

(Non-linear) dimensionality reduction

Jiří Kléma

Department of Computer Science,
Czech Technical University in Prague



<http://cw.felk.cvut.cz/wiki/courses/b4m36san/start>

A brief review of PCA

- any real symmetric matrix is diagonalized by a column matrix of its eigenvectors \mathbf{E} ,
- \mathbf{C}_X is real and symmetric, it follows that

$$\mathbf{C}_X = \mathbf{E}\mathbf{D}\mathbf{E}^T$$

- the only trick is to select \mathbf{P} to be a matrix where each column \mathbf{p}_i is an eigenvector of \mathbf{C}_X

$$\mathbf{P} = \mathbf{E}$$

- we also know that any orthogonal matrix has the same inverse and transpose,
- the above-selected \mathbf{P} is necessarily orthogonal

$$\mathbf{P}^T\mathbf{P} = \mathbf{I} \Rightarrow \mathbf{P}^{-1} = \mathbf{P}^T$$

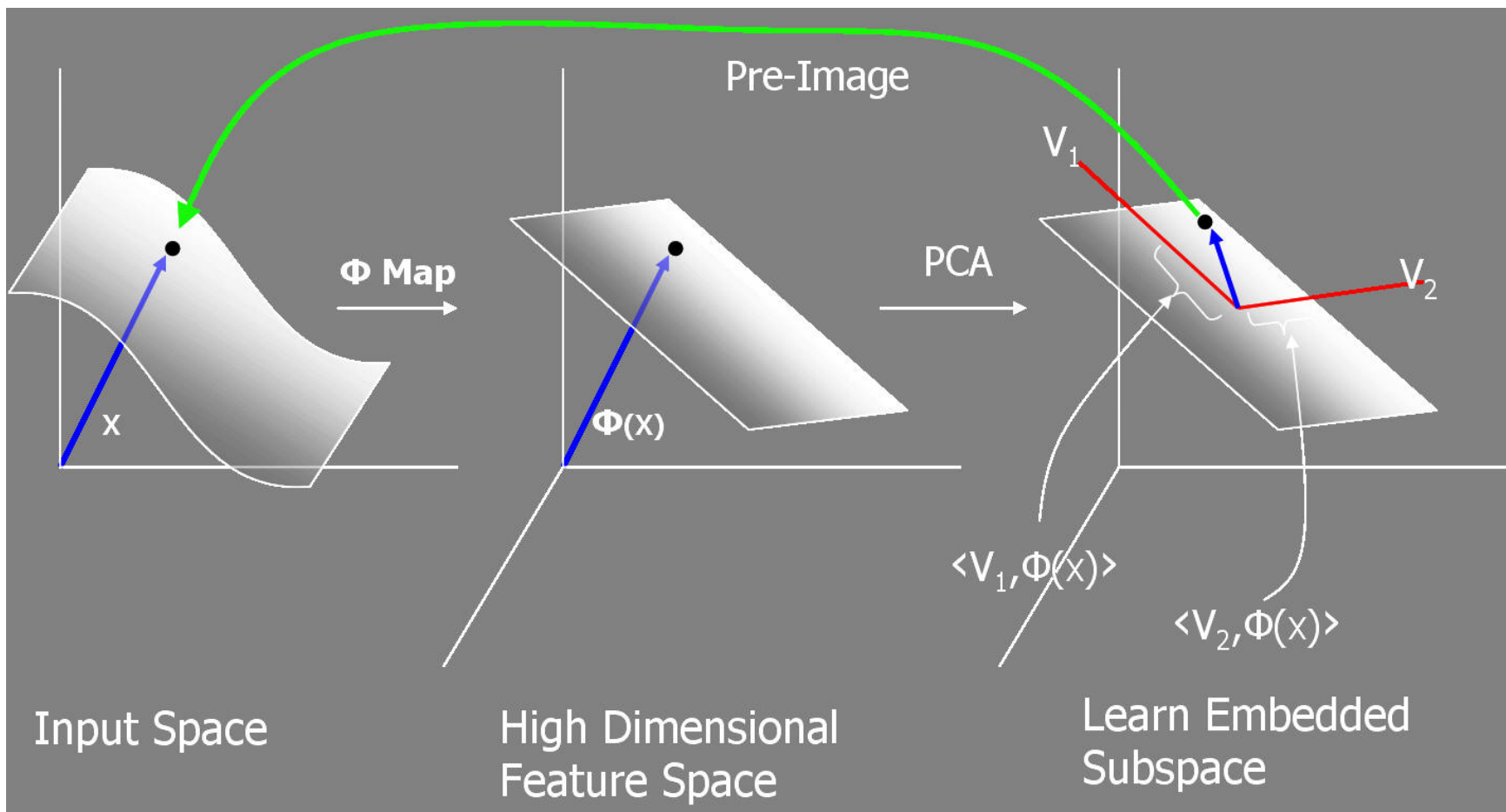
- then, it is easy to show that \mathbf{P} diagonalizes \mathbf{C}_T

$$\begin{aligned}\mathbf{C}_T &= \mathbf{P}^T\mathbf{C}_X\mathbf{P} = \mathbf{P}^T(\mathbf{E}\mathbf{D}\mathbf{E}^T)\mathbf{P} = \\ &= (\mathbf{P}^T\mathbf{P})\mathbf{D}(\mathbf{P}^T\mathbf{P}) = (\mathbf{P}^{-1}\mathbf{P})\mathbf{D}(\mathbf{P}^{-1}\mathbf{P}) = \mathbf{D}\end{aligned}$$

- PCA is solved by finding the eigenvectors of \mathbf{C}_X .

kernel PCA – the idea behind

- Introduce an intermediate feature space \mathcal{U}
 - $\mathcal{X} \rightarrow \mathcal{U} \rightarrow \mathcal{T}$, \mathcal{U} linearizes the original manifold.



<http://www.research.rutgers.edu>

kernel PCA

- $\phi(\mathbf{x}_i)$ are not available, we need to replace them by \mathbf{K} ,
- for $\lambda \geq 0$, \mathbf{v} 's are in the span of $\phi(\mathbf{x}_i)$,
- they can be written as linear combination of the object images

$$\mathbf{v} = \sum_{i=1}^m \alpha_i \phi(\mathbf{x}_i)$$

- we will substitute for \mathbf{v} into the eigenvector formula (the last in the previous slide)

$$\lambda \sum_{j=1}^m \alpha_j \phi(\mathbf{x}_j) = \frac{1}{m} \sum_{i=1}^m \phi(\mathbf{x}_i) \phi(\mathbf{x}_i)^T \sum_{j=1}^m \alpha_j \phi(\mathbf{x}_j)$$

- and use the trick to introduce the dot product, we will multiply by $\phi(\mathbf{x}_k)^T$

$$\lambda \sum_{j=1}^m \alpha_j \phi(\mathbf{x}_k)^T \phi(\mathbf{x}_j) = \frac{1}{m} \sum_{j=1}^m \alpha_j \sum_{i=1}^m (\phi(\mathbf{x}_k)^T \phi(\mathbf{x}_i)) (\phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j))$$

- the kernel function replaces all the occurrences of ϕ , when iterating $\forall k = 1 \dots m$ we obtain

$$\lambda \mathbf{K} \alpha = \frac{1}{m} \mathbf{K}^2 \alpha$$

- to diagonalize \mathbf{C}_T , we solve the eigenvalue problem for the kernel matrix

$$m\lambda\alpha = \mathbf{K}\alpha$$

Locally Linear Embedding (LLE) [Roweis, Saul, 2000]

- The ultimate case of piecewise linear modelling
 - approximation of the manifold by a combination of linear models,
 - in here, we have a linear model for each object,
- a special case of kernel PCA constructing a data-dependent kernel matrix
 - for some problems it is difficult to find a kernel for kernel PCA,
- advantages
 - efficient for large datasets, optimization does not involve local minima,
 - single parameter to tune (K),
 - invariant to scaling, rotation and translation,
- disadvantages
 - improper for representing future data,
 - can be unstable in sparse areas of the input space,
 - tends to collapse a lot of instances near the origin of \mathcal{T} .



t-Distributed Stochastic Neighbor Embedding (t-SNE)

- The key ideas are in the design of the stress function driving gradient descent search
 - convert Euclidean distances in both spaces into joint probabilities,
 - p_{ij} in the original space \mathcal{X} and q_{ij} in the reduced space \mathcal{T} ,
 - minimize their Kullback-Leibler divergence

$$S = KL(P||Q) = \sum_i \sum_{j \neq i} p_{ij} \log \frac{p_{ij}}{q_{ij}}$$

- The first “trick” lies in asymmetry of KL divergence
 - large p_{ij} modeled by small q_{ij} \rightarrow big penalty!
 - small p_{ij} modeled by large q_{ij} \rightarrow small penalty!
 - tends to preserve large p_{ij} 's and thus small distances in the original space.

SOM learning algorithm

1. Initialize weight vectors of all the neurons $\mathbf{w}_i^0 \in \mathcal{X} (\mathbb{R}^D)$, $\forall i = 1, \dots, N$
 - randomly or rather sample evenly from the subspace spanned by the largest principal component eigenvectors,
2. resample the instance set \mathbf{X}
 - A: sample systematically, B: **permute** the instance set, C: randomly (bootstrap),
3. get the next instance $\mathbf{x}_i \in \mathbb{R}^D$ from the resample,
4. find the best matching unit (BMU), i.e., the neuron nearest to \mathbf{x}_i ,
5. change the weights of neurons in the neighborhood of BMU (including BMU)

$$\mathbf{w}_i^{t+1} = \mathbf{w}_i^t + \alpha^t e^{-\frac{d(w_i^t, BMU)}{2(\sigma^t)^2}} (\mathbf{x}_i - \mathbf{w}_i^t)$$

- both the neighborhood size σ^t and the learning rate α^t decrease in time,
 - d is the distance in terms of the neighborhood relationship (e.g., the grid distance),
6. go to step 3 (or 2 when the sample is finished) until the given number of cycles is reached.



