

# Barvinok's naive algorithm for dimensionality reduction and its use in Distance Geometry

Leo Liberti, CNRS LIX Ecole Polytechnique  
liberti@lix.polytechnique.fr

TAU18 Workshop 180424-26

*Joint work with Vu Khac Ky*



# Outline

**Barvinok's Naive Algorithm**

Distance Geometry

Noisy distances

Computational results

# Concentration of measure

From [Barvinok, 1997]

*The value of a “well behaved” function at a random point of a “big” probability space  $X$  is “very close” to the mean value of the function.*

and

*In a sense, measure concentration can be considered as an extension of the law of large numbers.*

# Concentration of measure

Given Lipschitz function  $f : X \rightarrow \mathbb{R}$  s.t.

$$\forall x, y \in X \quad |f(x) - f(y)| \leq L\|x - y\|_2$$

for some  $L \geq 0$ , there is *concentration of measure* if  $\exists$  constants  $c, C$  s.t.

$$\forall \varepsilon > 0 \quad \mathbb{P}_x(|f(x) - \mathbb{E}(f)| > \varepsilon) \leq c e^{-C\varepsilon^2/L^2}$$

where  $\mathbb{E}(\cdot)$  is w.r.t. given Borel measure  $\mu$  over  $X$

$\equiv$  “*discrepancy from mean is unlikely*”

# Barvinok's theorem

## Consider:

- ▶ for each  $k \leq m$ , manifolds  $\mathcal{X}_k = \{x \in \mathbb{R}^n \mid x^\top Q^k x = a_k\}$   
where  $m \leq \text{poly}(n)$
- ▶ feasibility problem  $F \equiv [\bigcap_{k \leq m} \mathcal{X}_k \stackrel{?}{\neq} \emptyset]$
- ▶ SDP relaxation  $\forall x \leq m (Q^k \bullet X = a_k) \wedge X \succeq 0$  with soln.  $\bar{X}$

## Find an approximate rank-1 solution of $F$

**Algorithm:**  $T \leftarrow \text{factor}(\bar{X}); \quad y \sim \mathcal{N}^n(0, 1); \quad x' \leftarrow Ty$

Then  $\exists c > 0, n_0 \in \mathbb{N}$  such that  $\forall n \geq n_0$

$$\text{Prob} \left( \forall k \leq m \quad \text{dist}(x', \mathcal{X}_k) \leq c \sqrt{\|\bar{X}\|_2 \ln n} \right) \geq 0.9.$$

**IDEA:** since  $x'$  is “close” to each  $\mathcal{X}_k$ , try local descent!

# Elements of Barvinok's formula

$$\text{Prob} \left( \forall k \leq m \quad \text{dist}(x', \mathcal{X}_k) \leq c \sqrt{\|\bar{X}\|_2 \ln n} \right) \geq 0.9.$$

- ▶  $\sqrt{\|\bar{X}\|_2}$  arises from  $T$  (a factor of  $\bar{X}$ )
- ▶  $\sqrt{\ln n}$  ensures concentration of measure
- ▶ 0.9 follows by adjusting parameter values in union bounds

# Outline

Barvinok's Naive Algorithm

**Distance Geometry**

Noisy distances

Computational results

# Distance Geometry Problem

- ▶ Given  $K \in \mathbb{N}_+$  and  $G = (V, E, d)$  with  $d : V \rightarrow \mathbb{Q}_+$ , determine if  $\exists$  realization  $x : V \rightarrow \mathbb{R}^K$  s.t.:

$$\forall \{i, j\} \in E \quad \|x_i - x_j\|_2^2 = d_{ij}^2$$

- ▶ **Inverse problem to:** given  $n$  pts in  $\mathbb{R}^K$ , determine some of their pairwise distances and their adjacencies
- ▶  $\Rightarrow$  “Draw graph in  $\mathbb{R}^K$ , where edges  $\equiv$  segments of corresp. length”
- ▶ **Applications:** clock synchronization protocols, sensor network localization, protein conformation, nanostructures, autonomous underwater vehicles, rigidity, statics and more
- ▶ *Strongly NP-hard even if  $K$  fixed and  $\text{ran}(d) = \{1, 2\}$*

See e.g. [L. et al, SIAM Review 2014], [Dokmanić et al., IEEE Sig. Proc. Mag. 2015], [L. and Lavor, Springer 2017]



# SDP formulation of the DGP

- ▶ A feasibility problem:

$$\begin{aligned} & \min F \bullet X \\ \forall \{i, j\} \in E & \quad X_{ii} + X_{jj} - 2X_{ij} = d_{ij}^2 \\ & \quad X \succeq 0 \end{aligned}$$

- ▶  $F$  depends on application; for protein conformation, try:

$$F = \sum_{\{i,j\} \in E} (X_{ii} + X_{jj} - 2X_{ij}) + \nu \text{Tr}(X)$$

where  $\nu$  is small (try 0.001)

- ▶ **Issue: obtain realization  $\bar{X}$  in  $\mathbb{R}^n$  not  $\mathbb{R}^K$**   
*need dimension reduction*

# Barvinok's alg. for the DGP

- ▶  $\forall \{i, j\} \in E \quad \mathcal{X}_{ij} = \{x \in \mathbb{R}^{nK} \mid \|x_i - x_j\|_2^2 = d_{ij}^2\}$
  - ▶ **DGP**  $\equiv$  “is  $\bigcap_{\{i,j\} \in E} \mathcal{X}_{ij}$  non-empty?”
  - ▶ **SDP rel.**  $X_{ii} + X_{jj} - 2X_{ij} = d_{ij}^2 \wedge X \succeq 0$  with soln.  $\bar{X}$
- 

- ▶ **Difference w.r.t. Barvinok:**  $x \in \mathbb{R}^{nK}$
- ▶ **IDEA:** sample  $y \sim \mathcal{N}^{nK}(0, \frac{1}{\sqrt{K}})$

- ▶ **Our result:** Extension of Barvinok's thm to rank  $K$

$$\text{Prob} \left( \forall k \leq m \quad \text{dist}(x', \mathcal{X}_k) \leq c \sqrt{\|\bar{X}\|_2 \ln(nK)} \right) \geq 0.9.$$

**Analysis improvement yields**

$$\text{Prob} \left( \forall k \leq m \quad \text{dist}(x', \mathcal{X}_k) \leq c \sqrt{\|\bar{X}\|_2 \ln n} \right) \geq 0.9.$$

# Proof structure

- ▶ **Show that, on average,  $\forall k \leq m (Ty)^\top Q^k(Ty) = Q^k \bullet \bar{X} = a_k$** 
  - ▶ compute multivariate integrals
  - ▶ bilinear terms disappear because  $y$  normally distributed
  - ▶ decompose multivariate int. to a sum of univariate int.
- ▶ **Exploit concentration of measure to show errors happen rarely**
  - ▶ a couple of technical lemmata yielding bounds
  - ▶  $\Rightarrow$  bound Gaussian measure  $\mu$  of  $\varepsilon$ -neighbourhoods of

$$A_i^- = \{y \in \mathbb{R}^{n \times K} \mid Q^i(Ty) \leq Q^i \bullet \bar{X}\}$$

$$A_i^+ = \{y \in \mathbb{R}^{n \times K} \mid Q^i(Ty) \geq Q^i \bullet \bar{X}\}$$

$$A_i = \{y \in \mathbb{R}^{n \times K} \mid Q^i(Ty) = Q^i \bullet \bar{X}\}.$$

- ▶ use union bound for measure of  $A_i^-(\varepsilon) \cap A_i^+(\varepsilon)$
- ▶ show  $A_i^-(\varepsilon) \cap A_i^+(\varepsilon) = A_i(\varepsilon)$
- ▶ use union bound for measure of intersections of  $A_i(\varepsilon)$
- ▶ appropriate values for some parameters  $\Rightarrow$  result

# The heuristic

1. Solve SDP relaxation of DGP, get soln.  $\bar{X}$
2. Barvinok's algorithm:  
 $T \leftarrow \text{factor}(\bar{X}), y \sim \mathcal{N}^{nK}(0, \frac{1}{\sqrt{K}}), x' \leftarrow Ty$
3. Use  $x'$  as starting point for a local NLP solver on formulation

$$\min_x \sum_{\{i,j\} \in E} (\|x_i - x_j\|_2^2 - d_{ij}^2)^2$$

and return improved solution  $x$

# Outline

Barvinok's Naive Algorithm

Distance Geometry

**Noisy distances**

Computational results

# The Interval DGP

- ▶ Distances never precise in applications
- ▶ DGP with constraints

$$\forall \{i, j\} \in E \quad L_{ij}^2 \leq \|x_i - x_j\|^2 \leq U_{ij}^2$$

- ▶ SDP formulation constraints

$$\forall \{i, j\} \in E \quad L_{ij}^2 \leq X_{ii} + X_{jj} - 2X_{ij} \leq U_{ij}^2$$

- ▶ NLP formulation

$$\min \sum_{\{i,j\} \in E} \left( (\max(L_{ij}^2 - \|x_i - x_j\|_2^2, 0))^2 + (\max(\|x_i - x_j\|_2^2 - U_{ij}^2, 0))^2 \right)$$

# Our result

- ▶ **Extension of Barvinok's thm to intervals**

$$\text{Prob} \left( \forall k \leq m \quad \text{dist}(x', \mathcal{X}_k) \leq c \sqrt{\|\bar{X}\|_2 \ln m} \right) \geq 0.9.$$

- ▶ *Unfortunately, a worse bound: for most applications,  $\ln m \geq \ln n$*

# Outline

Barvinok's Naive Algorithm

Distance Geometry

Noisy distances

**Computational results**



# Evaluation metrics

- ▶  $\text{LDE}(x) = \max_{\{i,j\} \in E} | \|x_i - x_j\| - d_{ij} |$
- ▶  $\text{MDE}(x) = \frac{1}{|E|} \sum_{\{i,j\} \in E} | \|x_i - x_j\| - d_{ij} |$
- ▶ CPU time
- ▶ **For intervals:** replace  $\|x_i - x_j\| - d_{ij}$  by

$$\max(L_{ij} - \|x_i - x_j\|_2, 0) + \max(\|x_i - x_j\|_2 - U_{ij}, 0)$$

# Comparison with PCA

- ▶ Barvinok's heuristic: **SDP + Barvinok + NLP**
- ▶ Best-known dimensionality reduction algorithm:  
*Principal Component Analysis* (PCA)
  1. Decompose SDP soln.  $\bar{X}$  into  $P^\top \text{diag}(\lambda)P$
  2.  $\forall k > K$  let  $\lambda_k \leftarrow 0$
  3.  $x' \leftarrow P^\top \text{diag}(\lambda)P$
- ▶ Compare against heuristic: **SDP + PCA + NLP**

# Results on DGP

<i>instance</i>	<i>MDE</i>		<i>LDE</i>		<i>CPU</i>	
	barvinok	pca	barvinok	pca	barvinok	pca
names	<b>0.00</b>	0.11	<b>0.07</b>	<b>1.00</b>	39.33	<b>22.44</b>
pept	<b>0.00</b>	0.10	<b>0.03</b>	1.81	83.91	<b>56.65</b>
C0020pdb	<b>0.00</b>	0.12	<b>0.01</b>	2.72	76.73	<b>49.39</b>
1guu-1	0.03	<b>0.00</b>	0.26	<b>0.08</b>	370.73	<b>322.66</b>
1guu-4000	<b>0.03</b>	0.12	<b>0.73</b>	1.15	415.66	<b>397.87</b>
1guu	0.02	<b>0.01</b>	<b>0.29</b>	0.33	305.54	<b>267.52</b>
res_5000	<b>0.00</b>	0.15	<b>0.00</b>	2.24	84.19	<b>53.36</b>
res_2000	<b>0.00</b>	0.07	<b>0.00</b>	1.46	85.26	<b>53.39</b>
res_0	<b>0.00</b>	<b>0.00</b>	<b>0.00</b>	0.01	93.08	<b>62.64</b>
res_3000	<b>0.00</b>	0.01	<b>0.00</b>	1.08	88.51	<b>53.43</b>
res_1000	<b>0.00</b>	0.10	<b>0.00</b>	3.05	87.88	<b>52.98</b>
res_2kxa	<b>0.00</b>	0.15	<b>0.00</b>	2.92	764.34	<b>713.35</b>
C0030pk1	<b>0.00</b>	0.11	<b>0.07</b>	2.19	1178.73	<b>1024.86</b>

*Most of the CPU time taken by the local NLP solver `scipy.optimize.root`*

# Results on iDGP

<i>instance</i>	<i>MDE</i>		<i>LDE</i>		<i>CPU</i>	
	barvinok	pca	barvinok	pca	barvinok	pca
names	<b>0.04</b>	<b>0.00</b>	2.11	<b>0.07</b>	53.91	<b>36.86</b>
pept	<b>0.01</b>	<b>0.00</b>	0.46	<b>0.40</b>	133.28	<b>99.60</b>
C0020pdb	<b>0.02</b>	<b>0.00</b>	1.64	<b>0.42</b>	112.38	<b>79.22</b>
1guu-1	<b>0.02</b>	<b>0.01</b>	1.09	<b>0.58</b>	500.64	<b>440.50</b>
1guu-4000	<b>0.03</b>	<b>0.02</b>	1.49	1.49	522.19	<b>461.53</b>
1guu	<i>memory overflow</i>					
res_5000	<b>0.01</b>	<b>0.00</b>	<b>0.69</b>	<b>0.08</b>	30764.21	<b>30465.16</b>
res_2000	<b>0.01</b>	<b>0.00</b>	1.78	<b>0.10</b>	33017.88	<b>32713.91</b>
res_0	<b>0.00</b>	<b>0.00</b>	<b>0.11</b>	<b>0.11</b>	22897.14	<b>22619.79</b>
res_3000	<b>0.00</b>	<b>0.00</b>	<b>0.05</b>	<b>0.08</b>	26095.91	<b>25846.81</b>
res_1000	<b>0.00</b>	<b>0.00</b>	<b>0.05</b>	<b>0.07</b>	27790.87	<b>27542.96</b>
res_2kxa	<i>memory overflow</i>					
C0030pk1	<i>memory overflow</i>					

*Most of the CPU time taken by the local NLP solver `scipy.optimize.root`*

# Future directions

- ▶ What kind of local descent algorithm would be most appropriate?
- ▶ Improve performance on iDGP  
*change SDP/NLP formulations*
- ▶ Test this technique on a wider range of SDPs

[L., Ky Vu, *Barvinok's Naive Algorithm in Distance Geometry*,  
Op. Res. Lett., in revision]