

Principal Component and Correspondence Analyses: Multivariate Approach

Optum Global Analytics

23rd Nov' 2020

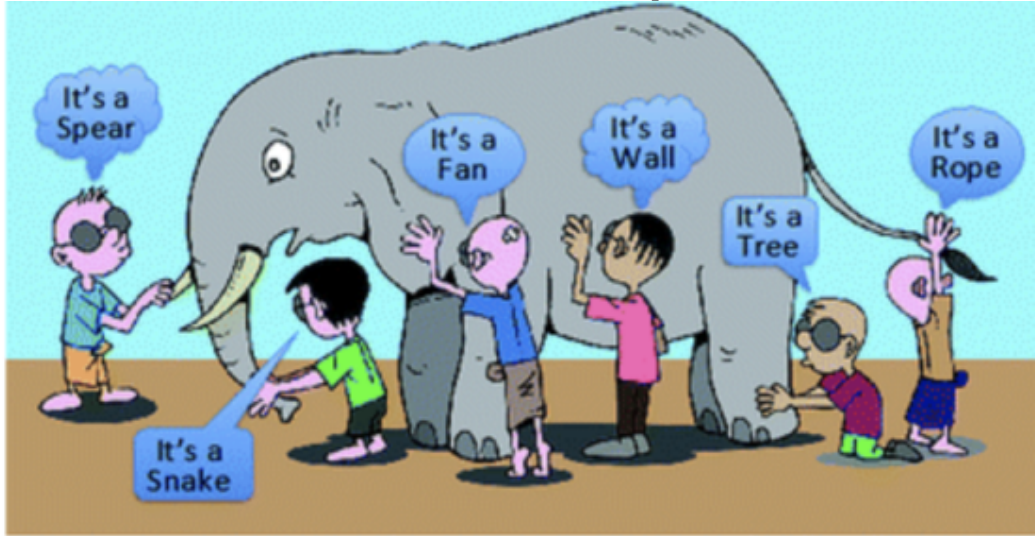


Agenda

- ❑ Intuitive way of thinking Multi Var. Statistics and PCA
- ❑ Concept of Dimensions Reduction
- ❑ Linear projections using PCA and its algebraic derivations
- ❑ PCA as maximizing variance and minimizing reconstruction loss
- ❑ Singular value decomposition (SVD)
- ❑ Non-Linear Projection /Kernelizing PCA
- ❑ Correspondence Analyses
- ❑ Use case Analyses using Python

Intuition: Multivariate Statistics

The blind men and the elephant

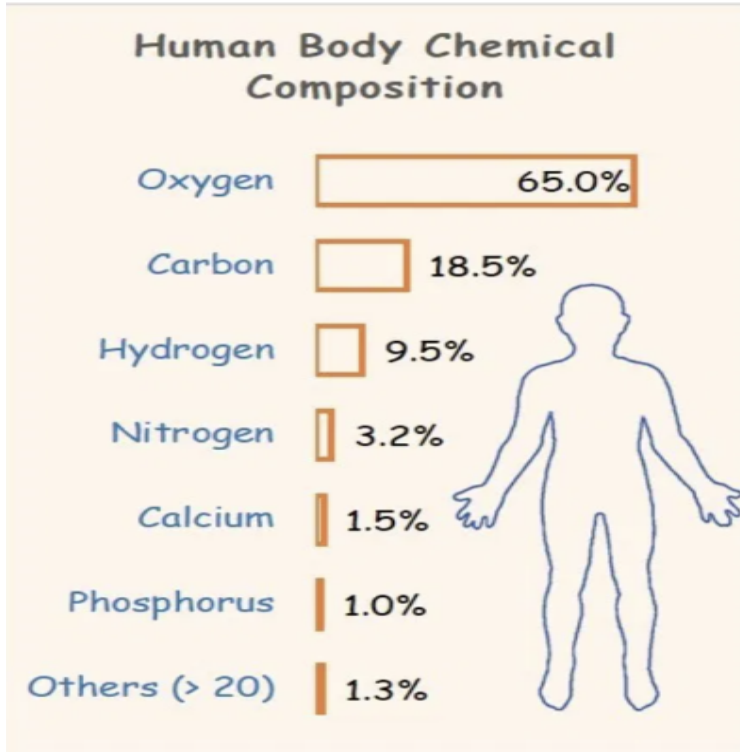


The Blind Men and the Elephant – What is it?

Six blind men who stumble across a strange creature and they try to understand what it is:

The first blind man feels the elephant's body and comes to the conclusion that the creature is in fact a wall, while his friend feels the tusk and declares that the elephant is a spear. Another shakes his head after feeling the trunk of the elephant and claims with some anxiety that it's actually a snake. The fourth blind man, whilst feeling the elephant's leg, states that they are incorrect and that it's indeed a tree. The next man has got hold of the elephant's large ear and announces that it's a fan, whilst the last blind man, who has hold of its tail, declares that his friends are all wrong and that this elephant is in fact a rope.

Intuition: Principal Component Analyses(PCA)

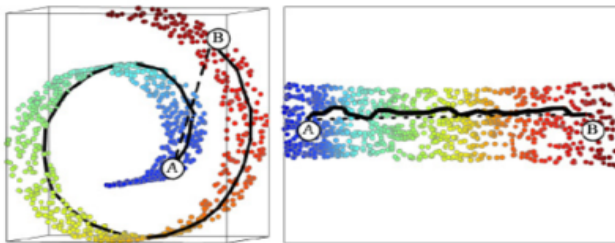
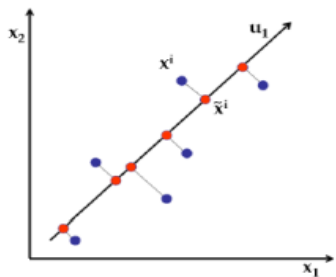


What happens when people die? Where do they go? All of us have pondered this at some point or other. There are several religious and a few non-religious theories to explain events after death. A theory which makes sense to me is: after death, we go back to our fundamental elements. Here, I am not talking about metaphysical or divine elements but chemistry. There are close to 120 known elements including lithium, oxygen, uranium, hydrogen, argon, and carbon. Despite the choice of roughly 120 elements, the human body does not have all of them in the same proportion. Incidentally, six elements, as shown in the chart, constitute close to 99% of the human body. Think of these elements as the principal components of the human body. The remaining ~1% of the human body has a little over 20 other elements. This means, around three-fourth of the elements in the periodic table are completely absent in the human body.

We noticed that about 1% of the human body has formed with ~ 20 elements. **The idea with PCA is that if you remove these 20 elements you will lose just 1% of the essence of the human body. This will also mean that your information load will decline by ~77% (20/26).** These ideas will form the basis of our understanding of principal component analysis

Dimensionality Reduction

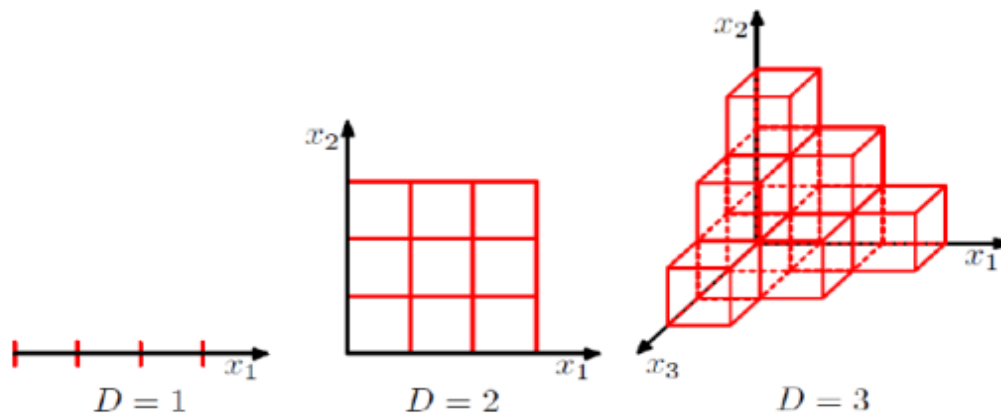
- Usually considered an **unsupervised learning** method
- Used for learning the **low-dimensional structures** in the data



- Also useful for “**feature learning**” or “**representation learning**” (learning a better, often smaller-dimensional, representation of the data), e.g.,
 - Documents using using topic vectors instead of bag-of-words vectors
 - Images using their constituent parts (faces - eigenfaces)
- Can be used for **speeding up** learning algorithms
- Can be used for **data compression**

Curse of Dimensionality

- Exponentially large # of examples required to “fill up” high-dim spaces



- Fewer dimensions \Rightarrow Less chances of overfitting \Rightarrow Better generalization
- Dimensionality reduction is a way to beat the curse of dimensionality

Linear Dimensionality Reduction

- A **projection matrix** $\mathbf{U} = [\mathbf{u}_1 \ \mathbf{u}_2 \ \dots \ \mathbf{u}_K]$ of size $D \times K$ defines K **linear projection directions**, each $\mathbf{u}_k \in \mathbb{R}^D$, for the D dim. data (assume $K < D$)
- Can use \mathbf{U} to transform $\mathbf{x}_n \in \mathbb{R}^D$ into $\mathbf{z}_n \in \mathbb{R}^K$ as shown below

$$\begin{matrix} K \times 1 \\ \mathbf{z}_n \end{matrix} = \begin{matrix} K \times D \\ \mathbf{U}^T \end{matrix} * \begin{matrix} D \times 1 \\ \mathbf{x}_n \end{matrix}$$

- Note that $\mathbf{z}_n = \mathbf{U}^T \mathbf{x}_n = [\mathbf{u}_1^T \mathbf{x}_n \ \mathbf{u}_2^T \mathbf{x}_n \ \dots \ \mathbf{u}_K^T \mathbf{x}_n]$ is a K -dim projection of \mathbf{x}_n
 - $\mathbf{z}_n \in \mathbb{R}^K$ is also called low-dimensional **“embedding”** of $\mathbf{x}_n \in \mathbb{R}^D$

Linear Dimensionality Reduction

- $\mathbf{X} = [\mathbf{x}_1 \ \mathbf{x}_2 \ \dots \ \mathbf{x}_N]$ is $D \times N$ matrix denoting all the N data points
- $\mathbf{Z} = [\mathbf{z}_1 \ \mathbf{z}_2 \ \dots \ \mathbf{z}_N]$ is $K \times N$ matrix denoting **embeddings** of data points
- With this notation, the figure on previous slide can be re-drawn as below

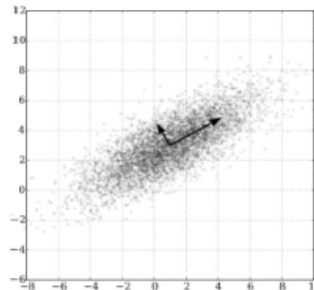
$$\begin{matrix} K \times N \\ \text{Z} \end{matrix} = \begin{matrix} K \times D \\ \mathbf{U}^T \end{matrix} * \begin{matrix} D \times N \\ \mathbf{X} \end{matrix}$$

- How do we learn the “best” projection matrix \mathbf{U} ?
- What criteria should we optimize for when learning \mathbf{U} ?
- Principal Component Analysis (PCA) is an algorithm for doing this

Principal Component Analysis (PCA)

Optum Global Analytics

- A classic linear dim. reduction method (Pearson, 1901; Hotelling, 1930)
- Can be seen as
 - Learning projection directions that capture **maximum variance** in data
 - Learning projection directions that result in **smallest reconstruction error**
- Can also be seen as **changing the basis** in which the data is represented (and transforming the features such that **new features become decorrelated**)

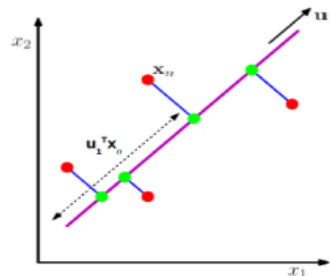


- Also related to other classic methods, e.g., **Factor Analysis** (Spearman, 1904)

PCA as Maximizing Variance

Variance Captured by Projections

- Consider projecting $\mathbf{x}_n \in \mathbb{R}^D$ on a **one-dim subspace** defined by $\mathbf{u}_1 \in \mathbb{R}^D$
- Projection/embedding of \mathbf{x}_n along a one-dim subspace $\mathbf{u}_1 = \mathbf{u}_1^\top \mathbf{x}_n$ (location of the green point along the purple line representing \mathbf{u}_1)

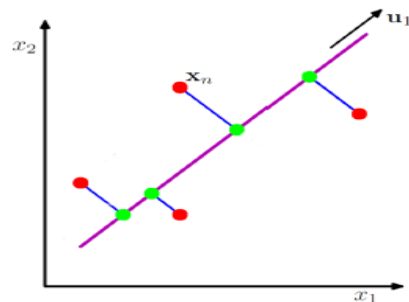


- Mean of projections of all the data: $\frac{1}{N} \sum_{n=1}^N \mathbf{u}_1^\top \mathbf{x}_n = \mathbf{u}_1^\top (\frac{1}{N} \sum_{n=1}^N \mathbf{x}_n) = \mathbf{u}_1^\top \boldsymbol{\mu}$
- **Variance** of the projected data (“spread” of the green points)

$$\frac{1}{N} \sum_{n=1}^N (\mathbf{u}_1^\top \mathbf{x}_n - \mathbf{u}_1^\top \boldsymbol{\mu})^2 = \frac{1}{N} \sum_{n=1}^N \{\mathbf{u}_1^\top (\mathbf{x}_n - \boldsymbol{\mu})\}^2 = \mathbf{u}_1^\top \mathbf{S} \mathbf{u}_1$$

- **S** is the $D \times D$ **data covariance matrix**: $\mathbf{s} = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})(\mathbf{x}_n - \boldsymbol{\mu})^\top$. If data already centered ($\boldsymbol{\mu} = 0$) then $\mathbf{s} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top = \frac{1}{N} \mathbf{X}^\top \mathbf{X}$

Direction of Maximum Variance



- We want \mathbf{u}_1 s.t. the variance of the projected data is maximized

$$\arg \max_{\mathbf{u}_1} \mathbf{u}_1^T \mathbf{S} \mathbf{u}_1$$

- To prevent trivial solution (max var. = infinite), assume $\|\mathbf{u}_1\| = 1 = \mathbf{u}_1^T \mathbf{u}_1$
- We will find \mathbf{u}_1 by solving the following constrained opt. problem

$$\arg \max_{\mathbf{u}_1} \mathbf{u}_1^T \mathbf{S} \mathbf{u}_1 + \lambda_1 (1 - \mathbf{u}_1^T \mathbf{u}_1)$$

where λ_1 is a Lagrange multiplier

Direction of Maximum Variance

- The objective function: $\arg \max_{\mathbf{u}_1} \mathbf{u}_1^\top \mathbf{S} \mathbf{u}_1 + \lambda_1(1 - \mathbf{u}_1^\top \mathbf{u}_1)$
- Taking the derivative w.r.t. \mathbf{u}_1 and setting to zero gives

$$\mathbf{S} \mathbf{u}_1 = \lambda_1 \mathbf{u}_1$$

- Thus \mathbf{u}_1 is an eigenvector of \mathbf{S} (with corresponding eigenvalue λ_1)
- But which of \mathbf{S} 's (D possible) eigenvectors it is?
- Note that since $\mathbf{u}_1^\top \mathbf{u}_1 = 1$, the variance of projected data is

$$\mathbf{u}_1^\top \mathbf{S} \mathbf{u}_1 = \lambda_1$$

- Var. is maximized when \mathbf{u}_1 is the (top) eigenvector with largest eigenvalue
- The top eigenvector \mathbf{u}_1 is also known as the first Principal Component (PC)
- Other directions can also be found likewise (with each being orthogonal to all previous ones) using the eigendecomposition of \mathbf{S} (this is PCA)

Steps in Principal Component Analysis

1. Standardize the d -dimensional dataset.
2. Construct the covariance matrix.
3. Decompose the covariance matrix into its eigenvectors and eigenvalues.
4. Sort the eigenvalues by decreasing order to rank the corresponding eigenvectors.
5. Select k eigenvectors which correspond to the k largest eigenvalues, where k is the dimensionality of the new feature subspace ($k \leq d$).
6. Construct a projection matrix \mathbf{W} from the "top" k eigenvectors.
7. Transform the d -dimensional input dataset \mathbf{X} using the projection matrix \mathbf{W} to obtain the new k -dimensional feature subspace.

PCA as Minimizing the Reconstruction Error

Data as Combination of Basis Vectors

- Assume *complete* orthonormal basis vectors $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_D$, each $\mathbf{u}_d \in \mathbb{R}^D$
- We can represent each data point $\mathbf{x}_n \in \mathbb{R}^D$ **exactly** using this new basis

$$\mathbf{x}_n = \sum_{k=1}^D z_{nk} \mathbf{u}_k$$
$$\begin{bmatrix} x_{n1} \\ x_{n2} \\ \vdots \\ x_{nD} \end{bmatrix} = \begin{bmatrix} | & | \\ \mathbf{u}_1 & \mathbf{u}_2 \\ | & | \end{bmatrix} \begin{bmatrix} | \\ \mathbf{u}_D \\ | \end{bmatrix} * \begin{bmatrix} z_{n1} \\ z_{n2} \\ \vdots \\ z_{nD} \end{bmatrix}$$

- Denoting $\mathbf{z}_n = [z_{n1} \ z_{n2} \ \dots \ z_{nD}]^\top$, $\mathbf{U} = [\mathbf{u}_1 \ \mathbf{u}_2 \ \dots \ \mathbf{u}_D]$, and using $\mathbf{U}^\top \mathbf{U} = \mathbf{I}_D$

$$\mathbf{x}_n = \mathbf{U} \mathbf{z}_n \quad \text{and} \quad \mathbf{z}_n = \mathbf{U}^\top \mathbf{x}_n$$

- Also note that each component of vector \mathbf{z}_n is $z_{nk} = \mathbf{u}_k^\top \mathbf{x}_n$

Reconstruction of Data from Projections

- Reconstruction of \mathbf{x}_n from \mathbf{z}_n will be **exact** if we use all the D basis vectors
- Will be **approximate** if we only use $K < D$ basis vectors: $\mathbf{x}_n \approx \sum_{k=1}^K z_{nk} \mathbf{u}_k$
- Let's use $K = 1$ basis vector. Then the **one-dim embedding** of \mathbf{x}_n is

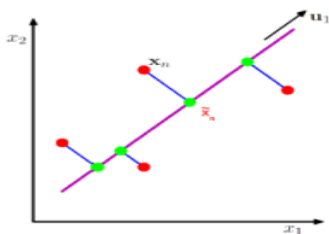
$$\mathbf{z}_n = \mathbf{u}_1^\top \mathbf{x}_n \quad (\text{note: this will just be a scalar})$$

- We can now try “reconstructing” \mathbf{x}_n from its embedding \mathbf{z}_n as follows

$$\tilde{\mathbf{x}}_n = \mathbf{u}_1 \mathbf{z}_n = \mathbf{u}_1 \mathbf{u}_1^\top \mathbf{x}_n$$

- Total error or “loss” in reconstructing all the data points

$$= \sum_{n=1}^N \|\mathbf{x}_n - \tilde{\mathbf{x}}_n\|^2 = \sum_{n=1}^N \|\mathbf{x}_n - \mathbf{u}_1 \mathbf{u}_1^\top \mathbf{x}_n\|^2$$



Direction with Best Reconstruction

- We want to find \mathbf{u}_1 that minimizes the reconstruction error

$$\begin{aligned}L(\mathbf{u}_1) &= \sum_{n=1}^N \|\mathbf{x}_n - \mathbf{u}_1 \mathbf{u}_1^\top \mathbf{x}_n\|^2 \\&= \sum_{n=1}^N \{\mathbf{x}_n^\top \mathbf{x}_n + (\mathbf{u}_1 \mathbf{u}_1^\top \mathbf{x}_n)^\top (\mathbf{u}_1 \mathbf{u}_1^\top \mathbf{x}_n) - 2\mathbf{x}_n^\top \mathbf{u}_1 \mathbf{u}_1^\top \mathbf{x}_n\} \\&= \sum_{n=1}^N -\mathbf{u}_1^\top \mathbf{x}_n \mathbf{x}_n^\top \mathbf{u}_1 \quad (\text{using } \mathbf{u}_1^\top \mathbf{u}_1 = 1 \text{ and ignoring constants w.r.t. } \mathbf{u}_1)\end{aligned}$$

- Thus the problem is equivalent to the following maximization

$$\arg \max_{\mathbf{u}_1: \|\mathbf{u}_1\|^2=1} \mathbf{u}_1^\top \left(\frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top \right) \mathbf{u}_1 = \arg \max_{\mathbf{u}_1: \|\mathbf{u}_1\|^2=1} \mathbf{u}_1^\top \mathbf{S} \mathbf{u}_1$$

where \mathbf{S} is the covariance matrix of the data (data assumed centered)

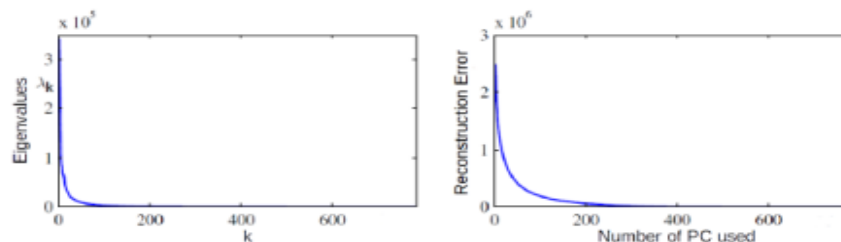
- It's the same objective that we had when we maximized the variance

How many Principal Components to Use?

- Eigenvalue λ_k measures the variance captured by the corresponding PC \mathbf{u}_k
- The “left-over” variance will therefore be

$$\sum_{k=K+1}^D \lambda_k$$

- Can choose K by looking at what fraction of variance is captured by the first K PCs
- Another direct way is to look at the spectrum of the eigenvalues plot, or the plot of reconstruction error vs number of PC



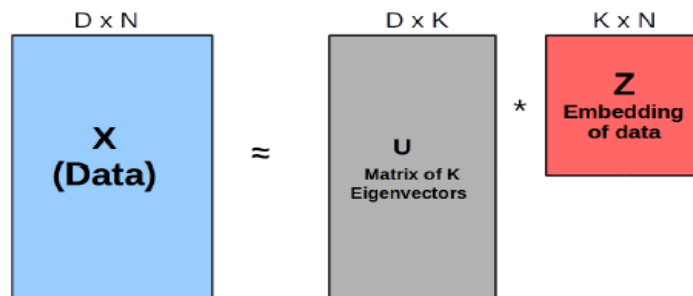
- Can also use other criteria such as AIC/BIC (or more advanced probabilistic approaches to PCA using nonparametric Bayesian methods)

PCA as Matrix Factorization

- Note that PCA represents each \mathbf{x}_n as $\mathbf{x}_n = \mathbf{U}\mathbf{z}_n$
- When using only $K < D$ components, $\mathbf{x}_n \approx \mathbf{U}\mathbf{z}_n$
- For all the N data points, we can write the same as

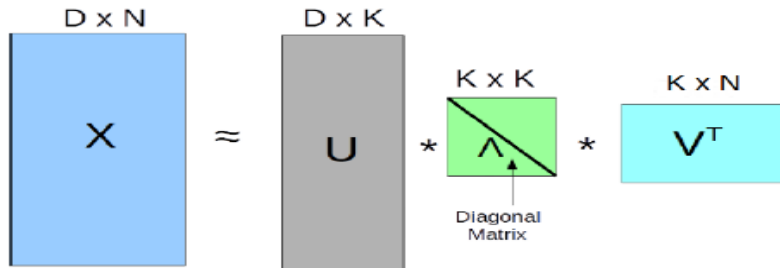
$$\mathbf{X} \approx \mathbf{U}\mathbf{Z}$$

where \mathbf{X} is $D \times N$, \mathbf{U} is $D \times K$ and \mathbf{Z} is $K \times N$



- The above approx. is equivalent to a **low-rank matrix factorization** of \mathbf{X}
 - Also closely related to Singular Value Decomposition (SVD); see next slide

- A rank- K SVD approximates a data matrix \mathbf{X} as follows: $\mathbf{X} \approx \mathbf{U}\mathbf{\Lambda}\mathbf{V}^T$



- \mathbf{U} is $D \times K$ matrix with top K left singular vectors of \mathbf{X}
- $\mathbf{\Lambda}$ is a $K \times K$ diagonal matrix (with top K singular values)
- \mathbf{V} is $N \times K$ matrix with top K right singular vectors of \mathbf{X}
- Rank- K SVD is based on minimizing the reconstruction error

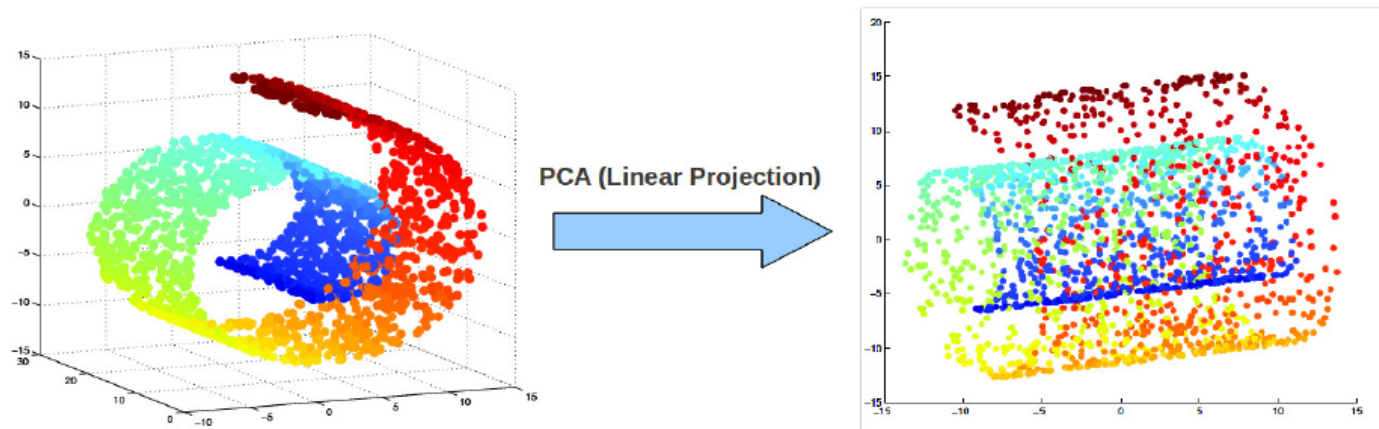
$$||\mathbf{X} - \mathbf{U}\mathbf{\Lambda}\mathbf{V}^T||$$

- PCA is equivalent to the **best rank- K SVD** after centering the data

- A linear projection method
 - Won't work well if data can't be approximated by a linear subspace
 - But PCA can be kernelized easily (Kernel PCA)

Beyond Linear Projections..

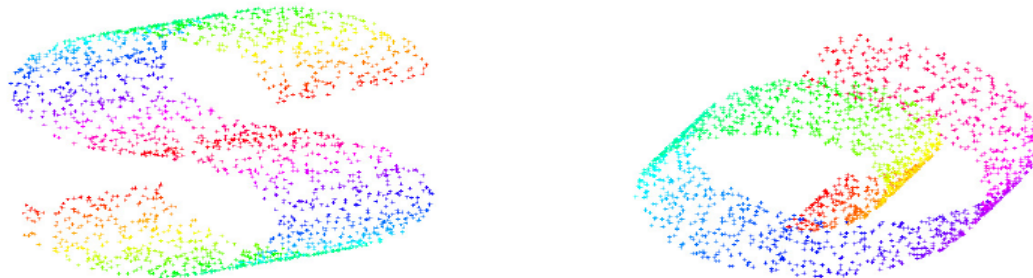
Consider the swiss-roll dataset (points lying close to a manifold):



Linear projection methods (e.g., PCA) can't capture intrinsic nonlinearities

Nonlinear Dimensionality Reduction

- Given: Low-dim. surface **embedded nonlinearly** in high-dim. space
 - Such a structure is called a **Manifold**



- Goal: Recover the low-dimensional surface

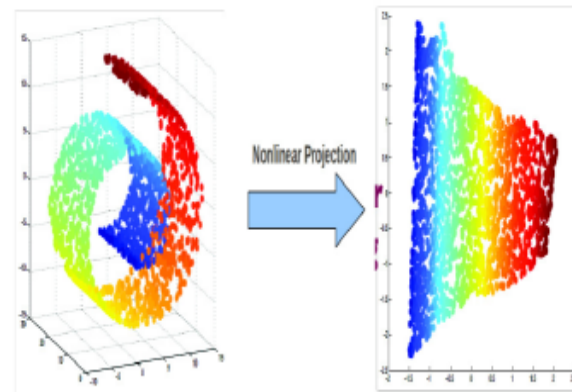


Nonlinear Dimensionality Reduction

Usually two ways of doing this

- **Nonlinearize a linear dimensionality reduction method.**
 - ✓ Kernel PCA (nonlinear PCA)
- **Using manifold based methods that intrinsically preserve nonlinear geometry**
 - ✓ Locally Linear Embedding (LLE)
 - ✓ Isomap
 - ✓ Maximum Variance Unfolding
 - ✓ Laplacian Eigenmaps
 - ✓ And others (tSNE [t-distributed stochastic neighbour embedding], Hessian LLE, etc.)

nonlinear low-dim projection:



Kernel PCA

- Recall PCA: Given N observations $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$, $\forall \mathbf{x}_n \in \mathbb{R}^D$, we define the $D \times D$ **covariance matrix** (assuming centered data $\sum_n \mathbf{x}_n = \mathbf{0}$)

$$\mathbf{S} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top$$

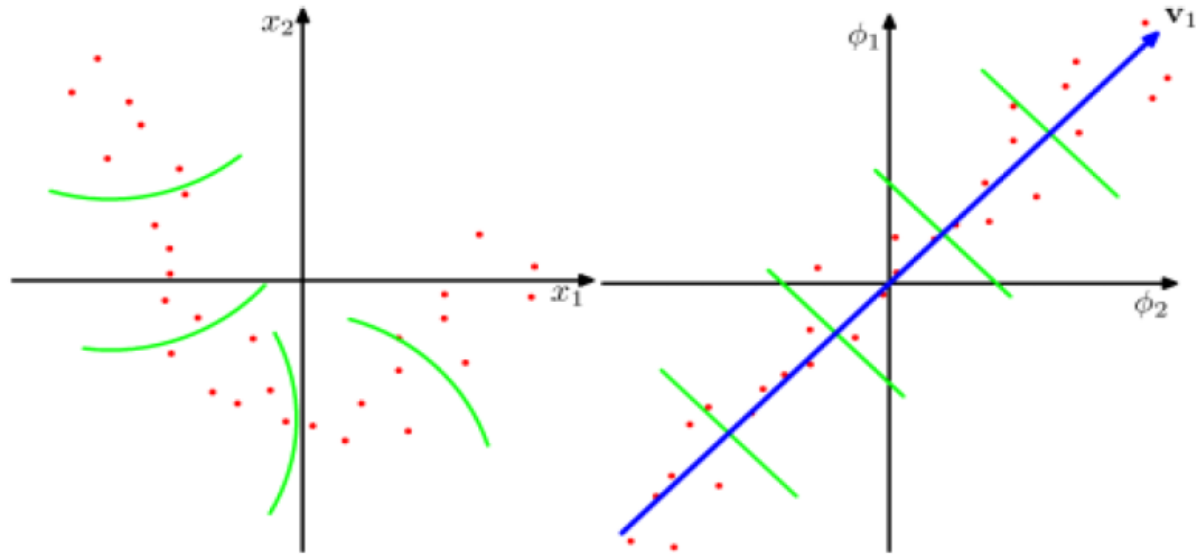
- PCA computes eigenvectors \mathbf{u}_i which satisfy $\mathbf{S}\mathbf{u}_i = \lambda_i \mathbf{u}_i \ \forall i = 1, \dots, D$
- Let's assume a kernel k with associated M dimensional nonlinear map ϕ
- $M \times M$ **covariance matrix in this space** (assume centered data $\sum_n \phi(\mathbf{x}_n) = \mathbf{0}$)

$$\mathbf{C} = \frac{1}{N} \sum_{n=1}^N \phi(\mathbf{x}_n) \phi(\mathbf{x}_n)^\top$$

- Kernel PCA:** Compute eigenvectors \mathbf{v}_i satisfying: $\mathbf{C}\mathbf{v}_i = \lambda_i \mathbf{v}_i \ \forall i = 1, \dots, M$
- We would like to do this **without having to compute \mathbf{C} or $\phi(\mathbf{x}_n)$'s**

** A “**kernel**” is usually used to refer to the **kernel** trick, a method of using a linear classifier to solve a non-linear problem.

Kernel PCA



Right figure: After mapping the data via ϕ , data is now close to a linear subspace

- Goal: Compute eigenvectors \mathbf{v}_i , i.e., $\mathbf{C}\mathbf{v}_i = \lambda_i\mathbf{v}_i$, each \mathbf{v}_i is M dimensional
- Plugging in the expression for \mathbf{C} , we have

$$\frac{1}{N} \sum_{n=1}^N \phi(\mathbf{x}_n) \phi(\mathbf{x}_n)^\top \mathbf{v}_i = \frac{1}{N} \sum_{n=1}^N \phi(\mathbf{x}_n) \phi(\mathbf{x}_n)^\top \mathbf{v}_i = \lambda_i \mathbf{v}_i$$

- Denoting $\alpha_{in} = \frac{1}{\lambda_i N} \phi(\mathbf{x}_n)^\top \mathbf{v}_i$, $\mathbf{v}_i = \sum_{n=1}^N \alpha_{in} \phi(\mathbf{x}_n)$ (also recall Rep. Thm.)
- Thus we can get \mathbf{v}_i by finding $\alpha_i = [\alpha_{i1} \dots \alpha_{iN}]$
- Plugging this back in the eigenvector equation $\mathbf{C}\mathbf{v}_i = \lambda_i \mathbf{v}_i$

$$\frac{1}{N} \sum_{n=1}^N \phi(\mathbf{x}_n) \phi(\mathbf{x}_n)^\top \sum_{m=1}^N \alpha_{im} \phi(\mathbf{x}_m) = \lambda_i \sum_{n=1}^N \alpha_{in} \phi(\mathbf{x}_n)$$

- Pre-multiplying both sides by $\phi(\mathbf{x}_\ell)^\top$ and re-arranging

$$\frac{1}{N} \sum_{n=1}^N \phi(\mathbf{x}_\ell)^\top \phi(\mathbf{x}_n) \sum_{m=1}^N \alpha_{im} \phi(\mathbf{x}_n)^\top \phi(\mathbf{x}_m) = \lambda_i \sum_{n=1}^N \alpha_{in} \phi(\mathbf{x}_\ell)^\top \phi(\mathbf{x}_n)$$

- Using $\phi(\mathbf{x}_n)^\top \phi(\mathbf{x}_m) = k(\mathbf{x}_n, \mathbf{x}_m)$, we get

$$\frac{1}{N} \sum_{n=1}^N k(\mathbf{x}_\ell, \mathbf{x}_n) \sum_{m=1}^N \alpha_{im} k(\mathbf{x}_n, \mathbf{x}_m) = \lambda_i \sum_{n=1}^N \alpha_{in} k(\mathbf{x}_\ell, \mathbf{x}_n)$$

- Define \mathbf{K} as the $N \times N$ **kernel matrix** with $K_{nm} = k(\mathbf{x}_n, \mathbf{x}_m)$
 - \mathbf{K} is the similarity of two examples \mathbf{x}_n and \mathbf{x}_m in the ϕ space
 - ϕ is implicitly defined by kernel function k (which can be, e.g., RBF kernel)
- Define α_i as the $N \times 1$ vector with elements α_{in}
- Using \mathbf{K} and α_i , the eigenvector equation becomes:

$$\mathbf{K}^2 \alpha_i = \lambda_i N \mathbf{K} \alpha_i \quad \Rightarrow \quad \boxed{\mathbf{K} \alpha_i = \lambda_i N \alpha_i}$$

- Thus α_i is an eigenvector of the $N \times N$ kernel matrix \mathbf{K}
- **Note:** Since $\mathbf{v}_i^\top \mathbf{v}_i = 1$ and $\mathbf{v}_i = \sum_{n=1}^N \alpha_{in} \phi(\mathbf{x}_n)$, we have $\alpha_i^\top \mathbf{K} \alpha_i = 1$, which means $\alpha_i^\top \lambda_i N \alpha_i = 1$ and $\alpha_i^\top \alpha_i = 1/(\lambda_i N)$. Thus the original solution with $\alpha_i^\top \alpha_i = 1$ (eigenvector with norm 1) needs to be re-scaled as $\tilde{\alpha}_{in} = \alpha_{in} / \sqrt{\lambda_i N}$

- In PCA, we centered the data before computing the covariance matrix
- For kernel PCA, we need to do the same

$$\tilde{\phi}(\mathbf{x}_n) = \phi(\mathbf{x}_n) - \frac{1}{N} \sum_{\ell=1}^N \phi(\mathbf{x}_\ell)$$

- Each element of the centered kernel matrix

$$\begin{aligned}\tilde{K}_{nm} &= \tilde{\phi}(\mathbf{x}_n)^\top \tilde{\phi}(\mathbf{x}_m) \\ &= \phi(\mathbf{x}_n)^\top \phi(\mathbf{x}_m) - \frac{1}{N} \sum_{\ell=1}^N \phi(\mathbf{x}_n)^\top \phi(\mathbf{x}_\ell) - \frac{1}{N} \sum_{\ell=1}^N \phi(\mathbf{x}_\ell)^\top \phi(\mathbf{x}_m) + \frac{1}{N^2} \sum_{j=1}^N \sum_{\ell=1}^N \phi(\mathbf{x}_j)^\top \phi(\mathbf{x}_\ell) \\ &= k(\mathbf{x}_n, \mathbf{x}_m) - \frac{1}{N} \sum_{\ell=1}^N k(\mathbf{x}_n, \mathbf{x}_\ell) - \frac{1}{N} \sum_{\ell=1}^N k(\mathbf{x}_\ell, \mathbf{x}_m) + \frac{1}{N^2} \sum_{j=1}^N \sum_{\ell=1}^N k(\mathbf{x}_\ell, \mathbf{x}_\ell)\end{aligned}$$

- In matrix notation, the centered $\tilde{\mathbf{K}} = \mathbf{K} - \mathbf{1}_N \mathbf{K} - \mathbf{K} \mathbf{1}_N + \mathbf{1}_N \mathbf{K} \mathbf{1}_N$
- $\mathbf{1}_N$ is the $N \times N$ matrix with every element $= 1/N$
- Eigen-decomposition is then done for the **centered kernel matrix $\tilde{\mathbf{K}}$**

- Suppose $\{\alpha_1, \dots, \alpha_L\}$ are the top L eigenvectors of kernel matrix $\tilde{\mathbf{K}}$
- The L -dimensional KPCA projection $\mathbf{z}_m = [z_{m1}, \dots, z_{mL}]$ of a point \mathbf{x}_m :

$$z_{m\ell} = \phi(\mathbf{x}_m)^\top \mathbf{v}_\ell \quad \forall \ell = 1, \dots, L$$

- Using the definition of \mathbf{v}_ℓ , i.e., $\mathbf{v}_\ell = \sum_{n=1}^N \alpha_{\ell n} \phi(\mathbf{x}_n)$, we have

$$z_{m\ell} = \phi(\mathbf{x}_m)^\top \mathbf{v}_\ell = \sum_{n=1}^N \alpha_{\ell n} k(\mathbf{x}_n, \mathbf{x}_m)$$

- Note: Cost of computing the embeddings scales in N
- Note: For linear kernel, KPCA reduces to PCA (but more efficient if $N < D$)

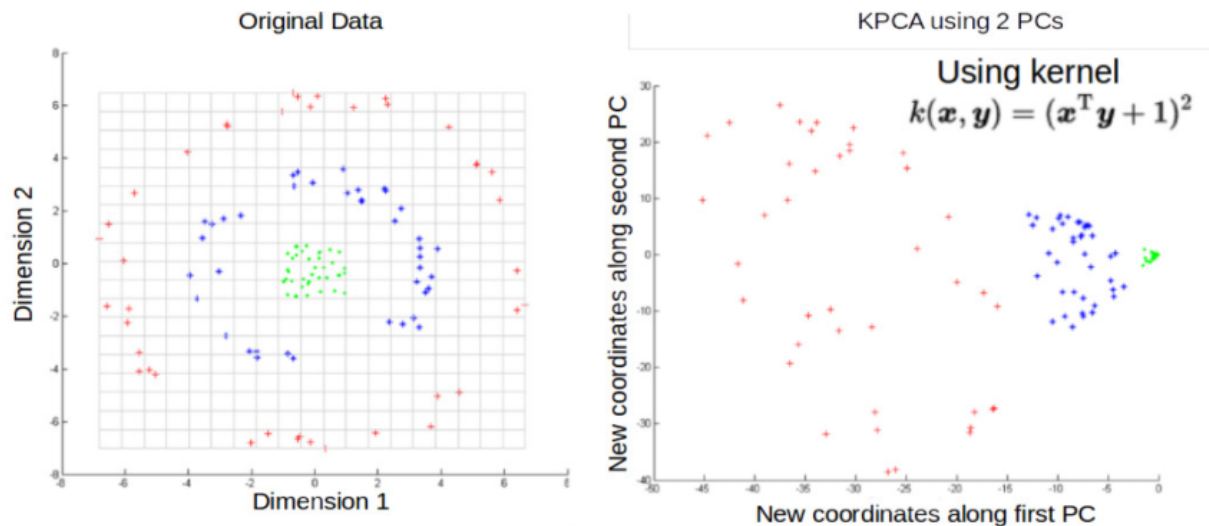
Kernel PCA: Summary of the algorithm

- Construct the $N \times N$ kernel matrix \mathbf{K}
- Center \mathbf{K} as follows $\tilde{\mathbf{K}} = \mathbf{K} - \mathbf{1}_N \mathbf{K} - \mathbf{K} \mathbf{1}_N + \mathbf{1}_N \mathbf{K} \mathbf{1}_N$, where $\mathbf{1}_N$ is an $N \times N$ matrix of all $1/N$
- Do eigen-decomposition of $\tilde{\mathbf{K}}$ and find top L eigenvcs $\alpha_1, \alpha_2, \dots, \alpha_L$ with eigenvals $\lambda_1, \lambda_2, \dots, \lambda_L$
- Re-scale each eigenvector as $\tilde{\alpha}_i = \alpha_i / \sqrt{\lambda_i N}$, $\forall i = 1, \dots, L$
- Finally, compute embedding $\mathbf{z}_m \in \mathbb{R}^L$ of any point \mathbf{x}_m as

$$z_{m\ell} = \sum_{n=1}^N \tilde{\alpha}_{\ell n} k(\mathbf{x}_n, \mathbf{x}_m)$$

- **Note:** For compactness, the $L \times N$ matrix of all L eigenvectors (each is N dimensional) can be written as $\tilde{\alpha} = [\tilde{\alpha}_1 \ \tilde{\alpha}_2 \ \dots \ \tilde{\alpha}_L]^\top$. Thus we can also write embedding $\mathbf{z}_m \in \mathbb{R}^L$ as $\mathbf{z}_m = \tilde{\alpha} \mathbf{k}_m$ where \mathbf{k}_m is $N \times 1$ vector of kernelized similarities of \mathbf{x}_m with all the N training data points

Kernel PCA: An Example



Note that even if we throw away the 2nd PC, we get a good 1-D embedding

Thank you

Optum Global Analytics

Swarupananda Adhikari

swarupananda_adhikari@optum.com

