

DISTANCE SHRINKAGE AND EUCLIDEAN EMBEDDING

Ming Yuan

Morgridge Institute for Research

and

Department of Statistics

University of Wisconsin-Madison

myuan@stat.wisc.edu

<http://www.stat.wisc.edu/~myuan>



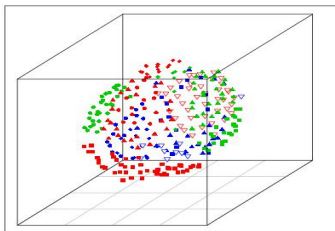
(Joint work with Luwan Zhang and Grace Wahba)

VPU SEQUENCE VARIATION

```

1  HQSLQLVALVVAIIIAI[VW]VVALEYRKL[RQK]IDRI[IR]RAEDSGMSEGGQEELSGLVENGHHPWDIDL  9
2  HQSLQLVALVVAIIIAI[VW]VVALEYRKL[RQK]IDRI[IR]RAEDSGMSEGGQEELSGLVENGHHPWDIDL  2
3  HQSLQLVALVVAIIIAI[VW]VVALEYRKL[RQK]IDRI[IR]RAEDSGMSEGGQEELSGLVENGHHPWDIDL  2
4  HQSLQLVALVVAIIIAI[VW]VVALEYRKL[RQK]IDRI[IR]RAEDSGMSEGGQEELSGLVENGHHPWDIDL  2
5  HQSLQLVALVVAIIIAI[VW]VVALEYRKL[RQK]IDRI[IR]RAEDSGMSEGGQEELSGLVENGHHPWDIDL  2
6  HQSLQLVALVVAIIIAI[VW]VVALEYRKL[RQK]IDRI[IR]RAEDSGMSEGGQEELSGLVENGHHPWDIDL  2
7  HQSLQLVALVVAIIIAI[VW]VVALEYRKL[RQK]IDRI[IR]RAEDSGMSEGGQEELSGLVENGHHPWDIDL  3
8  HQSLQLVALVVAIIIAI[VW]VVALEYRKL[RQK]IDRI[IR]RAEDSGMSEGGQEELSGLVENGHHPWDIDL  2
9  HQSLQLVALVVAIIIAI[VW]VVALEYRKL[RQK]IDRI[IR]RAEDSGMSEGGQEELSGLVENGHHPWDIDL  3
10 HQSLQLVALVVAIIIAI[VW]VVALEYRKL[RQK]IDRI[IR]RAEDSGMSEGGQEELSGLVENGHHPWDIDL  4
11 HQSLQLVALVVAIIIAI[VW]VVALEYRKL[RQK]IDRI[IR]RAEDSGMSEGGQEELSGLVENGHHPWDIDL  4
12 HQSLQLVALVVAIIIAI[VW]VVALEYRKL[RQK]IDRI[IR]RAEDSGMSEGGQEELSGLVENGHHPWDIDL  3
13 HQSLQLVALVVAIIIAI[VW]VVALEYRKL[RQK]IDRI[IR]RAEDSGMSEGGQEELSGLVENGHHPWDIDL  4
14 HQSLQLVALVVAIIIAI[VW]VVALEYRKL[RQK]IDRI[IR]RAEDSGMSEGGQEELSGLVENGHHPWDIDL  4
15 HQSLQLVALVVAIIIAI[VW]VVALEYRKL[RQK]IDRI[IR]RAEDSGMSEGGQEELSGLVENGHHPWDIDL  4
16 HQSLQLVALVVAIIIAI[VW]VVALEYRKL[RQK]IDRI[IR]RAEDSGMSEGGQEELSGLVENGHHPWDIDL  4
17 HQSLQLVALVVAIIIAI[VW]VVALEYRKL[RQK]IDRI[IR]RAEDSGMSEGGQEELSGLVENGHHPWDIDL  5
18 HQSLQLVALVVAIIIAI[VW]VVALEYRKL[RQK]IDRI[IR]RAEDSGMSEGGQEELSGLVENGHHPWDIDL  4
19 HQSLQLVALVVAIIIAI[VW]VVALEYRKL[RQK]IDRI[IR]RAEDSGMSEGGQEELSGLVENGHHPWDIDL  9
20 HQSLQLVALVVAIIIAI[VW]VVALEYRKL[RQK]IDRI[IR]RAEDSGMSEGGQEELSGLVENGHHPWDIDL  9
21 HQSLQLVALVVAIIIAI[VW]VVALEYRKL[RQK]IDRI[IR]RAEDSGMSEGGQEELSGLVENGHHPWDIDL  9
22 HQSLQLVALVVAIIIAI[VW]VVALEYRKL[RQK]IDRI[IR]RAEDSGMSEGGQEELSGLVENGHHPWDIDL  10
23 HQSLQLVALVVAIIIAI[VW]VVALEYRKL[RQK]IDRI[IR]RAEDSGMSEGGQEELSGLVENGHHPWDIDL  10
24 HQSLQLVALVVAIIIAI[VW]VVALEYRKL[RQK]IDRI[IR]RAEDSGMSEGGQEELSGLVENGHHPWDIDL  10
25 HQSLQLVALVVAIIIAI[VW]VVALEYRKL[RQK]IDRI[IR]RAEDSGMSEGGQEELSGLVENGHHPWDIDL  10
26 HQSLQLVALVVAIIIAI[VW]VVALEYRKL[RQK]IDRI[IR]RAEDSGMSEGGQEELSGLVENGHHPWDIDL  10
27 HQSLQLVALVVAIIIAI[VW]VVALEYRKL[RQK]IDRI[IR]RAEDSGMSEGGQEELSGLVENGHHPWDIDL  10
28 HQSLQLVALVVAIIIAI[VW]VVALEYRKL[RQK]IDRI[IR]RAEDSGMSEGGQEELSGLVENGHHPWDIDL  9
29 HQSLQLVALVVAIIIAI[VW]VVALEYRKL[RQK]IDRI[IR]RAEDSGMSEGGQEELSGLVENGHHPWDIDL  10
30 HQSLQLVALVVAIIIAI[VW]VVALEYRKL[RQK]IDRI[IR]RAEDSGMSEGGQEELSGLVENGHHPWDIDL  9
31 HQSLQLVALVVAIIIAI[VW]VVALEYRKL[RQK]IDRI[IR]RAEDSGMSEGGQEELSGLVENGHHPWDIDL  7
32 HQSLQLVALVVAIIIAI[VW]VVALEYRKL[RQK]IDRI[IR]RAEDSGMSEGGQEELSGLVENGHHPWDIDL  10
33 HQSLQLVALVVAIIIAI[VW]VVALEYRKL[RQK]IDRI[IR]RAEDSGMSEGGQEELSGLVENGHHPWDIDL  10
34 HQSLQLVALVVAIIIAI[VW]VVALEYRKL[RQK]IDRI[IR]RAEDSGMSEGGQEELSGLVENGHHPWDIDL  10
35 HQSLQLVALVVAIIIAI[VW]VVALEYRKL[RQK]IDRI[IR]RAEDSGMSEGGQEELSGLVENGHHPWDIDL  10
36 HQSLQLVALVVAIIIAI[VW]VVALEYRKL[RQK]IDRI[IR]RAEDSGMSEGGQEELSGLVENGHHPWDIDL  10
37 HQSLQLVALVVAIIIAI[VW]VVALEYRKL[RQK]IDRI[IR]RAEDSGMSEGGQEELSGLVENGHHPWDIDL  10
38 HQSLQLVALVVAIIIAI[VW]VVALEYRKL[RQK]IDRI[IR]RAEDSGMSEGGQEELSGLVENGHHPWDIDL  10
39 HQSLQLVALVVAIIIAI[VW]VVALEYRKL[RQK]IDRI[IR]RAEDSGMSEGGQEELSGLVENGHHPWDIDL  10
40 HQSLQLVALVVAIIIAI[VW]VVALEYRKL[RQK]IDRI[IR]RAEDSGMSEGGQEELSGLVENGHHPWDIDL  10
41 HQSLQLVALVVAIIIAI[VW]VVALEYRKL[RQK]IDRI[IR]RAEDSGMSEGGQEELSGLVENGHHPWDIDL  10
42 HQSLQLVALVVAIIIAI[VW]VVALEYRKL[RQK]IDRI[IR]RAEDSGMSEGGQEELSGLVENGHHPWDIDL  10
43 HQSLQLVALVVAIIIAI[VW]VVALEYRKL[RQK]IDRI[IR]RAEDSGMSEGGQEELSGLVENGHHPWDIDL  10
44 HQSLQLVALVVAIIIAI[VW]VVALEYRKL[RQK]IDRI[IR]RAEDSGMSEGGQEELSGLVENGHHPWDIDL  10

```



[Pickering et al., 2014]

Multidimensional Scaling

To what extent, does multidimensional scaling (MDS) faithfully reflect features in the original data?

FROM DISSIMILARITY TO DISTANCE

- ▶ A set of objects $\{O_1, \dots, O_n\}$ from an arbitrary domain \mathcal{O}
- ▶ Observe pairwise dissimilarity scores x_{ij} s

$$x_{ij} \approx \text{dist}(O_i, O_j), \quad (i, j) \in \Omega$$

- ▶ “Closest” Euclidean embedding – $p_1, \dots, p_n \in \mathbb{R}^{n-1}$

$$\text{dist}(O_i, O_j) = \|p_i - p_j\|^2, \quad 1 \leq i < j \leq n$$

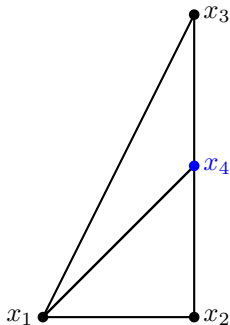
- ▶ Other applications – protein folding, chromosome conformation capture, graph drawing, ...

WHAT IS AN EDM

$$X = (x_{ij})_{1 \leq i, j \leq n} \implies D = (d_{ij} = \|p_i - p_j\|^2)_{1 \leq i, j \leq n} \in \mathcal{D}_n$$

- ▶ Nonnegativity – $d_{ij} \geq 0$
- ▶ Identity – $d_{ii} = 0$
- ▶ Symmetry – $d_{ij} = d_{ji}$
- ▶ Triangle inequality – $\sqrt{d_{ij}} + \sqrt{d_{jk}} \geq \sqrt{d_{ik}}$
- ▶ More than a metric

$$\begin{pmatrix} 0 & 1 & 5 & d_{14} \\ 1 & 0 & 4 & 1 \\ 5 & 4 & 0 & 1 \\ d_{14} & 1 & 1 & 0 \end{pmatrix} \implies (\sqrt{5} - 1)^2 \leq d_{14} \leq 4$$



GEOMETRY OF EDM

A symmetric matrix $D \in \mathbb{R}^{n \times n}$ is an Euclidean distance matrix iff

- ▶ It is hollow – $d_{ii} = 0$
- ▶ It is conditionally negative semi-definite on

$$\mathcal{X}_n = \{x \in \mathbb{R}^n : x^\top \mathbf{1} = 0\}$$

Embedding can be identified with the eigenstructure of [Schönberg transform](#) of D

$$\mathcal{R}(D) = -\frac{1}{2}J D J \quad \text{where } J = I - \mathbf{1}\mathbf{1}^\top/n$$

Important consequence:

- ▶ The set of $n \times n$ EDMs is a [convex cone](#) without interior

ESTIMATING AN EDM

$$x_{ij} = d_{ij} + \varepsilon_{ij} \implies D$$

- ▶ Projection type of estimate – \mathcal{D}_n is a closed convex hull

$$\mathcal{P}_{\mathcal{D}_n}(X) = \arg \min_{M \in \mathcal{D}_n} \|X - M\|_{\mathbb{F}}^2.$$

- ▶ Not working – at most $n(n-1)/2$ observations with $n(n-1)/2$ parameters
- ▶ Accounting for low embedding dimension – MDS

$$\mathcal{P}_{\mathcal{D}_n(r)}(X) = \arg \min_{M \in \mathcal{D}_n(r)} \|X - M\|_{\mathbb{F}}^2.$$

- Computationally challenging
- Statistically unstable

REGULARIZED KERNEL ESTIMATION

- ▶ From EDM to kernel $K = (k_{ij})$:

$$d_{ij} = \|p_i - p_j\|^2 = p_i^\top p_i + p_j^\top p_j - 2p_i^\top p_j =: k_{ii} + k_{jj} - 2k_{ij}$$

- D is an EDM iff $K \succeq 0$

- ▶ Regularized kernel estimate

$$\hat{K} = \arg \min_{M \succeq 0} \left\{ \sum_{(i,j) \in \Omega} \left(x_{ij} - \underbrace{\langle M, (e_i - e_j)(e_i - e_j)^\top \rangle}_{m_{ii} + m_{jj} - 2m_{ij}} \right)^2 + \lambda_n \text{trace}(M) \right\}$$

- ▶ Back to distance matrix

$$\hat{d}_{ij} = \hat{k}_{ii} + \hat{k}_{jj} - 2\hat{k}_{ij}$$

WHY DOES IT WORK?

- ▶ Kernel is **not** estimable from distance data

$\mathcal{T} : \mathcal{S}_n \rightarrow \mathcal{D}_n \quad M \mapsto (m_{ii} + m_{jj} - 2m_{ij})_{1 \leq i, j \leq n}$ is not injective

- ▶ What are we estimating – **Minimum trace kernel**

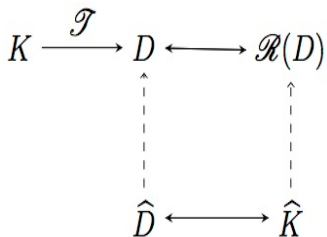
- the preimage $\mathcal{M}(D)$ of any $D \in \mathcal{D}_n$ under \mathcal{T} is **convex**
- there is a **unique** minimum trace kernel in $\mathcal{M}(D)$

$$\mathcal{R}(D) = \arg \min_{M \in \mathcal{M}(D)} \text{trace}(M)$$

- $\mathcal{R}(\cdot)$ is the **Schönberg transform**

$$\mathcal{R}(D) = -\frac{1}{2}JDJ$$

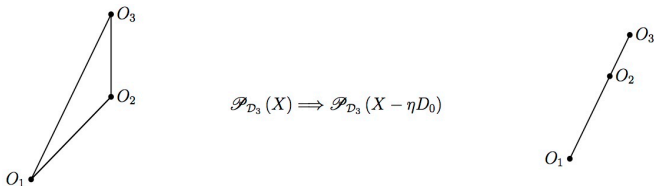
DISTANCE SHRINKAGE



$$\hat{D} = \mathcal{J}(\hat{K}) = \mathcal{P}_{\mathcal{D}_n} \left(X - \frac{\lambda_n}{2n} \mathbb{1}\mathbb{1}^\top \right)$$

$$\hat{D} = \mathcal{J}(\hat{K}) \text{ and } \hat{K} = \mathcal{R}(\hat{D})$$

EFFECT OF DISTANCE SHRINKAGE



- Embedding dim is one if

$$\frac{1}{3}(x_{12} + x_{13} + x_{23}) - \frac{\Delta_x}{3} \leq \eta < \frac{1}{3}(x_{12} + x_{13} + x_{23}) + \frac{2\Delta_x}{3},$$

where

$$\Delta_x := \sqrt{2[(x_{12} - x_{13})^2 + (x_{12} - x_{23})^2 + (x_{13} - x_{23})^2]}$$

HOW WELL DOES IT WORK?

Let $\hat{D} = \mathcal{P}_{\mathcal{D}_n} \left(X - \frac{\lambda}{2n} \mathbb{1}\mathbb{1}^\top \right)$. For any $\lambda \geq 2\|X - D\|$,

$$\|\hat{D} - D\|_F^2 \leq \inf_{M \in \mathcal{D}_n} \left\{ \|M - D\|_F^2 + \frac{9}{4} \lambda^2 (\dim(M) + 1) \right\}.$$

- ▶ If $\dim(D) = r$, then

$$\|\hat{D} - D\|_F^2 \lesssim r \|X - D\|^2$$

- ▶ For sub-Gaussian errors – $\lambda \sim \sqrt{n}$

$$\|\hat{D} - D\|_F^2 \lesssim_p rn \iff \text{Minimax optimal}$$

FIX DIMENSION EMBEDDING

- ▶ Eigenvalue decomposition of $\mathcal{R}(\hat{D})$
- ▶ Keep the leading r eigenvalues – W_r
- ▶ Getting back to Euclidean distance matrix $\hat{D}_r = \mathcal{I}(W_r)$

Let $D_r = \mathcal{P}_{\mathcal{D}_n(r)} D$. For any $\lambda \geq 2\|X - D\|$,

$$\|J(\hat{D}_r - D_r)J\|_{\text{F}}^2 \leq C \left(\min_{M \in \mathcal{D}_n(r)} \|J(D - M)J\|_{\text{F}}^2 + \lambda^2 r \right)$$

PROJECTION TO EDM

Recall from Schönberg (1935),

$$\mathcal{D}_n = S_1 \cap S_2$$

where

$$S_1 = \{M \in \mathbb{R}^{n \times n} : JMJ \preceq 0\},$$

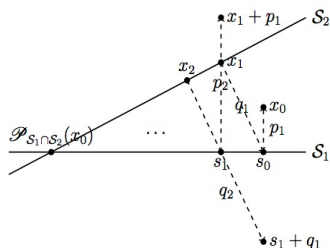
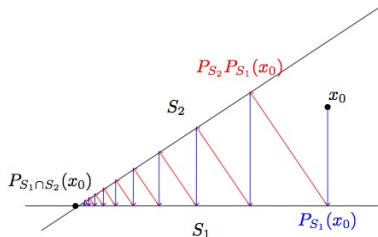
and

$$S_2 = \{M \in \mathbb{R}^{n \times n} : \text{diag}(M) = \mathbf{0}\}.$$

Both P_{S_1} and P_{S_2} have analytic forms amenable to computation.

[Glunt et al. (1990)]

ALTERNATING PROJECTION METHODS



- ▶ von Neumann's method of alternating projections (von Neumann, 1933)

$$S_1, S_2 \text{ closed} \implies \lim_{n \rightarrow \infty} (P_{S_1} P_{S_2})^n x_0 = P_{S_1 \cap S_2} x_0$$

- ▶ Dykstra's algorithm (1986)

$$x_n^0 := x_{n-1}^2, \quad x_n^i := P_{S_i}(x_n^{i-1} - y_{n-1}^i), \quad y_n^i = x_n^i - (x_n^{i-1} - y_{n-1}^i)$$

DEALING WITH MISSING OBSERVATIONS

- ▶ Observe

$$x_{ij} = d_{ij} + \varepsilon_{ij}, \quad (i, j) \in \Omega \subset \{(i, j) : 1 \leq i < j \leq n\}$$

- ▶ Goal

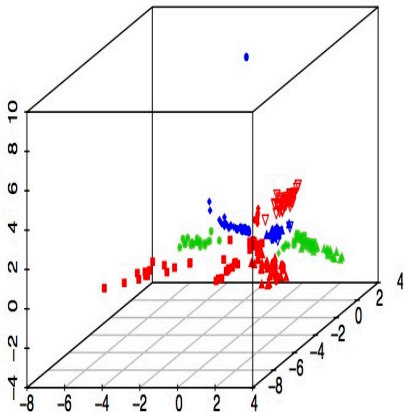
$$\min_{M \in \mathcal{D}_n} \sum_{(i, j) \in \Omega} [(x_{ij} - \eta) - m_{ij}]^2$$

- ▶ EM Algorithm

- Initialization x_{ij} for $(i, j) \in \Omega^c$
- M Step - $D^{(t+1)} \approx \mathcal{P}_{\mathcal{D}_n}(X^{(t)} - \eta D_0)$
- E Step - $x_{ij}^{(t+1)} = d_{ij}^{(t+1)}$ for $(i, j) \in \Omega^c$

- ▶ Cross-validation to choose η

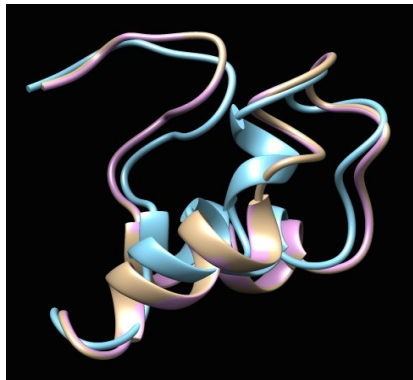
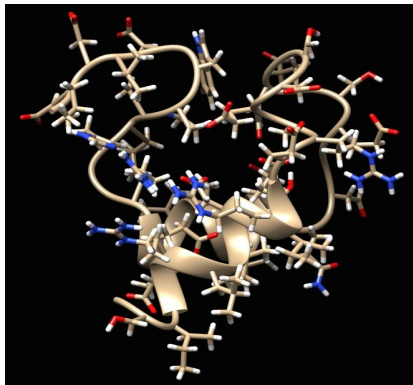
VPU SEQUENCE VARIATION



- ▶ 304 Vpu Sequences from 14 HIV-1 infected individuals
 - 5 long term non-progressors
 - 4 normal progressors
 - 5 Rapid progressors

PROTEIN SECONDARY STRUCTURE

- ▶ Coordinates of 671 atoms taken from PDB (ID: 2K7Y)
- ▶ Simulate pairwise distances with noise



SUMMARY

- ▶ The problem of reconstructing EDM from pairwise dissimilarity scores arises naturally in many applications
- ▶ Motivated in particular by biological problems:
 - Notion of distance between genomic sequences
 - Molecular structure determination
 - Chromosome conformation
- ▶ Distance shrinkage
 - Encourages low dimensional embedding
 - Leads to improved estimation risk
 - Efficient to compute
- ▶ Looking ahead – clustering, tree, ...