- Hard Thresholding for Compressed Sensing
  - New Family of Algorithms with Guarantees
- Hard Thresholding for Matrix Completion
- Digression at End
  - Fast, Memory-Efficient Dimensionality Reduction of Massive Graphs

# Compressed Sensing and Rank Minimization

$$\begin{aligned} \text{(CS)} : \quad & \min_{\mathbf{x}} \|\mathbf{x}\|_0 \\ & \text{s.t. } A\mathbf{x} = b. \end{aligned}$$

$$\mathbf{x} \in \mathbb{R}^n, A: \mathbb{R}^n \rightarrow \mathbb{R}^d, b \in \mathbb{R}^d.$$

$$\begin{aligned} \text{(ARMP)} : \quad & \min_X \text{rank}(X) \\ & \text{s.t. } \mathcal{A}(X) = b. \end{aligned}$$

$$X \in \mathbb{R}^{m \times n}, \mathcal{A}: \mathbb{R}^{m \times n} \rightarrow \mathbb{R}^d, b \in \mathbb{R}^d.$$

- Can view CS as specific instance of ARMP with  $X = \text{Diag}(\mathbf{x})$ .

# An Example: Minimum Rank Matrix Completion

- Netflix Challenge:
  - Given a few user-movie ratings
  - **Goal:** complete ratings matrix
- Small number of latent factors  $\equiv$  low-rank
- Special case of ARMP:

$$\begin{aligned} \text{(MCP)} : \quad & \min_X \text{rank}(X) \\ & \text{s.t. } \text{tr}(X\mathbf{e}_j\mathbf{e}_i^T) = b_{ij}, \forall (i,j) \in \Omega. \end{aligned}$$

- Typically, number of samples very small: Netflix has 1% samples

# CS and ARMP

Technique	CS	ARMP
Convex relaxation	$\ell_1$ (Lasso)	Trace-norm (SVT)
Greedy approach	MP, OMP, CoSamp	ADMiRA
Hard Thresholding	IHT, GradeS	<b>SVP</b> , IHT

Table: CS vs ARMP

# Restricted Isometry Property (RIP)

- Most CS methods assume RIP:

$$(1 - \delta_k) \|\mathbf{x}\|^2 \leq \|\mathbf{A}\mathbf{x}\|^2 \leq (1 + \delta_k) \|\mathbf{x}\|^2, \quad \forall \mathbf{x} \text{ s.t. } \|\mathbf{x}\|_0 \leq k$$

- Generalization to ARMP:

$$(1 - \delta_k) \|X\|_F^2 \leq \|\mathcal{A}(X)\|_2^2 \leq (1 + \delta_k) \|X\|_F^2, \quad \forall X \text{ s.t. } \text{rank}(X) \leq k$$

- Families satisfying RIP:

$$\mathcal{A}(X) = A \text{vec}(X),$$

- $A_{ij} \sim \mathcal{N}(0, 1/d)$
- $A_{ij} = \begin{cases} 1/\sqrt{d} & \text{with probability } 1/2 \\ -1/\sqrt{d} & \text{with probability } 1/2 \end{cases}$

# Projected Gradient

- Consider

$$\begin{aligned} \min_x \psi(x) &= \frac{1}{2} \|Ax - b\|_2^2, \\ \text{s.t. } x &\in \mathcal{C}(k) = \{x : \text{supp}(x) \leq k\}. \end{aligned}$$

- Adapt classical projected gradient
- Efficient projection onto non-convex support constraint

# Iterative Hard Thresholding (IHT)

---

## Algorithm 1 IHT/GradeS Algorithm

---

Initialize  $x^0 = 0$ ,  $t = 0$

Set step size  $\eta_t$

**repeat**

$$x^{t+1} = P_k(x^t - \eta_t \underbrace{A^*(Ax^t - b)}_{\nabla\psi(x)})$$

$t = t + 1$

**until** Convergence

---



---

## Algorithm 2 IHT-Newton Algorithm

---

Initialize  $x^0 = 0$ ,  $t = 0$

Set step size  $\eta_t$

**repeat**

$$y^{t+1} = P_k(x^t - \eta_t \underbrace{A^*(Ax^t - b)}_{\nabla\psi(x)})$$

$$x^{t+1} = \operatorname{argmin}_{\operatorname{supp}(x)=\operatorname{supp}(y^{t+1})} \|Ax - b\|^2$$

$$t = t + 1$$

**until** Convergence

---

- Simple analysis—apply RIP twice and Projection property once
- $\psi(x) = \frac{1}{2} \|Ax - b\|_2^2$

$$\psi(x^{t+1}) - \psi(x^t) = \langle \nabla \psi(x^t), x^{t+1} - x^t \rangle + \frac{1}{2} \|A \overbrace{(x^{t+1} - x^t)}^{\text{supp } 2k}\|^2$$

- Simple analysis—apply RIP twice and Projection property once
- $\psi(x) = \frac{1}{2} \|Ax - b\|_2^2$

$$\begin{aligned} \psi(x^{t+1}) - \psi(x^t) &= \langle \nabla \psi(x^t), x^{t+1} - x^t \rangle + \frac{1}{2} \overbrace{\|A(x^{t+1} - x^t)\|_2^2}^{\text{supp } 2k} \\ &\leq \langle \nabla \psi(x^t), x^{t+1} - x^t \rangle + \frac{1}{2} \underbrace{(1 + \delta_{2k}) \|x^{t+1} - x^t\|_2^2}_{\text{Using RIP}}, \end{aligned}$$

- Simple analysis—apply RIP twice and Projection property once
- $\psi(x) = \frac{1}{2}\|Ax - b\|_2^2$

$$\begin{aligned}
 \psi(x^{t+1}) - \psi(x^t) &= \langle \nabla \psi(x^t), x^{t+1} - x^t \rangle + \frac{1}{2} \overbrace{\|A(x^{t+1} - x^t)\|_2^2}^{\text{supp } 2k} \\
 &\leq \langle \nabla \psi(x^t), x^{t+1} - x^t \rangle + \frac{1}{2} \underbrace{(1 + \delta_{2k}) \|x^{t+1} - x^t\|_2^2}_{\text{Using RIP}}, \\
 &= \frac{1}{2} (1 + \delta_{2k}) \|x^{t+1} - y^{t+1}\|_2^2 - \frac{1}{2(1 + \delta_{2k})} \|A^T(Ax^t - b)\|_2^2
 \end{aligned}$$

- Simple analysis—apply RIP twice and Projection property once
- $\psi(x) = \frac{1}{2}\|Ax - b\|_2^2$

$$\begin{aligned}
 \psi(x^{t+1}) - \psi(x^t) &= \langle \nabla \psi(x^t), x^{t+1} - x^t \rangle + \frac{1}{2} \overbrace{\|A(x^{t+1} - x^t)\|_2^2}^{\text{supp } 2k} \\
 &\leq \langle \nabla \psi(x^t), x^{t+1} - x^t \rangle + \frac{1}{2} \underbrace{(1 + \delta_{2k})\|x^{t+1} - x^t\|_2^2}_{\text{Using RIP}}, \\
 &= \frac{1}{2}(1 + \delta_{2k})\|x^{t+1} - y^{t+1}\|_2^2 - \frac{1}{2(1 + \delta_{2k})}\|A^T(Ax^t - b)\|_2^2
 \end{aligned}$$

$$\text{where } y^{t+1} = x^t - \frac{1}{1 + \delta_{2k}} \nabla \psi(x^t), \quad x^{t+1} = P_k(y^{t+1})$$

- Simple analysis—apply RIP twice and Projection property once
- $\psi(x) = \frac{1}{2} \|Ax - b\|_2^2$

$$\begin{aligned}
 \psi(x^{t+1}) - \psi(x^t) &= \langle \nabla \psi(x^t), x^{t+1} - x^t \rangle + \frac{1}{2} \overbrace{\|A(x^{t+1} - x^t)\|_2^2}^{\text{supp } 2k} \\
 &\leq \langle \nabla \psi(x^t), x^{t+1} - x^t \rangle + \frac{1}{2} \underbrace{(1 + \delta_{2k}) \|x^{t+1} - x^t\|_2^2}_{\text{Using RIP}}, \\
 &= \frac{1}{2} (1 + \delta_{2k}) \|x^{t+1} - y^{t+1}\|_2^2 - \frac{1}{2(1 + \delta_{2k})} \|A^T(Ax^t - b)\|_2^2 \\
 &\leq \frac{1}{2} (1 + \delta_{2k}) \underbrace{\|x^* - y^{t+1}\|_2^2}_{\text{Projection}} - \frac{1}{2(1 + \delta_{2k})} \|A^T(Ax^t - b)\|_2^2
 \end{aligned}$$

- Simple analysis—apply RIP twice and Projection property once

- $\psi(x) = \frac{1}{2} \|Ax - b\|_2^2$

$$\begin{aligned}\psi(x^{t+1}) - \psi(x^t) &\leq \frac{1}{2}(1 + \delta_{2k}) \underbrace{\|x^* - y^{t+1}\|_2^2}_{\text{Projection}} - \frac{1}{2(1 + \delta_{2k})} \|A^T(Ax^t - b)\|_2^2 \\ &= \langle \nabla \psi(x^t), x^* - x^t \rangle + \frac{1}{2}(1 + \delta_{2k}) \|x^* - x^t\|_2^2\end{aligned}$$

- Simple analysis—apply RIP twice and Projection property once

- $\psi(x) = \frac{1}{2} \|Ax - b\|_2^2$

$$\begin{aligned}
 \psi(x^{t+1}) - \psi(x^t) &\leq \frac{1}{2}(1 + \delta_{2k}) \underbrace{\|x^* - y^{t+1}\|^2}_{\text{Projection}} - \frac{1}{2(1 + \delta_{2k})} \|A^T(Ax^t - b)\|^2 \\
 &= \langle \nabla \psi(x^t), x^* - x^t \rangle + \frac{1}{2}(1 + \delta_{2k}) \|x^* - x^t\|^2 \\
 &\leq \langle \nabla \psi(x^t), x^* - x^t \rangle + \underbrace{\frac{1}{2} \frac{1 + \delta_{2k}}{1 - \delta_{2k}} \|A(x^* - x^t)\|^2}_{\text{Using RIP}}
 \end{aligned}$$



- Simple analysis—apply RIP twice and Projection property once

- $\psi(x) = \frac{1}{2} \|Ax - b\|_2^2$

$$\begin{aligned}
 \psi(x^{t+1}) - \psi(x^t) &\leq \frac{1}{2}(1 + \delta_{2k}) \underbrace{\|x^* - y^{t+1}\|^2}_{\text{Projection}} - \frac{1}{2(1 + \delta_{2k})} \|A^T(Ax^t - b)\|^2 \\
 &= \langle \nabla \psi(x^t), x^* - x^t \rangle + \frac{1}{2}(1 + \delta_{2k}) \|x^* - x^t\|^2 \\
 &\leq \langle \nabla \psi(x^t), x^* - x^t \rangle + \underbrace{\frac{1}{2} \frac{1 + \delta_{2k}}{1 - \delta_{2k}} \|A(x^* - x^t)\|^2}_{\text{Using RIP}} \\
 &= \psi(x^*) - \psi(x^t) + \frac{\delta_{2k}}{(1 - \delta_{2k})} \|A(x^* - x^t)\|^2,
 \end{aligned}$$

- Simple analysis—apply RIP twice and Projection property once

- $\psi(x) = \frac{1}{2} \|Ax - b\|_2^2$

$$\begin{aligned}
 \psi(x^{t+1}) - \psi(x^t) &\leq \frac{1}{2} (1 + \delta_{2k}) \underbrace{\|x^* - y^{t+1}\|^2}_{\text{Projection}} - \frac{1}{2(1 + \delta_{2k})} \|A^T(Ax^t - b)\|^2 \\
 &= \langle \nabla \psi(x^t), x^* - x^t \rangle + \frac{1}{2} (1 + \delta_{2k}) \|x^* - x^t\|^2 \\
 &\leq \langle \nabla \psi(x^t), x^* - x^t \rangle + \frac{1}{2} \underbrace{\frac{1 + \delta_{2k}}{1 - \delta_{2k}} \|A(x^* - x^t)\|^2}_{\text{Using RIP}} \\
 &= \psi(x^*) - \psi(x^t) + \frac{\delta_{2k}}{(1 - \delta_{2k})} \|A(x^* - x^t)\|^2,
 \end{aligned}$$

For exact case,  $\psi(x^*) = 0$ ,  $A(x^*) = b$ . Hence,

$$\begin{aligned}
 \psi(x^{t+1}) &\leq \underbrace{\frac{2\delta_{2k}}{(1 - \delta_{2k})}}_{<1 \text{ for } \delta_{2k} < 1/3} \psi(x^t).
 \end{aligned}$$

# IHT: Recovery Guarantees

- Suppose  $b = Ax^*$ . When  $\delta_{2k} < 1/3$ , IHT outputs  $x$  st  $\|Ax - b\|_2^2 \leq \epsilon$  in  $\left\lceil C \log \frac{\|b\|_2^2}{2\epsilon} \right\rceil$  iterations [Garg & Khandekar, 2009].
- Similar geometric convergence for noisy case,  $b = Ax^* + e$ .
- Update: [Foucart, 2010] shows geometric convergence when  $\delta_{3k} < 1/2$  with step size  $\eta = 1$  (improved to  $\delta_{3k} < 1/\sqrt{3}$  in [Foucart, 2011]).
- Similar guarantees for IHT-Newton, which empirically works better.

# Structure of Gradient Step

- Consider IHT-Newton iterate  $x^t$
- Let  $\text{supp}(x^*) = S^*$ ,  $\text{supp}(x^t) = S^t$ , and  $J = (S^* \cup S^t)^c$
- Since  $A_{S^t}^*(Ax^t - b) = 0$ , gradient step:  $y_{t+1} = x^t - \eta_t A^*(Ax^t - b)$  has the form:

$$y_{t+1} = \begin{bmatrix} x_t \\ -\eta_t A_{S^* - S^t}^*(Ax^t - b) \\ -\eta_t A_J^*(Ax^t - b) \end{bmatrix}$$

# New Algorithm — OMP( $\ell$ ) with Replacement

---

## Algorithm 3 OMPR( $\ell$ ) Algorithm

---

Initialize  $x^0 = A^*b$ ,  $t = 0$

Set step size  $\eta_t$

**repeat**

$$y^{t+1} = x^t - \eta_t \underbrace{A^*(Ax^t - b)}_{\nabla\psi(x)}$$

$$C = \text{supp}(x^t) \cup \text{Top } \ell \text{ indices of } y_{S_t^c}^{t+1}$$

$$y^{t+1} = P_k(y_C^{t+1})$$

$$x^{t+1} = \operatorname{argmin}_{\text{supp}(x)=\text{supp}(y^{t+1})} \|Ax - b\|^2$$

$$t = t + 1$$

**until** Convergence

---

- When  $\ell = 1$ , OMPR(1) replaces one “working-set” index at a time.
- Note: OMPR( $k$ ) is IHT-Newton.

# Recovery Guarantees for OMP( $\ell$ )

- Property 1:  $\|y_D^{t+1}\| > \|x_U^t\|$ , where “desired” index set  $D = S^* - S_t$ , and “undesired” index set  $U = S_t - S^*$ .
- Property 2:  $\psi(x^t) - \psi(x^{t+1}) > c$
- Recovery Guarantee: Suppose  $b = Ax^*$ . When  $\delta_{2k} < 1/2$  and  $\eta = .9999$ , OMP(1) converges to  $x^*$  in  $O(k)$  iterations.
- Similar guarantee for noisy case
- Gives “best-known” RIP guarantee
- Time complexity =  $O(k * n * d)$ , but least squares update at each iteration is cheap.

# Matrix Completion

- Complete a low-rank matrix from few sampled entries
- Minimum rank matrix completion problem:

$$\begin{aligned}(\mathbf{MCP}) : \min_X \text{rank}(X), \\ \text{s.t. } \mathcal{P}_\Omega(X) = \mathcal{P}_\Omega(X^*).\end{aligned}$$

- $\mathcal{P}_\Omega : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}^{m \times n}$ —projection onto index set  $\Omega$ , i.e.,

$$(\mathcal{P}_\Omega(X))_{ij} = \begin{cases} X_{ij} & \text{for } (i, j) \in \Omega \\ 0 & \text{otherwise} \end{cases}$$

- Special case of ARMP: SVP can be applied directly
- **Problem:** MCP does not satisfy RIP in general

# SVP: Hard Thresholding for Matrix Completion

$$\begin{aligned} \min_X \psi(X) &= \frac{1}{2} \|P_\Omega(X - X^*)\|_F^2, \\ \text{s.t } X &\in \mathcal{C}(k) = \{X : \text{rank}(X) \leq k\}. \end{aligned}$$

---

## Algorithm 4 SVP for Matrix Completion

---

Initialize  $X^0 = 0$ ,  $t = 0$

Set step size  $\eta_t = 1/(1 + \delta)\rho$ ,  $\rho$ =sampling density,  $\delta$  is a parameter

**repeat**

$$X^{t+1} = P_k(X^t - \eta_t P_\Omega(X^t - X^*))$$

$$t = t + 1$$

**until** Convergence

---

- $P_k(X) = U_k \Sigma_k V_k^T$ : top  $k$  singular vectors of  $X$
- Computation of  $k$  singular vectors of:  $\underbrace{X^t}_{\text{low rank}} - \eta_t \underbrace{P_\Omega(X^t - X^*)}_{\text{sparse}}$
- Matrix-vector multiplication:  $O((m + n)k + |\Omega|)$



# SVP-Newton for Matrix Completion

---

## Algorithm 5 SVP-Newton

---

Initialize  $X^0 = 0$ ,  $t = 0$

Set step size  $\eta_t = 1/(1 + \delta)p$ ,  $p$ =sampling density,  $\delta$  is a parameter

**repeat**

$Y^{t+1} = U_k \Sigma_k V_k^T$ , where  $\text{svd}(X^t - \eta_t P_\Omega(X^t - X^*)) = U \Sigma V^T$

Given  $U_k$  and  $V_k$ , compute:

$$S_k = \underset{S}{\operatorname{argmin}} \|P_\Omega(U_k S V_k^T - X^*)\|_F^2$$

$X^{t+1} = U_k S_k V_k^T$

$t = t + 1$

**until** Convergence

---

# OMPR( $\ell$ ) for Matrix Completion

---

**Algorithm 6** OMPR( $\ell$ ) for Matrix Completion

---

Initialize  $X^0 = U_k \Sigma_k V_k^T$ , where  $\text{svd}(P_\Omega(X^*)) = U \Sigma V^T$ . Set  $t = 0$   
Set step size  $\eta_t = 1/(1 + \delta)p$ ,  $p$ =sampling density,  $\delta$  is a parameter

**repeat**

    Given  $X_t = U_k \Sigma_k V_k^T$ , compute:

$$S_k = \underset{S}{\operatorname{argmin}} \|P_\Omega(U_k S V_k^T - X^*)\|_F^2$$

    Compute svd of  $S_k$  to get svd of  $U_k S_k V_k^T = \bar{U}_k \bar{\Sigma}_k \bar{V}_k^T$

    Compute top  $\ell$  singular triplets of  $P_\Omega(\bar{U}_k \bar{\Sigma}_k \bar{V}_k^T - X^*)$  and take union with  $\{\bar{U}_k, \bar{\Sigma}_k, \bar{V}_k^T\}$

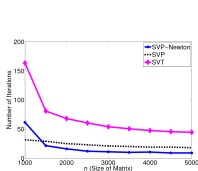
    Drop bottom  $\ell$  singular triplets of the above set to obtain  $X^{t+1}$

$t = t + 1$

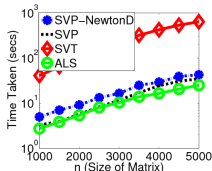
**until** Convergence

---

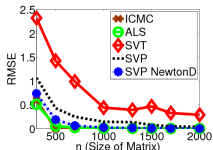
# Results: Matrix Completion for Synthetic Datasets



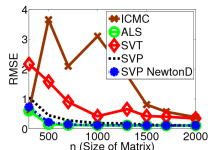
(a)



(b)



(c)



(d)

(a) Number of iterations required

(b) Time required for Noisy data

(c),(d) RMSE for Power-Law Sampling without/with Noise

# Results: Matrix Completion for MovieLens Dataset

$k$	SVP-NewtonD	SVP	ALS	SVT
2	0.90	1.15	0.88	1.06
3	0.89	1.14	0.87	0.98
5	0.89	1.09	0.86	0.95
7	0.89	1.08	0.86	0.93
10	0.90	1.07	0.87	0.91
12	0.91	1.08	0.88	0.90

Table: RMSE obtained by various methods

- **Problem:** Ratings matrix is not sampled uniformly

# Conclusions and Future Work

- New OMP with Replacement Algorithm
  - Guarantees recovery
  - Gives “best-known” RIP guarantee
- Future Work
  - Reduce Time Complexity (don’t compute full gradient)
  - Explore IHT-Newton or OMPR( $\ell$ ) in regression setting (when RIP does not hold)
  - Hard thresholding algorithms for other problems, e.g., sparse+low-rank matrix decomposition?

*Paper will soon be available at: <http://arxiv.org>*

# Massive Social Networks

# Testbed of Social Networks

Network	Date	# nodes	# links	# added links	% added links
Flickr	4/14/2007	1,990,149	41,302,536	–	–
	4/25/2007	1,990,149	42,056,754	754,218	1.8%
	5/6/2007	1,990,149	42,879,714	822,960	1.9%
LiveJournal	02/16/2009	1,770,961	83,663,478	–	–
	03/4/2009	1,770,961	84,413,542	750,064	0.8%
	04/03/2009	1,770,961	85,713,766	1,300,224	1.5%
MySpace	12/11/2008	2,137,264	90,333,122	–	–
	1/11/2009	2,137,264	90,979,264	646,142	0.7%
	2/14/2009	2,137,264	91,648,716	669,452	0.7%

	# clusters	avg size	% intra links	% inter links
Flickr	18	110,563	71.8%	28.2%
LiveJournal	17	106,241	72.5%	27.5%
MySpace	17	125,721	51.9%	48.1%

# Motivation

- Need to compute matrix spectral functions, e.g. Katz measure for predicting future friendships

$$(I - \beta A)^{-1} = V(I - \beta \Lambda)^{-1} V^T$$

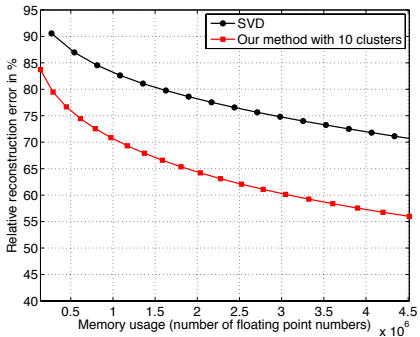
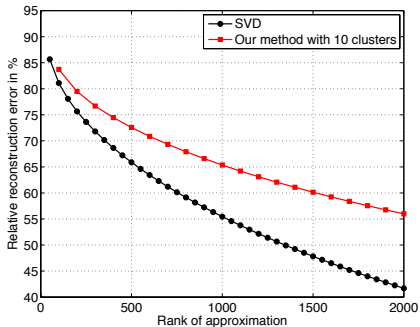
- Too expensive to compute for massive networks
- One solution: spectral approximation

$$f(A) \approx V f(\Lambda) V^T$$

- But SVD/PCA is wasteful (especially in terms of memory). Same for recent stochastic algorithms.



# Empirical Results



# Algorithm: Clustered low rank approximation

**Require:** An  $m \times m$  adjacency matrix  $A$  of a graph, number of clusters  $c$

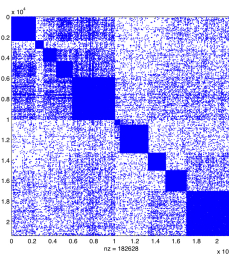
**Ensure:** Clustered low rank approximation of  $A$

- 1: Cluster the graph into  $c$  clusters
- 2: Compute a low rank approximation of each cluster

$$U_i S_i V_i \approx A_{ij}$$

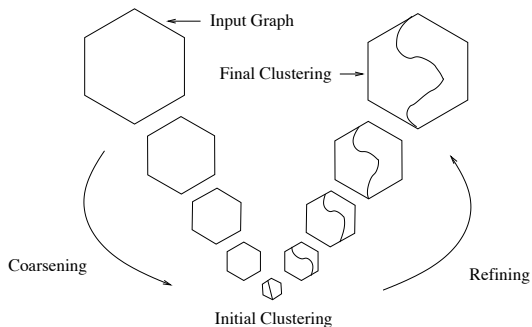
- 3: Extend the cluster-wise approximations, into an approximation for the entire matrix  $A$

$$S_{ij} = U_i^T A_{ij} V_j$$



# Graph Clustering: Multilevel Approach

- Overview:



[CHACO, Hendrickson & Leland, 1994]

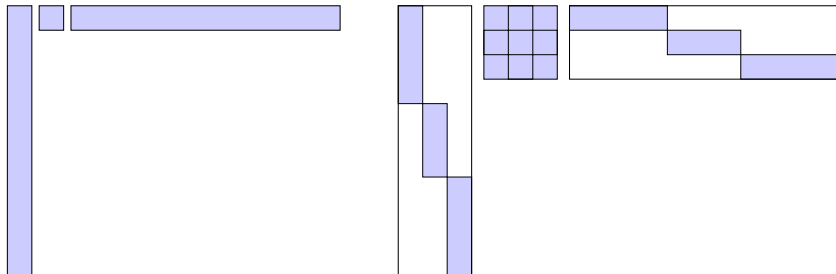
[METIS, Karypis & Kumar, 1999]

[GRACLUS, Dhillon, Guan & Kulis, 2005]

# Low rank vs clustered low-rank

Low rank:  $A \approx U\Sigma V^T$

Clustered low rank:  $A \approx \begin{bmatrix} U_1 & 0 \\ 0 & U_2 \end{bmatrix} \begin{bmatrix} S_{11} & S_{12} \\ S_{21} & S_{22} \end{bmatrix} \begin{bmatrix} V_1 & 0 \\ 0 & V_2 \end{bmatrix}^T$



- Observe  $\text{diag}(U_1, U_2, \dots, U_c)$  and has the same memory usage as  $U$