

# Unsupervised Dimension Reduction

## (Chapter 15)<sup>1</sup>

David Barber

(adapted by Miguel Palomino)

---

<sup>1</sup>These slides accompany the book *Bayesian Reasoning and Machine Learning*. The book and demos can be downloaded from [www.cs.ucl.ac.uk/staff/D.Barber/brml](http://www.cs.ucl.ac.uk/staff/D.Barber/brml). Feedback and corrections are also available on the site. Feel free to adapt these slides for your own purposes, but please include a link the above website.

# Principal Components Analysis

Data is often high dimensional. Provided there is some 'structure', data will typically lie close to a much lower dimensional 'manifold'. Here we concentrate on computationally efficient linear dimension reduction techniques.

$$\mathbf{y} = \mathbf{E}\mathbf{x} + \mathbf{c} \quad \text{or} \quad \mathbf{y} = \mathbf{E}\mathbf{x} \quad (\text{for centered data})$$

We express the approximation for (centered) datapoint  $\mathbf{x}^n$  as

$$\mathbf{x}^n \approx \sum_{j=1}^M y_j^n \mathbf{b}^j = \mathbf{B}\mathbf{y} \equiv \tilde{\mathbf{x}}^n$$

---

## Basis

The  $\mathbf{b}^j$  are 'basis' vectors that span the subspace. Collectively we can write  $\mathbf{B} = [\mathbf{b}^1, \dots, \mathbf{b}^M]$ .

---

## Low dimensional coordinates

The  $y_i^n$  are the low dimensional co-ordinates of the data.

# Minimal square loss approximation

To determine the best lower dimensional representation it is convenient to use the square distance error between  $\mathbf{x}$  and its reconstruction  $\tilde{\mathbf{x}}$ :

$$E(\mathbf{B}, \mathbf{Y}, \mathbf{c}) = \sum_{n=1}^N \sum_{i=1}^D [x_i^n - \tilde{x}_i^n]^2$$

The optimal basis  $\mathbf{B}$  is then obtained from the eigenvectors of the covariance matrix of the centred data:

$$\mathbf{S} = \mathbf{X}\mathbf{X}^T$$

The computation of the eigendecomposition of a  $D \times D$  matrix is  $O(D^3)$ . When the number  $N$  of datapoints is  $N \ll D$ , an alternative is to compute the eigendecomposition of  $\mathbf{X}^T\mathbf{X}$  in  $O(N^3)$  and use them to obtain those of  $\mathbf{S}$ .

# PCA via Singular value decomposition

Consider the SVD decomposition of the data matrix:

$$\mathbf{X} = \mathbf{U}\mathbf{D}\mathbf{V}^T$$

where  $\mathbf{U}^T\mathbf{U} = \mathbf{I}_D$  and  $\mathbf{V}^T\mathbf{V} = \mathbf{I}_N$  and  $\mathbf{D}$  is a diagonal matrix of the (positive) singular values. We assume that the decomposition has ordered the singular values so that the upper left diagonal element of  $\mathbf{D}$  contains the largest singular value. The matrix  $\mathbf{X}\mathbf{X}^T$  can then be written as

$$\mathbf{X}\mathbf{X}^T = \mathbf{U}\mathbf{D}\mathbf{V}^T\mathbf{V}\mathbf{D}\mathbf{U}^T = \mathbf{U}\mathbf{D}^2\mathbf{U}^T$$

Since  $\mathbf{U}\mathbf{D}^2\mathbf{U}^T$  is in the form of an eigen-decomposition, the PCA solution is equivalently given by performing the SVD decomposition of  $\mathbf{X}$ , for which the eigenvectors are then given by  $\mathbf{U}$ , and corresponding eigenvalues by the square of the singular values.

$$\mathbf{X} = \mathbf{U}\mathbf{D}\mathbf{V}^T \approx \underbrace{\mathbf{U}_M}_{D \times M} \underbrace{\mathbf{D}_M}_{M \times M} \underbrace{\mathbf{V}_M^T}_{M \times N} = \underbrace{\mathbf{B}}_{D \times M} \underbrace{\mathbf{Y}}_{M \times N}$$

where  $\mathbf{U}_M$ ,  $\mathbf{D}_M$ ,  $\mathbf{V}_M$  correspond to taking only the first  $M$  singular values of the full matrices.

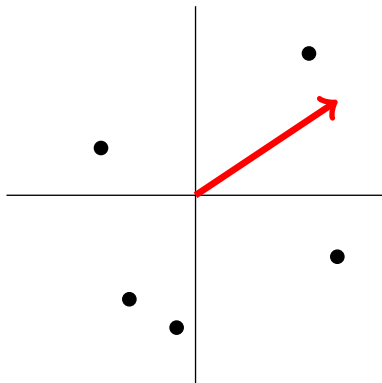
# Matrix Decompositions

Given a data matrix  $\mathbf{X}$  for which each column represents a datapoint, an approximate matrix decomposition is of the form  $\mathbf{X} \approx \mathbf{B}\mathbf{Y}$  into a basis matrix  $\mathbf{B}$  and weight (or coordinate) matrix  $\mathbf{Y}$ . Symbolically, matrix decompositions are of the form

$$\underbrace{\left( \begin{array}{c} X : \text{Data} \end{array} \right)}_{D \times N} \approx \underbrace{\left( \begin{array}{c} B : \text{Basis} \end{array} \right)}_{D \times M} \underbrace{\left( \begin{array}{c} Y : \text{Weights/Components} \end{array} \right)}_{M \times N}$$

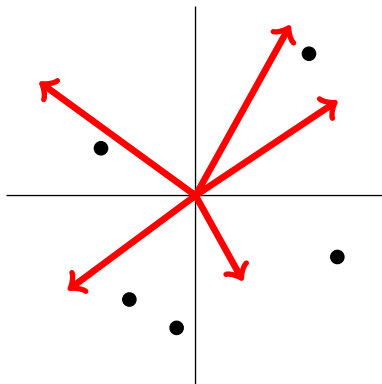
Based on the SVD of the data matrix, we see that PCA is in this class. Many methods can be considered as matrix decompositions under specific constraints.

# Under-complete decompositions



When  $M < D$ , there are fewer basis vectors than dimensions. The matrix  $\mathbf{B}$  is then called 'tall' or 'thin'. In this case the matrix  $\mathbf{Y}$  forms a lower dimensional approximate representation of the data  $\mathbf{X}$ , PCA being a classic example.

# Over-complete decompositions



For  $M > D$  the basis is over-complete, there being more basis vectors than dimensions. In such cases additional constraints are placed on either the basis or components. For example, one might require that only a small number of the large number of available basis vectors is used to form the representation for any given  $\mathbf{x}$ . Such sparse-representations are common in theoretical neurobiology where issues of energy efficiency, rapidity of processing and robustness are of interest.

# Probabilistic latent semantic analysis

Consider two objects,  $x$  and  $y$ , where  $\text{dom}(x) = \{1, \dots, I\}$  and  $\text{dom}(y) = \{1, \dots, J\}$  and a dataset  $(x^n, y^n, n = 1, \dots, N)$ .

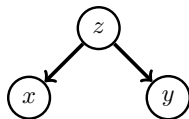
We have a count matrix with elements  $C_{ij}$  which describes the number of times the joint state  $x = i, y = j$  was observed in the dataset. We can transform this count matrix into a frequency matrix  $p$  with elements

$$p(x = i, y = j) = \frac{C_{ij}}{\sum_{ij} C_{ij}}$$



## Probabilistic latent semantic analysis (II)

Our interest is to find a decomposition of this frequency matrix of the form



That is,

$$\begin{aligned} \underbrace{p(x = i, y = j)}_{X_{ij}} &\approx \tilde{p}(x = i, y = j) \\ &= \sum_k \tilde{p}(x = i, y = j, z = k) \\ &= \sum_k \underbrace{\tilde{p}(x = i | z = k)}_{B_{ik}} \underbrace{\tilde{p}(y = j | z = k) \tilde{p}(z = k)}_{Y_{kj}} \end{aligned}$$

where all quantities  $\tilde{p}$  are distributions. This is then a form of matrix decomposition into **positive** basis  $\mathbf{B}$  and **positive** coordinates  $\mathbf{Y}$ . This has the interpretation of discovering latent topics  $z$  that describe the joint behaviour of  $x$  and  $y$ .

# An EM style training algorithm

For probabilities, a useful measure of discrepancy is the Kullback-Leibler divergence

$$\text{KL}(p|\tilde{p}) = \langle \log p \rangle_p - \langle \log \tilde{p} \rangle_p$$

Since  $p$  is fixed, minimising the Kullback-Leibler divergence with respect to the approximation  $\tilde{p}$  is equivalent to maximising the 'likelihood' term  $\langle \log \tilde{p} \rangle_p$ . This is

$$L \equiv \sum_{x,y} p(x,y) \log \tilde{p}(x,y)$$

It's convenient to derive an EM style algorithm to learn the **parameters**; in this case,  $\tilde{p}(x|z)$ ,  $\tilde{p}(y|z)$  and  $\tilde{p}(z)$ .

Consider

$$\begin{aligned} \text{KL}(q(z|x, y)|\tilde{p}(z|x, y)) \\ = \sum_z q(z|x, y) \log q(z|x, y) - \sum_z q(z|x, y) \log \tilde{p}(z|x, y) \geq 0 \end{aligned}$$

where  $\sum_z$  implies summation over all states of the variable  $z$ . Using

$$\tilde{p}(z|x, y) = \frac{\tilde{p}(x, y, z)}{\tilde{p}(x, y)}$$

and rearranging, this gives the bound,

$$\log \tilde{p}(x, y) \geq - \sum_z q(z|x, y) \log q(z|x, y) + \sum_z q(z|x, y) \log \tilde{p}(z, x, y)$$

Plugging this into the ‘likelihood’ term above, we have the bound

$$\begin{aligned} L \geq & - \sum_{x, y} p(x, y) \sum_z q(z|x, y) \log q(z|x, y) \\ & + \sum_{x, y} p(x, y) \sum_z q(z|x, y) [\log \tilde{p}(x|z) + \log \tilde{p}(y|z) + \log \tilde{p}(z)] \end{aligned}$$

# M-step

For fixed  $\tilde{p}(x|z), \tilde{p}(y|z)$ , the contribution to the bound from  $\tilde{p}(z)$  is

$$\sum_{x,y} p(x,y) \sum_z q(z|x,y) \log \tilde{p}(z)$$

Up to a constant, this is  $-\text{KL}\left(\sum_{x,y} q(z|x,y)p(x,y) \middle| \tilde{p}(z)\right)$  so that, optimally,

$$\tilde{p}(z) = \sum_{x,y} q(z|x,y)p(x,y)$$

Similarly, optimally

$$\tilde{p}(x|z) \propto \sum_y p(x,y)q(z|x,y)$$

and

$$\tilde{p}(y|z) \propto \sum_x p(x,y)q(z|x,y)$$

# E-step

The optimal setting for the  $q$  distribution at each iteration is

$$q(z|x, y) = \tilde{p}(z|x, y)$$

which is fixed throughout the M-step.

---

## Convergence

$L$  is guaranteed to increase (and the Kullback-Leibler divergence decrease) under iterating between the E and M-steps, since the method is analogous to an EM procedure.

# Conditional PLSA

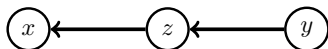
In some cases it is more natural to consider a conditional frequency matrix

$$p(x = i | y = j)$$

and seek an approximate decomposition

$$\underbrace{p(x = i | y = j)}_{X_{ij}} \approx \sum_k \underbrace{\tilde{p}(x = i | z = k)}_{B_{ik}} \underbrace{\tilde{p}(z = k | y = j)}_{Y_{kj}}$$

The model in this case is:



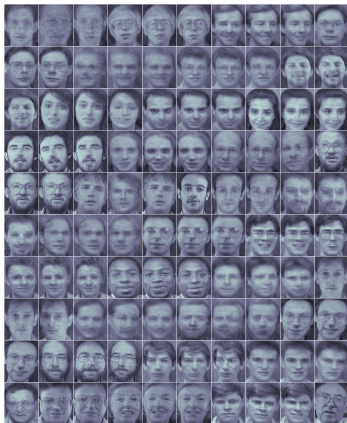
Deriving an EM style algorithm for this is also straightforward.

# Eigenfaces

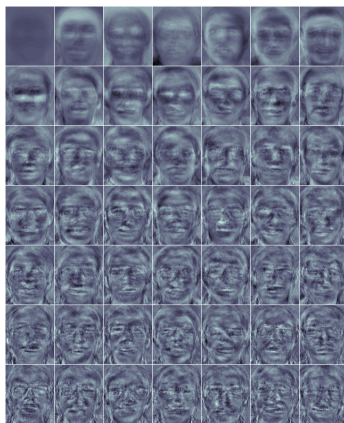


**Figure:** 100 of the 120 training images (40 people, with 3 images of each person). Each image consists of  $92 \times 112 = 10304$  non-negative greyscale pixels. The data is scaled so that, represented as an image, the components of each image sum to 1. The average value of each pixel across all images is  $9.70 \times 10^{-5}$ . This is a subset of the 400 images in the full Olivetti Research Face Database

# Eigenfaces



(a)

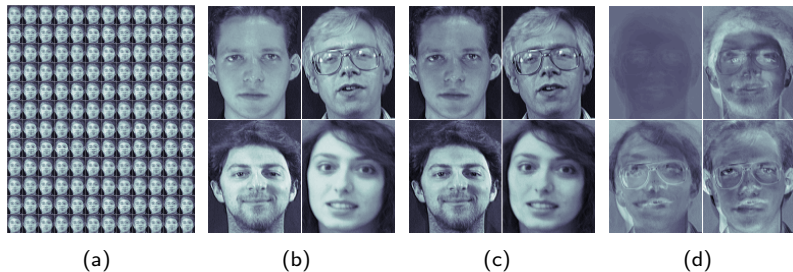


(b)

**Figure:** (a): SVD reconstruction of the images using a combination of the 49 eigen-images. (b): The eigen-images are found using SVD of the above data and taking the 49 eigenvectors with largest eigenvalue. The images corresponding to the largest eigenvalues are contained in the first row, and the next 7 in the row below, etc.



# Learning a positive basis

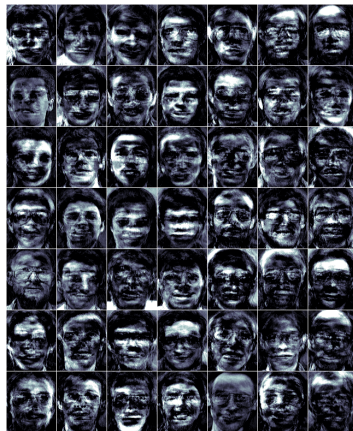


**Figure:** (a): Training data, consisting of a positive (convex) combination of the base images. (b): The chosen base images from which the training data is derived. (c): Basis learned using conditional PLSA on the training data. This is virtually indistinguishable from the true basis. (d): Eigenbasis (sometimes called 'eigenfaces').

# Positive reconstruction



(a)



(b)

**Figure:** (a): Conditional PLSA reconstruction of the images using a positive convex combination of the 49 positive base images in (b). The root mean square reconstruction error is  $1.391 \times 10^{-5}$ . The base images tend to be more 'localised' than the corresponding eigen-images. Here one sees local structure such as foreheads, chins, etc.

# Modelling citations

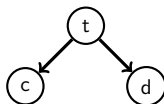
We have a collection of research documents which cite other documents. For example, document 1 might cite documents 3, 2, 10, etc. Given only the list of citations for each document, can we identify key research papers and the communities that cite them?

---

## A probabilistic formulation

We use the variable  $d \in \{1, \dots, D\}$  to index documents and  $c \in \{1, \dots, D\}$  to index citations (both  $d$  and  $c$  have the same domain, namely the index of a research article). If document  $d = i$  cites article  $c = j$  then we set the entry of the matrix  $C_{ij} = 1$ . If there is no citation,  $C_{ij}$  is set to zero. We can form a 'distribution' over documents and citations using

$$p(d = i, c = j) = \frac{C_{ij}}{\sum_{ij} C_{ij}}$$



c: citation  
t: topic  
d: document

and use PLSA to decompose this matrix into citation-topics.

# Modelling citations

The Cora corpus contains an archive of around 30,000 computer science research papers. From this archive the papers in the machine learning category are extracted, consisting of 4220 documents and 38,372 citations.

---

## Using PLSA

The joint PLSA method is fitted to the data using  $z = 7$  topics. From the trained model the expression  $p(c = j | z = k)$  defines how authoritative paper  $j$  is according to community  $z = k$ .

# Modelling citations

factor 1 0.0108 0.0066 0.0065	(Reinforcement Learning) Learning to predict by the methods of temporal differences. Sutton. Neuronlike adaptive elements that can solve difficult learning control problems. Barto et al. Practical Issues in Temporal Difference Learning. Tesauro.
factor 2 0.0038 0.0037 0.0036	(Rule Learning) Explanation-based generalization: a unifying view. Mitchell et al. Learning internal representations by error propagation. Rumelhart et al. Explanation-Based Learning: An Alternative View. DeJong et al.
factor 3 0.0120 0.0061 0.0049	(Neural Networks) Learning internal representations by error propagation. Rumelhart et al. Neural networks and the bias-variance dilemma. Geman et al. The Cascade-Correlation learning architecture. Fahlman et al.
factor 4 0.0093 0.0066 0.0055	(Theory) Classification and Regression Trees. Breiman et al. Learnability and the Vapnik-Chervonenkis dimension. Blumer et al. Learning Quickly when Irrelevant Attributes Abound. Littlestone.
factor 5 0.0118 0.0094 0.0056	(Probabilistic Reasoning) Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference. Pearl. Maximum likelihood from incomplete data via the em algorithm. Dempster et al. Local computations with probabilities on graphical structures. Lauritzen et al.
factor 6 0.0157 0.0132 0.0096	(Genetic Algorithms) Genetic Algorithms in Search, Optimization, and Machine Learning. Goldberg. Adaptation in Natural and Artificial Systems. Holland. Genetic Programming: On the Programming of Computers by Means of Natural Selection. Koza.
factor 7 0.0063 0.0054 0.0033	(Logic) Efficient induction of logic programs. Muggleton et al. Learning logical definitions from relations. Quinlan. Inductive Logic Programming Techniques and Applications. Lavrac et al.

**Table:** Highest ranked documents according to  $p(c|z)$ . The factor topic labels are manual assignments based on similarity to the Cora topics.