

Advanced course in machine learning
582744
Lecture 6

Arto Klami

Outline

Mixtures of linear models

Matrix factorization and recommender systems

- Recommender systems

- Matrix factorization

- Non-negative matrix factorization

Non-linear dimensionality reduction / manifold learning

- Multidimensional scaling

- Isomap

- Stochastic neighbor embedding

- Other methods

Pop-up course on linear models

Want to learn more about the previous lecture?

“Special course in unsupervised machine learning: Probabilistic factor analysis methods” in May 17-19

<https://www.cs.helsinki.fi/en/courses/582762/2016/k/k/1>

Organized by Suleiman Khan, Institute of Molecular Medicine (FIMM)

Mixtures of linear models

All of the FA models have marginal density that is a multivariate normal distribution with some low-rank covariance matrix

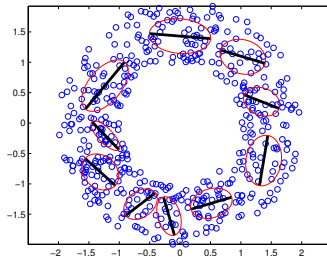
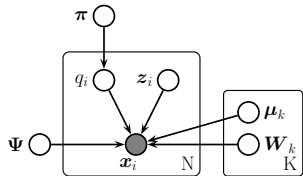
Hence, we can use them as component densities in a mixture model

Can either be interpreted as constrained mixture of Gaussians, or as piecewise-linear dimensionality reduction (though we do not get global coordinates)

All mixtures of Gaussians are actually special cases of mixture of factor analyzers; with zero factors we get spherical/diagonal covariance, etc.

...and mixtures with one cluster naturally give PCA/FA/CCA

Mixtures of linear models



Section 12.1.5. presents the EM algorithm for this general case

Recommender systems: collaborative filtering

Given triplets (user, item, rating), predict the rating for new (user, item) -pairs

Prototypical application for matrix factorization with missing data
 $\mathbf{X} \approx \mathbf{UV}^T$: Learn factorization based on the observed entries, predict the missing ones

The Netflix challenge: 100M ratings for 480k users and 18k movies (99% missing)

Alternative solutions based on user- or item-similarity:

- ▶ Find other users with similar preferences, recommend items they liked
- ▶ Find items similar to the ones the user liked (Amazon et al.)

(Outside collaborative filtering: Content-based recommendation)

Probabilistic matrix factorization

$$\mathbf{x}_{ij} = p_X\left(\sum_k \mathbf{u}_{ik} \mathbf{v}_{jk} \mid \theta_X\right)$$

$$\mathbf{U} = p_U(\theta_U)$$

$$\mathbf{V} = p_V(\theta_V)$$

Probabilistic PCA/FA is one example, using normal distributions as the densities above

For recommender systems we usually need to replace the likelihood with Bernoulli (for binary ratings) or, for example, Poisson distribution (“how many times one listened this song”)

Non-negative matrix factorization

If the data is positive, then we might want to model it with positive factors as well

Easier to interpret and encourages sparsity

Probabilistic formulation for example using Poisson likelihood and Gamma priors for the factors (Section 27.2.4)

Non-negative matrix factorization (Section 13.8.1)

Non-probabilistic formulation by minimizing $\|\mathbf{X} - \mathbf{WH}\|^2$ subject to the non-negativity constraint

Assuming we start with positive \mathbf{W} and \mathbf{H} , we can use gradient updates with small steps:

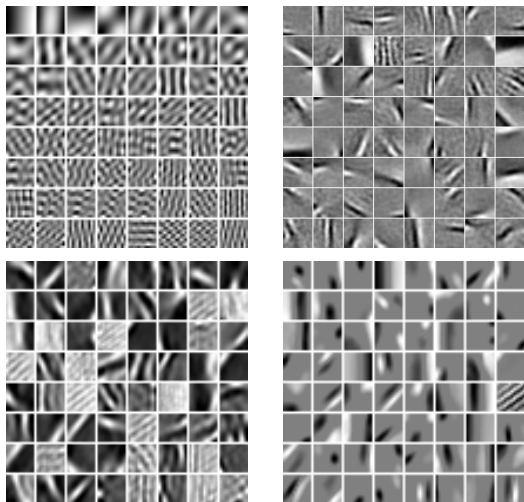
$$\begin{aligned}\mathbf{H}_{t+1} &= \mathbf{H}_t + \eta_H [(\mathbf{W}^T \mathbf{X}) - (\mathbf{W}^T \mathbf{W} \mathbf{H}_t)] \\ \mathbf{W}_{t+1} &= \mathbf{W}_t + \eta_W [(\mathbf{X} \mathbf{H}_t^T) - (\mathbf{W}_t \mathbf{H} \mathbf{H}_t^T)]\end{aligned}$$

Element-wise multiplicative updates by setting $\eta_H = \frac{\mathbf{H}_t}{(\mathbf{W}^T \mathbf{W} \mathbf{H}_t)}$:

$$\begin{aligned}\mathbf{H}_{t+1} &= \mathbf{H}_t \frac{(\mathbf{W}^T \mathbf{X})}{(\mathbf{W}^T \mathbf{W} \mathbf{H})} \\ \mathbf{W}_{t+1} &= \mathbf{W} \frac{\mathbf{X} \mathbf{H}^T}{\mathbf{W} \mathbf{H} \mathbf{H}^T}\end{aligned}$$

Can be shown to converge to a local optimum

Sparse dictionaries



Top row: PCA and ICA

Bottom row: NMF and sparse PCA

Non-linear dimensionality reduction

Linear models are efficient and often easy to interpret, but linearity is a strong assumption

Plenty of techniques for non-linear dimensionality reduction as well

We will go through a few alternatives not discussed in the book

Laplacian/spectral eigenmaps

We already know one method for this: Compute the graph laplacian of some graph connecting the observations, and represent the data with the first M eigenvectors

Then use, for example, PCA to drop into even smaller dimensionality

Emphasizes cluster structure due to the properties of the graph Laplacian

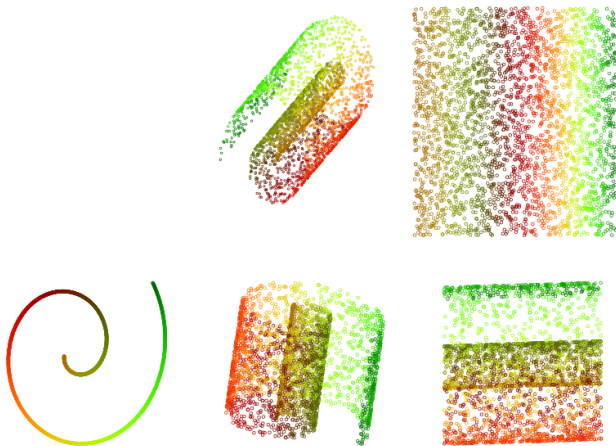
Manifold structure in data



Manifold structure in data



Manifold structure in data



Multidimensional scaling (MDS)

MDS is a general framework for non-linear dimensionality reduction techniques, starting with pairwise distances d_{ij}

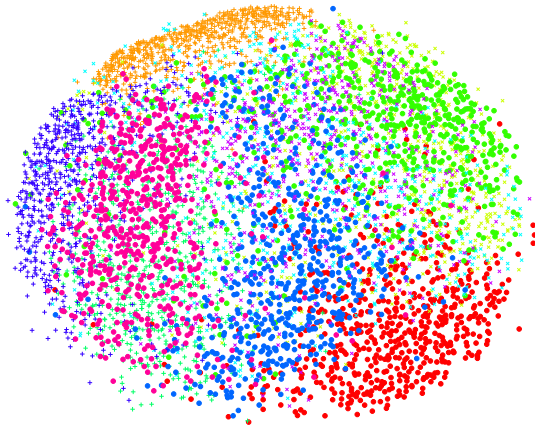
Construct a lower-dimensional representation \mathbf{x} for which we can compute \hat{d}_{ij} using Euclidean distance between \mathbf{x}_i and \mathbf{x}_j

Sammons's mapping: Loss $L(\hat{d}_{ij}) = \sum_{i < j} \frac{(\hat{d}_{ij} - d_{ij})^2}{d_{ij}}$ – we want representation for which the distances are as close to the original ones as possible

Gradient-based optimization over \mathbf{x} possible, but the loss function is not very nice

Computed for a fixed collection of N data points, and does not define a closed-form mapping for new points

Sammon's mapping



(b) Visualization by Sammon mapping.

From van der Maaten & Hinton, JMLR 2008

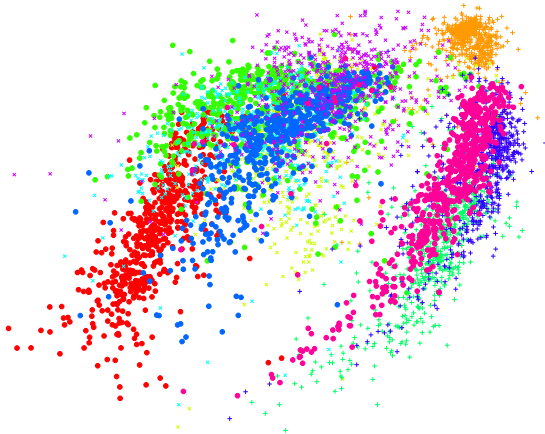
Isomap

One practical MDS algorithm is Isomap

- ▶ Compute the k nearest neighbor graph (the same we used in spectral clustering)
- ▶ Compute geodesic distances between all pairs of samples using the Floyd's algorithm
- ▶ Perform MDS on the resulting graph to find the low-dimensional representation; in practice PCA of the similarity matrix is often used

Attempts to preserve geodesic distances along the manifold, solving the swiss roll example

Isomap



(a) Visualization by Isomap.

From van der Maaten & Hinton, JMLR 2008

Stochastic neighbor embedding (SNE)

Another method that is commonly used today is SNE (or its more robust variant t-SNE)

Uses the loss

$$L = KL(p, q) = \sum_{i,j} p_{ij} \log \frac{p_{ij}}{q_{ij}},$$

where

$$q_{ij} = \frac{e^{-\|z_i - z_j\|^2}}{\sum_{k \neq j} e^{-\|z_i - z_k\|^2}}$$

defines a neighborhood in the low-dimensional space that should match the original neighborhood

$$p_{ij} = \frac{e^{-\|x_i - x_j\|^2}}{\sum_{k \neq j} e^{-\|x_i - x_k\|^2}}$$

Again we can use gradient-based optimization:

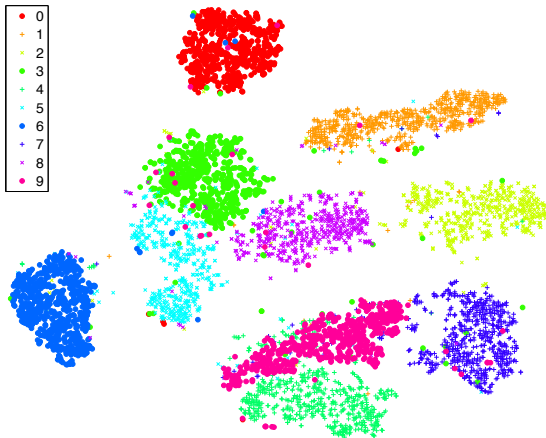
$$\frac{\partial L}{\partial z_i} = 2 \sum_j (z_i - z_j)(p_{ij} - q_{ij} + p_{ji} - q_{ji})$$

Stochastic neighbor embedding (SNE)

Attempts to capture the local structure as well as possible, using PCA to initialize the \mathbf{z} helps somewhat

Even better version, t-SNE, uses t-distributions in the embedding space to further discard global structure, often resulting in rather pretty visualizations

Stochastic neighbor embedding



(a) Visualization by t-SNE.

From van der Maaten & Hinton, JMLR 2008

Other non-linear dimensionality reduction methods

Kernel-PCA:

- ▶ PCA after a non-linear data transformation represented as a kernel; we will talk about kernel methods when discussing supervised learning
- ▶ Quite close explicitly transforming the input data and then performing PCA
- ▶ See Section 14.4.4, but note that this is not very widely used

Self-organizing map:

- ▶ Approximately k-means clustering where the clusters are organized into a grid
- ▶ Iterative updates where the closest centroid is pulled towards the sample and some of its closest neighbor in the grid are pulled with it
- ▶ We can then visualize the grid itself

Self-organizing map

