# Machine Learning
## Section 13: ISOMAP, LLE, MVU, MDS

Stefan Harmeling

WS 2021/22

# What we have seen so far?

Sections:

# Machine Learning
## Section 13: ISOMAP, LLE, MVU, MDS

Stefan Harmeling

WS 2021/22

### Definition 13.1 (Dimensionality reduction)

*Given n data points $X = [x_1, \ldots, x_n] \in \mathbb{R}^{D \times n}$. Find a low dimensional representation $Z = [z_1, \ldots, z_n] \in \mathbb{R}^{d \times n}$ with $d \ll D$ that keeps most of the properties of the higher dimensional data.*

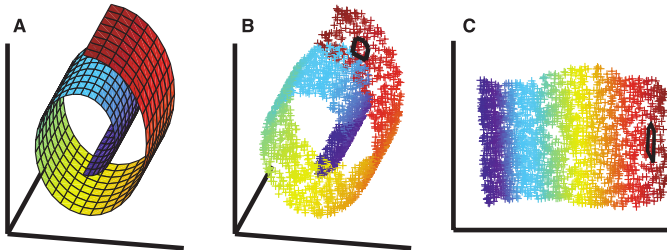What properties should be kept?

- variance (PCA, linear methods)
- neighborhood relationships (ISOMAP, LLE, MVU, nonlinear methods)

## Linear dimensionality reduction

Data points are assumed to lie on an linear subspace, i.e. on a line, plane, hyperplane or some $\mathbb{R}^d$ somewhere inside the $\mathbb{R}^D$ for $d < D$.

## Nonlinear dimensionality reduction

Data points are assumed to lie on an curved subspace, i.e. on a curve or more generally a *manifold* inside the $\mathbb{R}^D$. Simply put, a *manifold* is a subset of $\mathbb{R}^D$ that locally looks like the $\mathbb{R}^d$ with $d < D$.
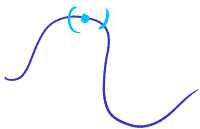


from Roweis/Saul's paper on LLE

Manifold:



Locally like $\mathbb{R}^2$

"Manifold hypothesis": Our data lies on a manifold

# Nonlinear dimensionality reduction

from http://en.wikipedia.org/wiki/Nonlinear_dimensionality_reduction

- Sammon's mapping, 1969
- Self-organizing map (SOM, aka Kohonen map, based on neural networks), 1982
- Principal curves and manifolds, 1984
- Autoencoders (some neural networks), 19xx
- Generative topographic map (GTM, probabilistic version of SOM), 1996
- Curvilinear component/distance analysis (CCA, CDA), 1997
- Kernel PCA (kPCA), 1998
- ISOMAP, 2000
- Locally-linear embedding (LLE), 2000
- Laplacian Eigenmaps, 2001
- Hessian LLE, 2003
- Gaussian process latent variable models (GPLVM), 2004
- Maximum variance unfolding (MVA, aka semidefinite embedding), 2004
- Relational perspective map, 2004
- Nonlinear PCA (based on neural networks), 2005
- Local tangent space alignment (LTSA), 2005
- Modified LLE, 2006
- Diffusion maps, 2006
- Local multidimensional scaling, 2006
- Manifold alignment, 2008
- Manifold sculpting, 2008
- Diffeomorphic Dimensionality Reduction (Diffeomap), 2009
- Rank visu, 2009
- Topologically Constrained Isometric Embedding (TCIE), 2010
- t-distributed stochastic neighbor embedding (t-SNE), 2008

# Two important papers both in a single issue of Science

ISOMAP

LLE

## Success of ISOMAP and LLE (checked Dec 2021)

- ISOMAP paper 15000 citations
  (Tenenbaum 75000, De Silva ?, Langford 41000)
- LLE paper 16000 citations
  (Roweis 38000, Saul 41000)
- generated a lot of follow-up papers
- published in Science journal

(numbers from `http://scholar.google.de` on 2021-29-11)

## Why are ISOMAP and LLE so successful?

- easy to use
- easy to understand
- code freely available
- almost no parameter tweaking
- inspiring examples
- published in Science journal

*Do this with your work!*

# ISOMAP: manifold of faces (synthetic)



Up-down pose

Left-right pose

Lighting direction

from Tenenbaum, da Silva, Langford's paper on ISOMAP

# ISOMAP: manifold of digits



from Tenenbaum, da Silva, Langford's paper on ISOMAP

# LLE: manifold of faces (real)



**Fig. 3.** Images of faces (*11*) mapped into the embedding space described by the first two coordinates of LLE. Representative faces are shown next to circled points in different parts of the space. The bottom images correspond to points along the top-right path (linked by solid line), illustrating one particular mode of variability in pose and expression.

from Roweis/Saul's paper on LLE

Both methods (ISOMAP and LLE) capture the manifold by defining a...

## Neighborhood graph, aka proximity graph

Given a $D \times n$ data matrix $X = [x_1, \ldots, x_n] \in \mathbb{R}^{D \times n}$.

- vertices $x_1, \ldots, x_n$.
- edges between neighbors, i.e. close-by points

## Examples

*[handwritten: choose k]*

- $k$ nearest neighbor graph (knn-graph) *[handwritten: euclidean nearest in ambient space.]*
  - the neighbors of $x_i$ are its $k$-nearest data points
  - $k$ is the parameter (what happens for $k = 0$, $k = 1$, $\ldots$, $k = n$)?
- $\varepsilon$ graph
  - the neighbors of $x_i$ are data points closer than $\varepsilon$
  - $\varepsilon$ is a parameter (what happens for $\varepsilon \in [0, \infty)$?)

Which to choose?

- $k$ nearest neighbor is more flexible (does not depend on the scaling).
- However, both do not guarantee that all points are connected in one graph.

**B**

euclidean distance in the ambient space

"geodesic distance" walking on the manifold.

(jumping from pt. to pt. short dist.)

(edge weighted)

Idea: Graph with vertices $x_1, ..., x_n$ (data pts)
& edges = connection to close neighbors.

---

2NN-graph

ε-graph

### Basic idea of ISOMAP

Calculate (geodesic) distances along the graph and find low dimensional embedding using multi dimensional scaling. *geodesic* means "along the manifold"

### Basic idea of locally linear embedding (LLE)

Approximate the manifold locally linearly, then reconstruct a low dimensional embedding matching the local linear structure.

## Basic idea of ISOMAP

Calculate (geodesic) distances along the graph and find low dimensional embedding using multi dimensional scaling.
*geodesic* means "along the manifold"

## Algorithm 13.2 (ISOMAP, sketch)

1. *construct neighborhood graph*
2. *compute all-pairs-shortest paths along the graph*
3. *apply MDS to get low dimensional embedding*

## Basic idea of locally linear embedding (LLE)

Approximate the manifold locally linearly, then reconstruct a low dimensional embedding matching the local linear structure.

## Algorithm 13.3 (LLE, sketch)

1. *construct neighborhood graph*
2. *express data points as local linear combinations*
3. *solve eigenvalue problem to get low dimensional embedding*

k too big in relation to sample density

# ISOMAP

# ISOMAP: construct the neighborhood graph

Algorithm 13.4 (construct neighborhood graph)

*Given a $D \times n$ data matrix $X = [x_1, \ldots, x_n] \in \mathbb{R}^{D \times n}$.*

1. *calculate all distances $D_{ij} = \|x_i - x_j\|$*
2. *for each point $x_i$ make a list of its neighbors (either via $k$ or $\varepsilon$)*
3. *define a weights matrix $W$ with entries*

*(for these: geodesic = euclidean distance)*

$$W_{ij} = \begin{cases} 0 & \text{if } i = j \\ D_{ij} & \text{if } x_i \text{ and } x_j \text{ are neighbors} \\ \infty & \text{otherwise} \end{cases}$$

*For $0 < W_{ij} < \infty$ there is an edge from $x_i$ to $x_j$.*

- Question: is $W$ symmetric?
- Answer: yes, for $\varepsilon$-graphs, no for knn-graph
- also for knn-graph we might symmetrise...

# ISOMAP: calculate all-pairs-shortest paths (1)

Don't reinvent the wheel, instead look it up, e.g. in CLRS:

# ISOMAP: calculate all-pairs shortest paths (2)

## Algorithm 13.5 (Floyd-Warshall algorithm, all-pairs shortest paths)

*Given a graph with edge weights $W_{ij}$ (their lengths) calculate a matrix D of all lengths of all paths along the graph. Initialize $D = W$.*

1. *update all lengths in D wrt paths via $x_1$*
2. *update all lengths in D wrt paths via $x_2$*
3. *. . .*
n. *update all lengths in D wrt paths via $x_n$*

$$\begin{pmatrix} 0 & a & b & \infty & \dots & \infty \\ c & & & & & \\ & & \ddots & & & \\ & & & \infty & \ddots & \\ \infty & \infty & & & & 0 \end{pmatrix}$$

```python
def floyd_warshall(W):    # Python/numpy: 3 loops
    # assumes non-squared distances in W
    D = W.copy()
    n = D.shape[0]
    for k in range(n):
        for i in range(n):
            for j in range(n):
                D[i,j] = min(D[i,j], D[i,k]+D[k,j])
    return D
```

$O(n^3)$

# ISOMAP: calculate all-pairs shortest paths (2)

> **Algorithm 13.6 (Floyd-Warshall algorithm, all-pairs shortest paths)**
>
> *Given a graph with edge weights $W_{ij}$ (their lengths) calculate a matrix D of all lengths of all paths along the graph. Initialize $D = W$.*
>
> 1. *update all lengths in D wrt paths via $x_1$*
> 2. *update all lengths in D wrt paths via $x_2$*
> 3. *...*
> n. *update all lengths in D wrt paths via $x_n$*

```python
# check with %prun or %lprun which is faster!
def floyd_warshall(W):    # Python/numpy: 2 loops
    # assumes non-squared distances in W
    D = W.copy()
    n = D.shape[0]
    for k in range(n):
        for i in range(n):
            D[i,:] = np.minimum(D[i,:], D[i,k]+D[k,:])
    return D
```

$\mathcal{O}(n^3)$

Result: Matrix with geodesic distances

# ISOMAP: calculate all-pairs shortest paths (2)

**Algorithm 13.7 (Floyd-Warshall algorithm, all-pairs shortest paths)**

*Given a graph with edge weights $W_{ij}$ (their lengths) calculate a matrix D of all lengths of all paths along the graph. Initialize $D = W$.*

1. *update all lengths in D wrt paths via $x_1$*
2. *update all lengths in D wrt paths via $x_2$*
3. *...*
n. *update all lengths in D wrt paths via $x_n$*

```
function fw(W)   % matlab version
    D = W # W[i,i]==0, W[i,j]=dij for neighbors, otherwise W[i,j]=inf
    n = size(W,1)
    for k in 1:n
        for i in 1:n
            for j in 1:n
                D[i,j] = min(D[i,j], D[i,k]+D[k,j])
            end
        end
    end
    return D
end
```

# ISOMAP: summary

### Algorithm 13.8 (ISOMAP)

*Given a $D \times n$ data matrix $X = [x_1, \ldots, x_n] \in \mathbb{R}^{D \times n}$.*

1. *construct neighborhood graph represented by $n \times n$ matrix $W$*
2. *compute all-pairs-shortest paths along the graph, ie. calc $D$*
3. *apply MDS transform distances $D$ into embedding $Z = [z_1, \ldots, z_n]$*

What is MDS?

# MDS = Multi-dimensional scaling

**MDS in a nutshell**

Given a distance matrix, recover a data matrix with those distances.

Data matrix

$$X = [x_1, \ldots, x_n] \in \mathbb{R}^{d \times n} \text{ (with entries } X_{ki}) \quad \text{where we assume } \sum_{i=1}^{n} x_i = 0.$$

Gram matrix, aka inner product matrix, aka kernel matrix

$$G = X^\mathsf{T} X \in \mathbb{R}^{n \times n} \qquad \text{with entries } G_{ij} = x_i^\mathsf{T} x_j$$

Covariance matrix, aka outer product matrix (omitting $1/n$)

$$C = X X^\mathsf{T} = \sum_{i=1}^{n} x_i x_i^\mathsf{T} \in \mathbb{R}^{d \times d} \qquad \text{with entries } C_{kl} = \sum_{i=1}^{n} X_{ki} X_{li}$$

Squared distance matrix $\qquad \langle x_i, x_i \rangle + \langle x_j, x_j \rangle - 2 \langle x_i, x_j \rangle$

$$D \in \mathbb{R}^{n \times n} \qquad \text{with entries } D_{ij} = \|x_i - x_j\|^2 \overset{\shortparallel}{=} (x_i - x_j)^\mathsf{T}(x_i - x_j)$$

$$G = X^\mathsf{T} X$$

$$\Lambda^{\frac{1}{2}} = \begin{pmatrix} \sqrt{\lambda_1} & 0 \\ 0 & \sqrt{\lambda_n} \end{pmatrix}$$

Note:

▸ $X^\mathsf{T} X = V \Lambda^{1/2} \Lambda^{1/2} V^\mathsf{T} = G$

▸ the mean of $X$ is arbitrary, since it doesn't change the distances

▸ $X$ can be arbitrarily rotated, since it doesn't change the distances

Question:

▸ How can we calculate the Gram matrix from the distance matrix?

1. *Squared distances can be calculated from the Gram matrix*

$$D_{ij} = \underbrace{x_i^{\mathsf{T}} x_i} + \underbrace{x_j^{\mathsf{T}} x_j} - 2x_i^{\mathsf{T}} x_j = (x_i - x_j)^{\mathsf{T}} (x_i - x_j)$$

$$D = \underbrace{(G \odot I_n) 1_n 1_n^{\mathsf{T}}} + \underbrace{1_n 1_n^{\mathsf{T}} (G \odot I_n)} - 2G$$

*with $G \odot I_n = \texttt{diagm(diag(G))}$ being the matrix with the squared norms $x_1^{\mathsf{T}} x_1, \ldots, x_n^{\mathsf{T}} x_n$ along the diagonal.*

2. *Removing the mean from the ones matrix results in the zero matrix*

$$1_n 1_n^{\mathsf{T}} H = 1_n 1_n^{\mathsf{T}} - 1_n (1_n^{\mathsf{T}} 1_n) 1_n^{\mathsf{T}} / n = 0_{n \times n} = H 1_n 1_n^{\mathsf{T}}$$

*where $H = I_n - 1_n 1_n^{\mathsf{T}} / n$ is the centering matrix, note $1_n^{\mathsf{T}} 1_n = n$*

3. *Assuming the mean of dataset $X$ is zero, i.e. $X 1_n = 0_n$, we have*

$$XH = X - X 1_n 1_n^{\mathsf{T}} / n = X - 0_n 1_n^{\mathsf{T}} / n = X$$

*Furthermore we have for the Gram matrix $G = X^{\mathsf{T}} X$*

$$HGH = HX^{\mathsf{T}} XH = X^{\mathsf{T}} X = G$$

$$A \odot B = \begin{pmatrix} a_{11} & \cdots & a_{1n} \\ & \ddots & \vdots \\ & & a_{nn} \end{pmatrix} \odot \begin{pmatrix} b_{11} & \cdots & b_{1n} \\ & \ddots & \\ & & b_{nn} \end{pmatrix}$$

$$\parallel$$

$$\begin{pmatrix} a_{11} \cdot b_{11} & & \\ & \ddots & \\ & & a_{nn} \cdot b_{nn} \end{pmatrix}$$

### Theorem 13.11

*Assuming the mean of dataset X is zero, we have*

$$G = -\frac{1}{2} HDH$$

*Thus we can calculate the Gram matrix from the squared distance matrix under the assumption of zero mean.*

Proof:

$$
\begin{aligned}
HDH &= H(G \odot I_n)1_n 1_n^\mathsf{T} H + H 1_n 1_n^\mathsf{T}(G \odot I_n)H - 2HGH \\
&= H(G \odot I_n)0_{n \times n} + 0_{n \times n}(G \odot I_n)H - 2G \\
&= -2G
\end{aligned}
$$

**Algorithm 13.12 (Multi-dimensional scaling, MDS)**

*Given a matrix D of the squared distances.*

1. *Calculate Gram matrix $G = -\frac{1}{2}HDH$*
2. *Calculate EVD of $G = V\Lambda V^{\mathsf{T}}$*
3. *Calculate data matrix $X = \Lambda^{1/2}V^{\mathsf{T}}$*

$\left.\begin{array}{c}\\ \\\end{array}\right\}$ take "square root" of $G$

the columns of this
will be our
dimension reduced
data pts.

`MDSDemo.ipynb`

# ISOMAP: summary

### Algorithm 13.13 (ISOMAP)

*Given a $D \times n$ data matrix $X = [x_1, \ldots, x_n] \in \mathbb{R}^{D \times n}$.*

1. *construct neighborhood graph represented by $n \times n$ matrix $W$*
2. *compute all-pairs-shortest paths along the graph, ie. calc $D$*
3. *apply MDS transform distances $D$ into embedding $Z = [z_1, \ldots, z_n]$*

$n$-dim points, but flat in $\mathbb{R}^n$, then do PCA for example

**LLE**

# LLE: sketch (1)



$$X_i = \sum_k W_{ik} Y_k$$

① Select neighbors

② Reconstruct with linear weights

③ Map to embedded coordinates

from Roweis/Saul's paper on LLE

# LLE: sketch (2)

1. for each data point $x_i$ select its neighbors, e.g. its $k$ nearest neighbors
2. find weights $W$ that locally reconstruct the data linearly (recall definition of manifold), i.e. minimize

$$\epsilon(W) = \sum_i \left\| x_i - \sum_j W_{ij} x_j \right\|^2$$

   with $W_{ij} = 0$ if $x_j$ is not a neighbor of $x_i$ and $\sum_j W_{ij} = 1$
3. find low dimensional embedding $Z = [z_1, \ldots, z_n]$ that minimizes

$$\phi(Z) = \sum_i \left\| z_i - \sum_j W_{ij} z_j \right\|^2$$

# LLE: finding the weights (1)

Constraint optimization problem:

$$\min_W \quad \sum_i \left\| x_i - \sum_j W_{ij} x_j \right\|^2$$

$$\text{s.t.} \quad \sum_j W_{ij} = 1 \text{ for } i = 1, \ldots, n$$

Can be solved for each data point $x$ with neighbors $\eta_j$ separately

$$\min_w \quad \left\| x - \sum_j w_j \eta_j \right\|^2$$

$$\text{s.t.} \quad \sum_j w_j = 1$$

# LLE: finding the weights (2)

Rewrite objective

$$\frac{1}{2}\left\|x - \sum_j w_j \eta_j\right\|^2 = \frac{1}{2}\left\|\sum_j w_j(x - \eta_j)\right\|^2 \quad \text{because } \sum_j w_j = 1$$

$$= \frac{1}{2}\sum_{jk} w_j \underbrace{(x - \eta_j)^\mathsf{T}(x - \eta_k)}_{=:C_{jk}} w_k$$

$$= \frac{1}{2}w^\mathsf{T} C w \quad \text{where matrix } C \text{ has entries } C_{jk}$$

Lagrangian

$$L(w, \lambda) = \frac{1}{2}w^\mathsf{T} C w - \lambda(1^\mathsf{T} w - 1)$$

$$dL = (w^\mathsf{T} C - \lambda 1^\mathsf{T})dw$$

$$w = \lambda C^{-1} 1$$

$$w_j = \lambda \sum_k (C^{-1})_{jk}$$

# LLE: finding the weights (3)

How to choose $\lambda$?

$\lambda$ is the normalizing constant, i.e. choose it s.t. $1^\mathsf{T} w = 1$:

$$\text{solve } Cw = 1 \qquad\qquad \text{then enforce } 1^\mathsf{T} w = 1$$

Summary (how to find a weight vector for a single data point)

1. find nearest neighbors $\eta_1, \ldots, \eta_k$
2. calculate local covariance matrix $C$ with entries $(x - \eta_j)^\mathsf{T}(x - \eta_k)$
3. define weights $w$ according to formula above

Summary (how to combine the weight vectors into a $n \times n$ matrix $W$)

1. calculate weight vectors for all data points
2. for a data point $x_i$ store the weights for each neighbor in the appropriate entries in the $i$th column of $W$, i.e. fill lots of zeros for non-neighbors

# LLE: finding the embedding

Constraint optimization problem

$$\min_Y \quad \sum_i \left\| y_i - \sum_j W_{ij} y_j \right\|^2$$
$$\text{s.t.} \quad Y1 = 0 \quad \text{zero mean}$$
$$YY^\mathsf{T} = I \quad \text{covariance identity}$$

Can be rewritten as:

$$\min_Y \quad \text{tr } YMY^\mathsf{T}$$
$$\text{s.t.} \quad Y1 = 0$$
$$YY^\mathsf{T} = I$$
$$M = I - W^\mathsf{T} - W + W^\mathsf{T}W = (I - W)^\mathsf{T}(I - W)$$

This can be solved: one can show that the rows $Y$ are the eigenvectors to the smallest eigenvalues. The math is almost identical to the simultaneous PCA problem (see appendix of PCA section).

# LLE: summary

almost exactly copied from `https://www.cs.nyu.edu/~roweis/lle/algorithm.html`

## Algorithm 13.14 (Locally linear embedding (LLE))

*Given a $D \times n$ data matrix $X = [x_1, \ldots, x_n] \in \mathbb{R}^{D \times n}$.*

1. *find nearest neighbors*

   ```
   for i=1:N
     compute the distance from Xi to every other point Xj
     find the K smallest distances
     assign the corresponding points to be neighbours of Xi
   end
   ```

2. *solve for reconstruction weights*

   ```
   for i=1:N
     create matrix Z consisting of all neighbours of Xi
     subtract Xi from every column of Z
     compute the local covariance C=Z'*Z
     solve linear system C*w = 1 for w
     set Wij=0 if j is not a neighbor of i
     set the remaining elements in the ith row of W equal to w/sum(w);
   end
   ```

3. *compute embedding coordinates*

   ```
   create sparse matrix M = (I-W)'*(I-W)
   find bottom d+1 eigenvectors of M
     (corresponding to the d+1 smallest eigenvalues)
   set the qth ROW of Y to be the q+1 smallest eigenvector
     (discard the bottom eigenvector [1,1,1,1...] with eigenvalue zero)
   ```

# MVU

Basic idea of maximum variance unfolding (MVU, aka semidefinite embedding)

Unfold the manifold by keeping local distances constant and maximizing all other distances.

Algorithm 13.15 (MVU, sketch)

1. *construct neighborhood graph*
2. *solve semidefinite programming problem*

# Maximum variance unfolding (1)

> Maximize the variance while keeping local distances constant:
>
> $$\max_Y \quad \frac{1}{2n} \sum_{ij} \|y_i - y_j\|^2$$
>
> $$\text{s.t.} \quad \|y_i - y_j\|^2 = \|x_i - x_j\|^2 \text{ if } x_i \text{ and } x_j \text{ are neighbors}$$
>
> $$\sum_i y_i = 0 \quad \text{i.e. embedding has mean zero}$$
>
> Unfortunately, this problem is tough to solve. Is it convex? No.

Let's rewrite the problem with the Gram matrix $G = Y^\mathsf{T} Y$.

# Maximum variance unfolding (2)

see e.g. http://en.wikipedia.org/wiki/Semidefinite_embedding

Assuming mean zero for the embedding implies:

$$\sum_{ij} G_{ij} = 1_n^\mathsf{T} G 1_n = (Y1_n)^\mathsf{T} Y1_n = 0 \ \text{ since } Y1_n = 0$$

Rewrite distances:

$$\|y_i - y_j\|^2 = G_{ii} + G_{jj} - G_{ij} - G_{ji}$$

Rewrite objective function:

$$\begin{aligned}
\sum_{ij} \|y_i - y_j\|^2 &= \sum_{ij} (y_i^\mathsf{T} y_i + y_j^\mathsf{T} y_j - y_i^\mathsf{T} y_j - y_j^\mathsf{T} y_i) \\
&= \sum_{ij} y_i^\mathsf{T} y_i + \sum_{ij} y_j^\mathsf{T} y_j - \sum_{ij} y_i^\mathsf{T} y_j - \sum_{ij} y_j^\mathsf{T} y_i \\
&= 2(\sum_{ij} G_{ii} - \sum_{ij} G_{ij}) = 2n \operatorname{tr} G
\end{aligned}$$

# Maximum variance unfolding (2a) — matrix version

see e.g. http://en.wikipedia.org/wiki/Semidefinite_embedding

Assuming mean zero for the embedding implies:

$$1_n^\top G 1_n = 1_n^\top Y^\top Y 1_n = 0 \text{ since } Y 1_n = 0$$

Rewrite matrix of squared distances:

$$D_{ij} = \|y_i - y_j\|^2$$
$$D = (G \odot I) 1_n 1_n^\top + 1_n 1_n^\top (G \odot I) - 2G$$

Rewrite objective function, sum of squared distances:

$$\begin{aligned}
1_n^\top D 1_n &= 1_n^\top (G \odot I) 1_n 1_n^\top 1_n + 1_n^\top 1_n 1_n^\top (G \odot I) 1_n - 2(1_n^\top G 1_n) \\
&= 1_n^\top \operatorname{diag}(G) n + n \operatorname{diag}(G)^\top 1_n - 0 \\
&= 2n \operatorname{tr} G
\end{aligned}$$

# Maximum variance unfolding (3)

see e.g. http://en.wikipedia.org/wiki/Semidefinite_embedding

> **Maximum variance unfolding as semi-definite programming**
>
> $$\max_G \quad \operatorname{tr} G$$
> $$\text{s.t.} \quad G_{ii} + G_{jj} - G_{ij} - G_{ji} = \|x_i - x_j\|^2 \text{ if } x_i \text{ and } x_j \text{ are neighbors}$$
> $$\sum_{ij} G_{ij} = 0$$
> $$G \geq 0$$
>
> The constraint $G \geq 0$ means that $G$ must be positive (semi-)definite. This ensures that we can recover $Y$ s.t. $G = Y^\mathsf{T} Y$.

Notes:

▸ now it is a convex optimization problem (see Boyd's book)

▸ called *semi-definite programming*, fast solvers available

▸ however, still expensive since $G$ has many variables

# How to map back and forth?

# Summary of ISOMAP, LLE, MVU

Words of the day

- nonlinear dimensionality reduction
- manifold, swissroll
- proximity graph, knn graph, $\varepsilon$ graph
- Euclidean vs geodesic distances
- ISOMAP, LLE, MVU, (and many more)
- semidefinite programming

How to estimate the dimension of the manifold?

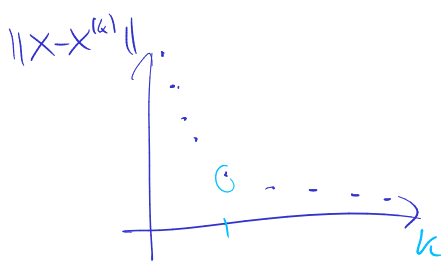$X$ data matrix $(x_1 | \dots | x_n) = \left( \frac{a_1}{a_D} \right)$

$(D < n)$

has rank $D$ if we really need all of the $D$ dim.s

Can try to approx. $X$ by a rank $d$-matrix $X^{(d)}$

$\left( X = U \Sigma V^\top \quad \text{SVD of } X \quad , \quad \Sigma = \begin{pmatrix} \sigma_1 \dots \sigma_r & 0 \\ 0 & 0 \end{pmatrix} \right)$

$X^{(k)} := U \Sigma^{(k)} V^\top$

$\Sigma^{(k)} = \begin{pmatrix} \sigma_1 \dots \sigma_k & 0 \\ 0 & 0 \dots 0 \end{pmatrix}$

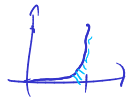Check $\| X - X^{(k)} \|$ for many $k$, then take last $k$ where the distance is "small" $=: d$.

Another option:  2NN estimator (Facco et.al. 2017)



$$\mu_i = \frac{r_{i,2}}{r_{i,1}}$$

Curse of dimensionality:



Prob. distr. of the random var $\mu$:

$$P(\mu) = d \cdot \frac{1}{\mu^{d+1}}$$

Mean: $\int_1^\infty \mu \cdot d \cdot \frac{1}{\mu^{d+1}} \, d\mu = \int_1^\infty d \cdot \mu^{-d} \, d\mu$

$$= \left[ -\frac{d}{d-1} \mu^{-(d-1)} \right]_1^\infty = \frac{d}{d-1}$$

Estimate $E\mu = \frac{d}{d-1}$, then solve for $d$.