

Graph Algorithms with Sparse Linear Algebra

Tim Baer

C3.ai Platform Engineering
University of Illinois at Urbana-Champaign

July 19, 2021

Overview

Introduction

All Pairs Shortest Paths

Sparse Linear Algebra

Minimum Spanning Trees

Expanders

TODO:

- ▶ Introduce notation better/add notation slide?
- ▶ add pauses?
- ▶ Fix footnote numbers

Motivation

Graph algorithms can be rewritten with linear algebra. Adjacency matrices connect graphs with linear algebra.

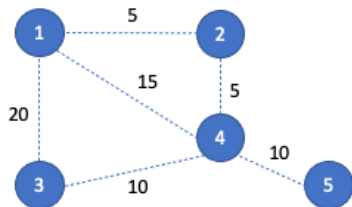
Not until sparse linear algebra was it practical

Leverage techniques developed for fast matrix multiplication

Adjacency Matrices

Given a graph $G = (V, E)$ with weights $w : E \rightarrow \mathbb{R}$, construct the **adjacency matrix**¹ $\mathbf{A} \in \mathbb{R}^{n \times n}$ where

$$a_{ij} \leftarrow \begin{cases} w(i,j) & \text{if } (i,j) \in E \\ 0 & \text{otherwise} \end{cases}$$



$$\mathbf{A} = \begin{bmatrix} 0 & 5 & 20 & 15 & 0 \\ 5 & 0 & 0 & 5 & 0 \\ 20 & 0 & 0 & 10 & 0 \\ 15 & 5 & 10 & 0 & 10 \\ 0 & 0 & 0 & 10 & 0 \end{bmatrix}$$

¹If G is undirected, store only upper-triangular part of \mathbf{A} .

Algebraic Structures

Monoids (\oplus) generalize *addition*.

$$(\text{"C3"}) + (\text{"ai"}) \rightarrow (\text{"C3.ai"})$$

Semirings (\oplus, \otimes) generalize *addition and multiplication*.

$$x_1 \cdot y_1 + \cdots + x_n \cdot y_n \rightarrow \min\{(x_1 + y_1), \dots, (x_n + y_n)\}$$

Monoids and **semirings** have several nice properties²:

- ▶ Identity elements
- ▶ Associativity
- ▶ ...

²However, neither structure enforces the existence of inverses.

Matrix Multiplication

Matrix-vector multiplication:

$$y \leftarrow \mathbf{A}x$$

Writing as an elementwise expression:

$$y_i \leftarrow \sum_k a_{ik} x_k$$

Generalizing to (\oplus, \otimes) :

$$y_i \leftarrow \bigoplus_k a_{ik} \otimes x_k$$

Matrix-matrix multiplication:

$$\mathbf{C} \leftarrow \mathbf{A}\mathbf{B}$$

Writing as an elementwise expression:

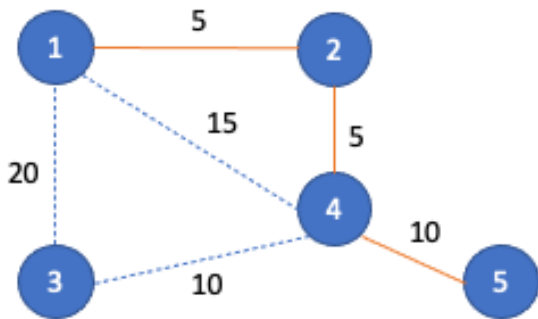
$$c_{ij} \leftarrow \sum_k a_{ik} b_{kj}$$

Generalizing to (\oplus, \otimes) :

$$c_{ij} \leftarrow \bigoplus_k a_{ik} \otimes b_{kj}$$

All Pairs Shortest Paths (APSP)

For each pair of vertices $i, j \in V$, compute the length of a shortest path from i to j .

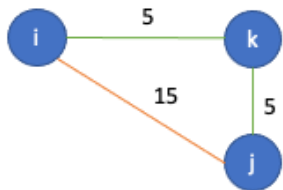


Sample application: **driving directions**, where road junctions are vertices and road distances are edge weights.

Bellman-Ford Algorithm

Fix a vertex i . Let $d_{ij}^{(\ell)}$ denote the length of a shortest path from i to another vertex j consisting of at most ℓ edges.

An edge (i, j) is **tense** if $d_{ij}^{(\ell)} > w(i, k) + d_{kj}^{(\ell)}$.



For $\ell \leftarrow 1$ to $n - 1$, **relax** all tense edges:

$$d_{ij}^{(\ell+1)} \leftarrow \min_k \{ w(i, k) + d_{kj}^{(\ell)} \}$$

A Linear Algebraic Algorithm

We store tentative shortest path distances in $\mathbf{D}^{(\ell)} \in \mathbb{R}^{n \times n}$ where

$$\mathbf{D}^{(\ell+1)} \leftarrow \mathbf{D}^{(\ell)} \oplus \mathbf{A}\mathbf{D}^{(\ell)} \quad \text{on } (\min, +)$$

Writing as an elementwise expression:

$$d_{ij}^{(\ell+1)} \leftarrow d_{ij}^{(\ell)} \oplus \left(\bigoplus_k a_{ik} \otimes d_{kj}^{(\ell)} \right)$$

Replacing the generic (\oplus, \otimes) with the **tropical semiring** $(\min, +)$:

$$d_{ij}^{(\ell+1)} \leftarrow \min\{d_{ij}^{(\ell)}, \min_k \{w(i, k) + d_{kj}^{(\ell)}\}\}$$

We may accelerate via matrix squaring $\mathbf{D}^{(\ell)} \leftarrow \mathbf{D}^{(\ell/2)} \oplus \mathbf{D}^{(\ell/2)} \mathbf{D}^{(\ell/2)}$.

Matrix Multiplication is Fast

Matrix multiplication has been studied from a variety of algorithmic lenses:

- ▶ Time complexity³
- ▶ Parallel and distributed
- ▶ With sparse data

³The current best time complexity is $O(n^{2.3728596})$ by Alman and Vassilevska via the laser method.

Strassen's Algorithm

$$\begin{bmatrix} \mathbf{C}_{11} & \mathbf{C}_{12} \\ \mathbf{C}_{21} & \mathbf{C}_{22} \end{bmatrix} \leftarrow \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix} \begin{bmatrix} \mathbf{B}_{11} & \mathbf{B}_{12} \\ \mathbf{B}_{21} & \mathbf{B}_{22} \end{bmatrix}$$

$$\mathbf{M}_1 \leftarrow (\mathbf{A}_{11} + \mathbf{A}_{22}) \cdot (\mathbf{B}_{11} + \mathbf{B}_{22})$$

$$\mathbf{M}_2 \leftarrow (\mathbf{A}_{21} + \mathbf{A}_{22}) \cdot \mathbf{B}_{11}$$

$$\mathbf{M}_3 \leftarrow \mathbf{A}_{11} \cdot (\mathbf{B}_{12} - \mathbf{B}_{22})$$

$$\mathbf{M}_4 \leftarrow \mathbf{A}_{22} \cdot (\mathbf{B}_{21} - \mathbf{B}_{11})$$

$$\mathbf{M}_5 \leftarrow (\mathbf{A}_{11} + \mathbf{A}_{12}) \cdot \mathbf{B}_{22}$$

$$\mathbf{M}_6 \leftarrow (\mathbf{A}_{21} - \mathbf{A}_{11}) \cdot (\mathbf{B}_{11} + \mathbf{B}_{12})$$

$$\mathbf{M}_7 \leftarrow (\mathbf{A}_{12} - \mathbf{A}_{22}) \cdot (\mathbf{B}_{21} + \mathbf{B}_{22})$$

$$\mathbf{C}_{11} \leftarrow \mathbf{M}_1 + \mathbf{M}_4 - \mathbf{M}_5 + \mathbf{M}_7$$

$$\mathbf{C}_{12} \leftarrow \mathbf{M}_3 + \mathbf{M}_5$$

$$\mathbf{C}_{21} \leftarrow \mathbf{M}_2 + \mathbf{M}_4$$

$$\mathbf{C}_{22} \leftarrow \mathbf{M}_1 - \mathbf{M}_2 + \mathbf{M}_3 + \mathbf{M}_6$$

The time complexity recurrence is given by

$$T(n) = 7T(n/2) + O(n^2) = O(n^{\log_2 7}) \approx O(n^{2.807})$$

Many Real-World Graphs are Large

Real-world graphs are **too large** to process on a single machine.

Graph	n	m
2002 crawl of the .uk domain	18.5M	261.8M
Friendster	65.6M	1.8B
Full USA road network	23.9M	28.9M
Orkut social network	3.1M	117M

Table: Publicly available large graphs

We may **distribute** computations across multiple machines³.

³TODO: It is cheaper to use multiple commodity machines than building specialized machines

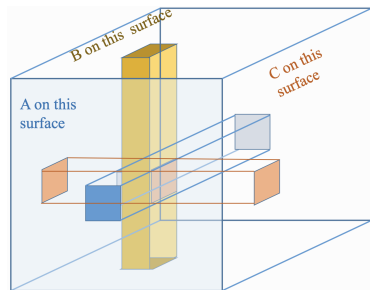
Parallel Matrix Multiplication

3D Matrix Multiplication was poorly introduced... introduce data layout in another slide (alex what does process owns ... mean?) or introduce parallel/distribute matrix algorithms? Would be nice to show trivial matrix-matrix distribution and why it's bad for communication

3D Matrix Multiplication

Organize processes into a 3D cube.

```
for i in range(n):  
    for j in range(n):  
        for k in range(n):  
            c[i,j] += a[i,k] * b[k,j]
```



Process (r, s, t) owns $\mathbf{A}^{(r,t)}$ and $\mathbf{B}^{(t,s)}$ and contributes $\mathbf{C}^{(r,s)}$
where $\mathbf{C}^{(r,s)} \leftarrow \mathbf{A}^{(r,t)} \mathbf{B}^{(t,s)}$.

Figure from CS 484: Parallel Programming.

See *Communication-Optimal Parallel 2.5D Matrix Multiplication and LU Factorization Algorithms* by Solomonik and Demmel (2011).

Many Real-World Graphs are Very Sparse

Many real-world graphs have **relatively few edges**.

Graph	n	m	sp
2002 crawl of the .uk domain	18.5M	261.8M	1.5e-6
Friendster	65.6M	1.8B	8.4e-9
Full USA road network	23.9M	28.9M	1.0e-7
Orkut social network	3.1M	117M	2.4e-5

Table: Publicly available large graphs

However, storing an adjacency matrix requires $O(n^2)$ **entries**.

Compressed Sparse Matrix Formats

We may store only the *nonzero* entries for **sparse** matrices.

Compressed sparse row⁴ (CSR) stores 3 lists:

- ▶ Values (length nnz)
- ▶ Column indices (length nnz): column of value
- ▶ Row offsets (length $n + 1$): number of values above this row

$$\begin{bmatrix} 0 & 1 & 0 & 2 \\ 0 & 0 & 3 & 0 \\ 4 & 0 & 0 & 0 \\ 5 & 6 & 0 & 0 \end{bmatrix}$$

Values: 1, 2, 3, 4, 5, 6

Column indices: 1, 3, 2, 0, 0, 1

Row offset: 0, 2, 3, 4, 6

⁴See also coordinate list and compressed sparse column (CSC) formats.

Compressed Sparse Matrix Algorithms

Sparse matrix-vector multiplication (SpMV):

$$y_i \leftarrow \sum_k a_{ik} x_k$$

```
for i in range(n):  
    for t in range(row_off[i], row_off[i+1]):  
        y[i] += val[t] * x[col_idx[t]]
```

Sparse matrix-matrix multiplication (SpMSpM):

- Communication-optimal algorithm is equivalent to partitioning a hypergraph⁵.

⁵*Hypergraph Partitioning for Sparse Matrix-Matrix Multiplication* by Ballard et al (2016).

Cyclops Tensor Framework

Matrix multiplications can be generalized to **tensors**⁶:

$$u_{i_1 r}^{(1)} \leftarrow \sum_{i_2 \dots i_d} t_{i_1 \dots i_d} u_{i_2 r}^{(2)} \cdots u_{i_d r}^{(d)}$$

The **Cyclops Tensor Framework**⁷ automatically parallelizes custom operations on tensors:

- ▶ Distributed memory
- ▶ Sparse tensors and operations
- ▶ Arbitrary semirings

Tensor algebra can be an interface for expressing *graph algorithms*.

⁶For example, MTTKRP is the main kernel for tensor decomposition.

⁷A *Massively Parallel Tensor Contraction Framework for Coupled-Cluster Computations* by Solomonik et. al (2014).

Hybrid Programming (OpenMP + MPI)

The Cyclops Tensor Framework uses **hybrid programming** to leverage parallelism at the *thread* and *process* level.

OpenMP

- ▶ Shared memory
- ▶ Fork, exec, wait model

```
#pragma omp parallel for  
for (i = 0; i < n; i++) {  
    // do something  
}
```

MPI⁸

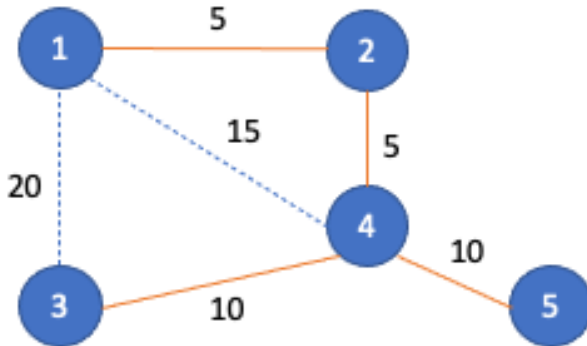
- ▶ Distributed memory
- ▶ Point-to-point and collectives
- ▶ Standard for HPC

```
// first node  
MPI_Send(&send_data, ...)  
  
// second node  
MPI_Recv(&recv_data, ...)
```

⁸In distributed deep learning, the *AllReduce* model often communicates via NCCL, which is modeled after MPI.

Minimum Spanning Tree (MST)

Find a spanning tree whose total weight is as small as possible.

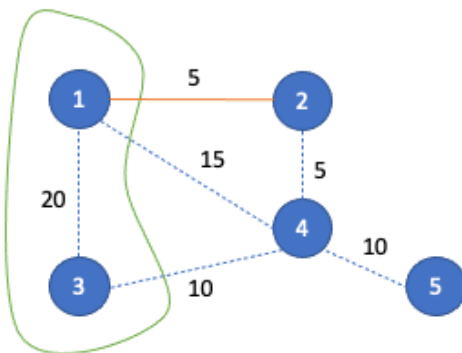


Sample application: a phone company laying cable in a neighborhood, where houses are vertices and cable lengths are edge weights.

Minimum Spanning Tree Algorithms

Most⁹ algorithms rely only on the **cut property** of MSTs.

Given a cut, the edge with the smallest weight that crosses the cut belongs to the mst.



⁹Karger et. al's $O(m)$ expected time algorithm also uses the **cycle property**.

Awerbuch-Shiloach Parallel MST Algorithm

TODO: improve?

Maintain a forest of **stars**¹⁰, where each vertex has a parent.

At the first iteration, the forest consists of n vertices, each with a self-loop.

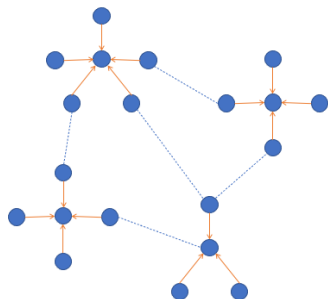
Repeats the following steps until convergence of the forest:

- ▶ Hooking: grows the MST by joining two stars together.
- ▶ Shortcutting: compresses trees into stars.

Can we rewrite this algorithm with **sparse linear algebra**?

¹⁰A **star** is a directed tree of height at most 1 and is intuitively a part of the minimum spanning tree.

Forest of Stars



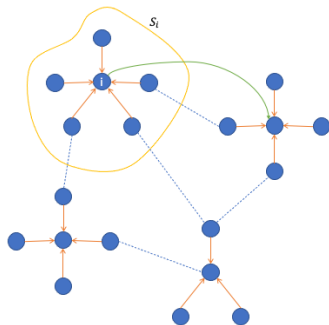
The **parent vector** $p \in \mathbb{R}^n$ stores the parent of each vertex.

The **parent matrix** $P \in \mathbb{R}^{n \times n}$ is defined

$$p_{ij} \leftarrow \begin{cases} 1 & \text{if } p_i = j \\ 0 & \text{otherwise} \end{cases}$$

Hooking

Each star attempts¹¹ to hook with its **smallest cut edge**.



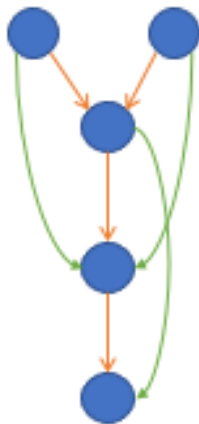
Hooking: joins two stars together.

$$p_i \leftarrow p \left[\underset{\substack{e \in \text{cut}(S_i, V \setminus S_i) \\ \text{edges that cross cut}}}{\text{argmin}} \quad w(e) \right]$$

¹¹Need to break ties to prevent cycles.

Shortcut

Hooking may generate trees of arbitrary height.



Shortcutting¹²: reduces the height of each tree in the forest by a factor of nearly two. $p_i \leftarrow p_{p_i}$

¹²Until completion so that all trees become stars.

Hooking as a Multilinear Operation

TODO: since we restrict trees in the forest to be stars, this is easy

f computes for each pair of nodes (i, j) , whether they share a parent.

$$f(p_i, a_{ij}, p_j) \leftarrow \begin{cases} a_{ij} : p_i \neq p_j \\ \infty : \text{otherwise} \end{cases}$$

Accumulating $f(p_i, a_{ij}, p_j)$ over j computes for each node i , its lightest edge that would leave its star.

$$q_i \leftarrow \bigoplus_j f(p_i, a_{ij}, p_j)$$

For each star root i , we **contract** its children's lightest “leaving” edges into itself.

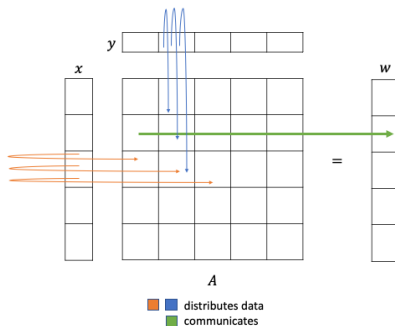
$$r_{p_j} \xleftarrow{\text{MINWEIGHT}} q_j$$

Data Distribution

How can we leverage the additional structure of $f(p_i, a_{ij}, p_j)$?

$$w_i \leftarrow \bigoplus_j f(x_i, a_{ij}, y_j)$$

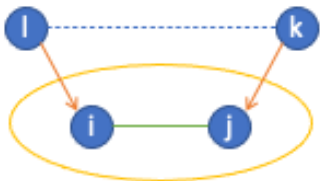
Process (r, s) owns $\mathbf{A}^{(r,s)}$ and needs $x^{(r)}$ and $y^{(s)}$. Process (r, s) contributes $w^{(r)}$ where $w^{(r)} \leftarrow f(x^{(r)}, \mathbf{A}^{(r,s)}, y^{(s)})$.



Contraction

To reduce the size of the graph in the next iteration, we may merge vertices via **contraction**.

$\mathbf{C} \leftarrow \mathbf{P}^T \mathbf{A} \mathbf{P}$ on the (\min, \cdot) semiring performs a **contraction** of \mathbf{A} .



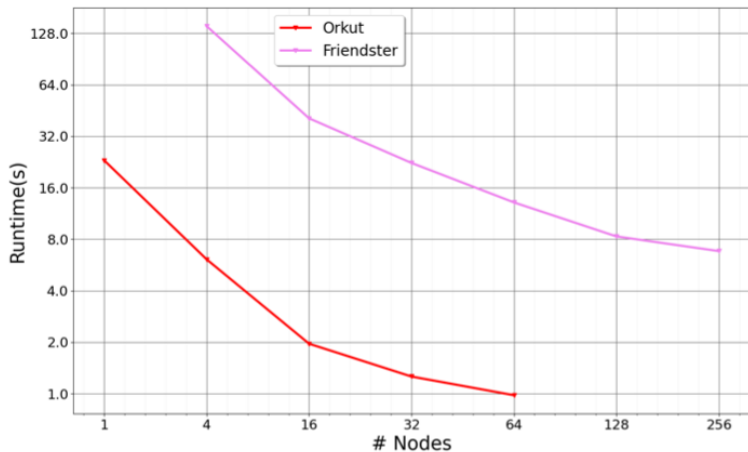
$$\begin{aligned} c_{ij} &\leftarrow \sum_l p_{li}^T \left(\sum_k a_{lk} p_{kj} \right) \\ &= \sum_l p_{li} \left(\sum_k a_{lk} p_{kj} \right) \\ &= \sum_{l,k} p_{li} a_{lk} p_{kj} \end{aligned}$$

$w(l, k)$ is projected onto $w(p_l, p_k)$ if:

- ▶ The edge (l, k) exists
- ▶ i is l 's parent ($p_{li} = 1$)
- ▶ j is k 's parent ($p_{kj} = 1$)

MST Strong Scaling Data

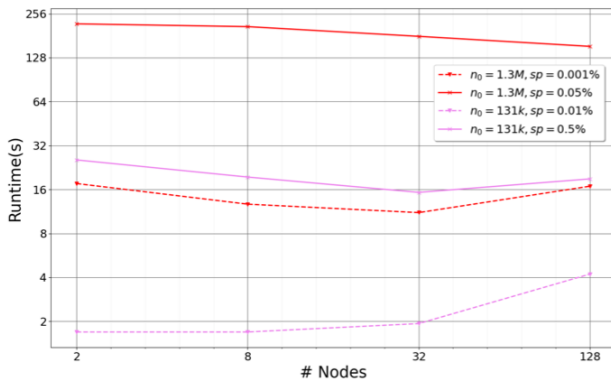
Strong scaling: fix the *total* problem size and vary the # of nodes.



(a) Social network graphs

MST Weak Scaling Data

Weak scaling: fix the problem size on *each* node and vary the # of nodes & total problem size accordingly.

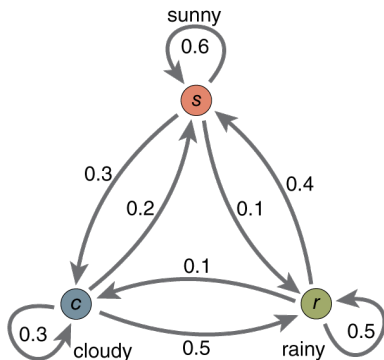


(a) Uniform random graph

Figure 2: Edge weak scaling of uniform random graphs.
Constant $n^2/p = n_0^2$ and edge percentage $f = 100 * m/n^2$

Markov Chains

Markov chain: Given a set of states S , if the current state is i then go to state j with probability p_{ij} .



$$P = \begin{bmatrix} 0.6 & 0.3 & 0.1 \\ 0.2 & 0.3 & 0.5 \\ 0.4 & 0.1 & 0.5 \end{bmatrix}$$

Random Walks

We prove expectations of the below via analyzing **random walks**:

- ▶ **Hitting time**: number of steps to reach state j from state i ?
- ▶ **Cover time**: number of steps to reach all states from state i ?
- ▶ **Stationary distribution**: what is the probability of ending in state i in a walk with t steps as $t \rightarrow \infty$?
- ▶ **Mixing time**: number of steps to approach the stationary distribution?

A vector $\pi = (\pi)_{i \in \mathcal{S}}$ is the **stationary distribution** if

$$\pi = \pi \mathbf{P}$$

The stationary distribution is an eigenvector of \mathbf{P} with eigenvalue 1.

Mixing Time

From a starting distribution $\pi^{(0)}$, take a random walk with t steps.

$$\begin{aligned}\pi^{(t)} &\leftarrow \pi^{(t-1)} \mathbf{P} \\ &\vdots \\ &= \pi^{(0)} \mathbf{P}^t\end{aligned}$$

Mixing time: what is the expected number of steps to approach the stationary distribution?

The above algorithm is the same as **power iteration**¹³:

- ▶ Converges to the dominant eigenvector of \mathbf{P} (namely, π)
- ▶ Convergence ratio is given by $|\lambda_1/\lambda_2|$

¹³Why is norm 1?

Expanders

A d -regular graph is an α -expander¹⁴ if $|\lambda_2| \leq \dots \leq |\lambda_n| \leq \alpha$.

Rapid mixing lemma: For an α -expander,

$$(\pi^{(t)} - \pi) \leq \alpha^t$$

Do good expanders ($\alpha = O(\frac{1}{\sqrt{d}})$) exist?

- ▶ Random graph (probabilistic method)
- ▶ Graph products (zig-zag product)
- ▶ Ramanujan graphs ($\alpha = \frac{1}{\sqrt{d}}$)

¹⁴ $\alpha < 1$. Combinatorial definition: A d -regular graph is an α -expander if

$$\frac{\partial S}{|S|} \geq \Omega(d(1 - \alpha)) \quad \forall S \subseteq V \quad (|S| \leq n/2)$$

Applications of Expanders

Expanders are a very powerful tool in modern TCS:

- ▶ Derandomization (**randomness** is a *resource*)
- ▶ Error correcting codes
- ▶ Sorting networks
- ▶ ...

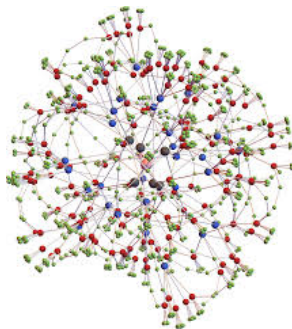


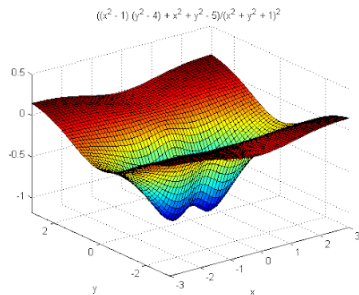
Figure from *Cutoff on All Ramanujan Graphs* by Lubetzky and Peres (2015).

Monte Carlo Algorithms

A **Monte Carlo algorithm** Π for a decision problem:

- ▶ Correctly outputs YES with probability $\geq 1/2$
- ▶ Correctly outputs NO with probability 1

Is the polynomial¹⁵ $P(x_1, \dots, x_k)$ over \mathbb{Z}_p **zero everywhere**?



P has at most degree d roots:

- ▶ Pick a random vector $x_1, \dots, x_k \in \mathbb{Z}_p$.
- ▶ $\Pr[P(x_1, \dots, x_k) = 0] \leq \frac{d}{p}$.
- ▶ If YES, repeat?
- ▶ If NO, output NO.

¹⁵For example, $((x^2 - 1)(y^2 - 4) + x^2 + y^2 - 5) / (x^2 + y^2 + 1)^2$.

Figure from *Visualizing Functions of Several Variables and Surfaces* by

Probability Amplification

Consider a Monte Carlo algorithm Π :

- ▶ Error probability $\leq 1/10$.
- ▶ m random bits.

How can we *lower* the error probability?

Repeat t times \rightarrow error probability $\leq (\frac{1}{10})^t$ and tm random bits.

Can we achieve a similar error probability with fewer random bits?

- ▶ Expanders!

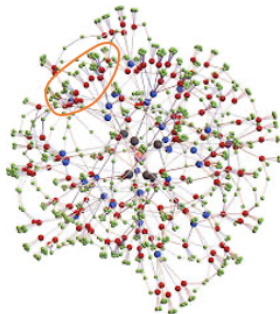
Derandomization

Construct a d -regular α -expander G with $n = 2^m$ vertices¹⁶.

Fix a subset $B \subseteq V$ with $|B| = \beta n$.

Take a random walk v_1, \dots, v_t from a random starting vertex.

$$\Pr[v_1 \in B \wedge \dots \wedge v_t \in B] \leq (\alpha + \beta)^t$$



¹⁶Efficient representations have been designed to construct G “on the fly”.

Derandomization (cont.)

If G is a good expander, we are *unlikely* to get stuck in B ¹⁷.

Probability amplification:

- ▶ Each vertex corresponds to a m bit string.
- ▶ When evaluating Π , we randomly pick a m bit string.

B is the set of m bit strings that output a wrong answer to Π .

$$|B| \leq \frac{1}{10}|n| \implies \beta \leq \frac{1}{10}$$

- ▶ Error probability $\leq (\alpha + \frac{1}{10})^t$.
- ▶ $m + t \log d$ random bits.

¹⁷Can also be viewed from the combinatorial definition of expanders.

Summary

Introduction

All Pairs Shortest Paths

Sparse Linear Algebra

Minimum Spanning Trees

Expanders

References

- ▶ *Graph Algorithms in the Language of Linear Algebra* by Kepner and Gilbert (2011)
- ▶ CS 554: Parallel Numerical Algorithms by Solomonik
- ▶ *New Connectivity and MSF Algorithms for Shuffle-Exchange Network and PRAM* by Awerbuch and Shiloach (1987)
- ▶ CS 574: Randomized Algorithms by Chan

Fin

Thanks. Questions?

Extra Slides

The following slides are *unfinished* but are cool related topics that I'd like to talk about at some point.

Ball

Given a vertex, compute its b closest neighbors.

Probabilistic Method

Prove existence with the **probabilistic method** then derandomize for a construction.

Dimensionality Reduction

Does there exist an embedding from ℓ_1 to Hamming space?

Take a random projection (ie multiply by a random matrix)

then use probabilistic method

Convolution as Matrix Multiplication

Convolution is a basic kernel for neural networks.

Can use toeplitz matrices to rewrite as matrix multiplication.

Also can generalize to tensors via matricization.

Tensor Decomposition

Tensor Completion