#### Applied Machine Learning (AML)

#### Non-Linear Dimensionality Reduction

Oisin Mac Aodha • Siddharth N.

### Outline

#### **Non-Linearity**

- Extensions for Linear Dimensionality Reduction
  - Kernel PCA
- Visualisation
  - Multi-Dimensional Scaling (MDS)
  - Isomap
  - Locally Linear Embeddings (LLE)
  - *t*-distributed Stochastic Neighbour Embedding (*t*-SNE)
  - Uniform Manifold Approximation & Projection (UMAP)

informatics

### Extensions for Linear Dimensionality Reduction

# PCA on Non-Linear Data

#### **Example: Shells**



1

### **Feature Transformation**

**Key Idea:** Transform inputs  $\boldsymbol{x}$  using a feature map  $\phi(\boldsymbol{x})$  $\boldsymbol{x} \in \mathbb{R}^{D}, \, \boldsymbol{\phi}(\boldsymbol{x}) \in \mathbb{R}^{C}, \, C > D!$ 



### **Kernel PCA**

 $X = [\boldsymbol{x}_1; \ldots; \boldsymbol{x}_N],$  $\kappa(\boldsymbol{x}_i, \boldsymbol{x}_j) = \phi(\boldsymbol{x}_i)^\top \phi(\boldsymbol{x}_j)$ (kernel function)

 $S = \frac{1}{N} X X^{\mathsf{T}}$  $S\boldsymbol{v} = \lambda \boldsymbol{v}$ 

### **Kernel Trick**

No need to compute  $\phi(\mathbf{x})$ —only  $\kappa(\mathbf{x}_i, \mathbf{x}_i)!$ 

informatics



 $S = \frac{1}{N} \Phi(X) \Phi(X)^{\top}$  $\boldsymbol{v} = \Phi(\boldsymbol{X})\boldsymbol{a} = \sum_{i=1}^{N} a_i \phi(\boldsymbol{x}_i)$ 

(...some algebra)

 $K\boldsymbol{a} = N\lambda\boldsymbol{a}$ 

# Kernel PCA: Example



### **Kernel PCA**

#### Choosing the right kernel

- Not an easy choice
- Construct by hand/eye where feasible
- Possible to learn kernel matrices *K* directly from data!

Several non-linear dimensionality reduction methods can be viewed as kernel PCA, with kernels learned from data [1]

1. J. Ham et al, A Kernel View of the Dimensionality Reduction of Manifolds, 2004

4

#### Visualisation

#### **Manifold Hypothesis**

*High-dimensional data in the real world really lies on low-dimensional manifolds within that high-dimensional space.* 



### Overview

#### **Key Ideas**

- Difficult to construct a single global transformation of data
- Focus instead on some *local* measure of closeness
- Project data *x* onto lower-dimensional manifold as *e*
- Question: can we *preserve* local measure of closeness?

Let X denote the high-dimensional data space, and  $\mathcal{E}$  denote the low-dimensional manifold space. We can define the following distance measures on the two spaces respectively

 $\mathcal{D}_{\mathcal{X}}(\boldsymbol{x}_i, \boldsymbol{x}_j)$ 

 $\mathcal{D}_{\mathcal{E}}(oldsymbol{e}_i,oldsymbol{e}_j)$ 

### Multi-Dimensional Scaling (MDS)

Project data from X to  $\mathcal{E}$  while preserving the distance between *every pair of* samples in original data X

$$\mathcal{D}_{X}(\boldsymbol{x}_{i}, \boldsymbol{x}_{j}) = \|\boldsymbol{x}_{i} - \boldsymbol{x}_{j}\| \qquad \mathcal{D}_{\mathcal{E}}(\boldsymbol{e}_{i}, \boldsymbol{e}_{j}) = \|\boldsymbol{e}_{i} - \boldsymbol{e}_{j}\| \qquad (\text{example})$$
Objective: min  $\sum (\mathcal{D}_{X}(\boldsymbol{x}_{i}, \boldsymbol{x}_{j}) - \mathcal{D}_{\mathcal{E}}(\boldsymbol{e}_{i}, \boldsymbol{e}_{j}))^{2}$ 

- distance metrics  $\mathcal{D}_X, \mathcal{D}_\mathcal{E}$  could really be anything
- choosing L<sub>2</sub> helps make optimisation simpler

informatics

6

7

### **MDS:** Example

#### **Swiss Roll**



Data

MDS (Euclidean)

- Projects down to 2D while preserving distances
- Preserving distances *outside* the manifold!

informatics

# Isomap

Project data from X to  $\mathcal{E}$  while preserving the distance between points on the *embedded manifold*, not arbitrary distance!



#### 

8

### Find the Geodesic

- Generate nearest-neighbour graph *G* on data
- Shortest distance between points in this graph
  - Floyd-Warshall algorithm (all pairs shortest path)
- Perform MDS with  $\mathcal{D}_{X}(\mathbf{x}_{i}, \mathbf{x}_{j}) = \mathcal{G}_{FW}(\mathbf{x}_{i}, \mathbf{x}_{j})$

Objective: 
$$\min \sum_{i,j} (\mathcal{G}_{\mathsf{FW}}(\mathbf{x}_i, \mathbf{x}_j) - \|\mathbf{e}_i - \mathbf{e}_j\|)^2$$

J.B. Tenenbaum et al, A Global Geometric Framework for Nonlinear Dimensionality Reduction, 2000

### Isomap: Manifold



+ Floyd-Warshall

**Isomap: Example and Issues** 



#### Issues

- requires *uniform* and *dense* sampling on manifold
- prone to topological instabilities (effect of noise, non-convexity)
- can get disconnected graphs!
- slow with size of data Floyd-Warshall is  $O(N^3)$ !



10

q

# Locally Linear Embeddings (LLE)

Project data from X to  $\mathcal{E}$  while preserving the *linear transform that reconstructs points* from the *K* nearest neighbours



### Algorithm

- **1**. Find the weights  $w_k$
- 2. Fix weights  $w_k$  and find optimal embeddings  $e_i$ solved using a (sparse) eigen-decomposition

informatics



## LLE: Example and Issues



### Advantages

- globally optimal
- only bottleneck—finding  $e_i$
- *not* preserving specific distance

#### Issues

- reliance on local smoothness
- sensitive to noise
- no theoretical guarantees about manifold

13

# *t*-distributed Stochastic Neighbour Embedding (*t*-SNE)



Project data from X to  $\mathcal{E}$  while preserving the *probability distribution over* pairwise similarities between points in the space.

$$p(\mathbf{x}_{j}|\mathbf{x}_{i}) = \frac{\exp(-||\mathbf{x}_{i} - \mathbf{x}_{j}||^{2}/2\sigma_{i}^{2})}{\sum_{k \neq i} \exp(-||\mathbf{x}_{i} - \mathbf{x}_{k}||^{2}/2\sigma_{i}^{2})} \quad (i \neq j) \qquad \mathcal{D}_{\mathcal{E}}(\mathbf{e}_{i}, \mathbf{e}_{j}) \coloneqq \operatorname{Student} - t(||\mathbf{e}_{i} - \mathbf{e}_{j}||^{2}, v = 1)$$

$$p(\mathbf{x}_{i}|\mathbf{x}_{i}) = 0 \qquad \sum_{i} p(\mathbf{x}_{j}|\mathbf{x}_{i}) = 1 \qquad \mathcal{D}_{\mathcal{E}}(\mathbf{e}_{i}, \mathbf{e}_{j}) = \frac{\left(1 + ||\mathbf{e}_{i} - \mathbf{e}_{j}||^{2}\right)^{-1}}{\sum_{k} \sum_{l \neq k} \left(1 + ||\mathbf{e}_{k} - \mathbf{e}_{l}||^{2}\right)^{-1}}$$

$$\mathcal{D}_{\mathcal{X}}(\mathbf{x}_{i}, \mathbf{x}_{j}) = \frac{1}{2N} \left( p(\mathbf{x}_{i}|\mathbf{x}_{j}) + p(\mathbf{x}_{j}|\mathbf{x}_{i}) \right) \qquad \mathcal{D}_{\mathcal{E}}(\mathbf{e}_{i}, \mathbf{e}_{i}) = 0$$

$$\mathcal{D}_{\mathcal{X}}(\mathbf{x}_{i}, \mathbf{x}_{i}) = 0 \qquad \sum_{ij} \mathcal{D}_{\mathcal{X}}(\mathbf{x}_{i}, \mathbf{x}_{j}) = 1 \qquad \text{Student} t \text{ allows dissimilar points in } \mathcal{X} \text{ to be modelled far away in } \mathcal{E}!$$

Objective:  $\min \operatorname{KL}(\mathcal{D}_{\chi}(\boldsymbol{x}_i, \boldsymbol{x}_j) \| \mathcal{D}_{\mathcal{E}}(\boldsymbol{e}_i, \boldsymbol{e}_j))$ 

L.J.P. van der Maaten & G.E. Hinton, Visualizing Data using t-SNE, 2008 14





*t*-SNE: Examples

### So Far

- Different (relatively simple) aspects may be preserved in mapping X to  $\mathcal{E}$
- Trade-off in terms of simplicity of assumption and ease of optimisation
- *t*-SNE quite popular
  - can incorporate both global and local structure (distributional)
- works on large scale and high-dimensional data

#### Issue

None of these methods learn an *explicit* metric

• 
$$\{x_1,\ldots,x_N\} \stackrel{f}{\longrightarrow} \{e_1,\ldots,e_N\}$$
 but  $x_t \stackrel{igstyle}{\longrightarrow} e_t$  for unseen  $x_t$ 

cannot project an unseen point without redoing the optimisation!

#### informatics

16

# **Uniform Manifold Approximation & Projection (UMAP)**

Leverage Riemannian geometry and topology to construct a general framework for manifold learning and dimensionality reduction.

#### **Overview**

- Use simplices as basic building block
- Estimate connectivity and distances between points
- Construct graph that captures topology of manifold!







informatics

### **UMAP:**Dimensionality Reduction

- Construct faithful topological representation of data
- Compute *cross-entropy* between topological structures of X and  $\mathcal{E}$  in terms of the simplices (building blocks)
- Optimise low-dimensional representation to have minimum cross-entropy

#### Advantages

- Can learn a metric, so computing  $x_t \xrightarrow{f} e_t$  for unseen  $x_t$  is feasible (when labels available)!
- Is fast and scalable
- Can be run unsupervised, supervised, or even weakly supervised!

### **UMAP: Examples and Comparisons**



18

informatics

### Summary

- Non-linear dimensionality reduction helps visualise complex data in low dimensions—relies on the manifold hypothesis
- Can explore ways to transform data to run linear versions of algorithms (e.g. kernel PCA)
- Most methods exploit the nearest neighbour graph in some form (e.g. MDS, Isomap, LLE, etc.)
- Data is typically required to be clean (not much noise) and dense ...not too strong a requirement, especially with vision or language
- Many approaches to choose from—*t*-SNE and UMAP most popular