

Dimensionality Reduction

Part 1: PCA and KPCA

Ng Yen Kaow

Dimensionality Reduction

□ Linear methods

- **PCA** (Principal Component Analysis)
- **cMDS** (Classical Multidimensional Scaling)

□ Non-linear methods

- **KPCA** (Kernel PCA)
- **mMDS** (Metric MDS)
- **Isomap**
- **LLE** (Locally Linear Embedding)
- **t-SNE** (t-distributed Stochastic Neighbor Embedding)
- **UMAP** (Uniform Manifold Approximation and Projection)

Principal Component Analysis

- Let X be an $n \times m$ matrix where each row represents a vector in an m -D space
 - X is like a spreadsheet with features in column
- What do we ideally expect to be the “principal components” of X
 1. The components **form a basis**
 2. The components are **orthogonal**
 3. The **first component accounts for the most variation**, the second component accounts for the most variation after removing the first component, and so on

Principal Component Analysis

- Let X be an $n \times m$ matrix where each row represents a datapoint in an m -D space
 - X represents n datapoints in m -D
- Assume that the rows in X are generated by a random vector $\mathbf{X} \in \mathbb{R}^m$
 - Note the difference between \mathbf{X} and X
 - The theory of PCA is based on X (and its $m \times m$ covariance matrix M)
- For the first component, we want to find unit vector $u \in \mathbb{R}^m$ such that $\text{var}(u^T X)$ is maximized

Principal Component Analysis

- The eigenvector u of the covariance matrix of \mathbf{X} with the largest eigenvalue maximizes $\text{var}(u^T \mathbf{X})$
-

Let $\mathbf{X} \in \mathbb{R}^n$ be a random vector with

- mean $\mu \in \mathbb{R}^n$ and
- covariance matrix $M = \mathbb{E}[(\mathbf{X} - \mu)(\mathbf{X} - \mu)^T]$

For any $u \in \mathbb{R}^n$, the projection of $u^T \mathbf{X}$ has

- $\mathbb{E}[u^T \mathbf{X}] = u^T \mu$ and
- $\text{var}(u^T \mathbf{X}) = \mathbb{E}[(u^T \mathbf{X} - u^T \mu)^2]$
 $= \mathbb{E}[u^T (\mathbf{X} - \mu)(\mathbf{X} - \mu)^T u] = u^T M u$

From min-max theorem, $u^T M u$ is maximized when u is the eigenvector of M with the largest eigenvalue

Principal Component Analysis

- Extend to k principal components, we want
 - k -D subspace of \mathbf{X} that is defined by orthogonal basis $p_1, \dots, p_k \in \mathbb{R}^m$ and displacement $p_0 \in \mathbb{R}^m$
 - Distance from \mathbf{X} to this subspace is minimized
-
- Projection of \mathbf{X} onto subspace is $P^T \mathbf{X} + p_0$, where P is matrix whose rows are p_1, \dots, p_k
 - Squared distance to subspace is $\mathbb{E} \|\mathbf{X} - (P^T \mathbf{X} + p_0)\|^2$
 - By calculus, $p_0 = \mathbb{E} \|\mathbf{X} - P^T \mathbf{X}\| = (1 - P^T)\mu$, hence
$$\mathbb{E} \|\mathbf{X} - (P^T \mathbf{X} + p_0)\|^2 = \mathbb{E} \|\mathbf{X} - \mu\|^2 - \mathbb{E} \|P^T (\mathbf{X} - \mu)\|^2$$
 - To maximize that, need to maximize $\mathbb{E} \|P^T (\mathbf{X} - \mu)\|^2 = \text{var}(P^T \mathbf{X})$
 - Finally, same as in previous slide, p_1, \dots, p_k are eigenvectors of M

Principal Component Analysis

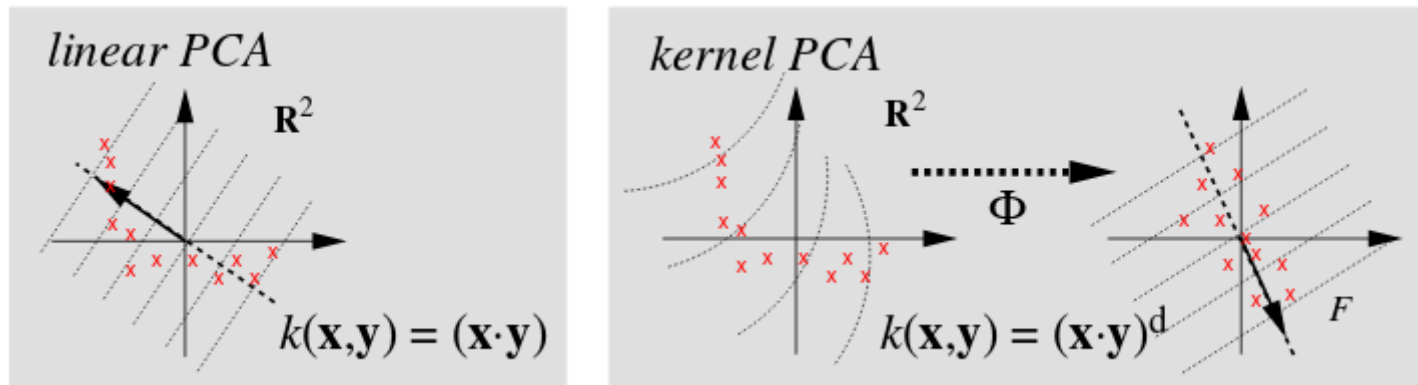
- If X is normalized (values along each feature deducted by the feature's mean), an unbiased estimator of M can be obtained as

$$M = \frac{1}{n} X^T X \quad (\text{or } M = \frac{1}{n} \sum_i x_i'^T x_i)$$

- Since SVD of X eigendecomposes $X^T X$
 - We can solve PCA through either
 1. Eigendecompose M , or
 2. Solve SVD for X

Kernel PCA motivation

- Datapoints that do not lie on a linear manifold in the coordinate space may lie on one after some non-linear transformation ϕ



Scholkopf, Smola, and Muller. Kernel Principal Component Analysis, 1999

- Principal components in the transformed space may be more meaningful in such a case

Kernel PCA idea

- Extract principal components in a ϕ -mapped feature space:
 1. $x' = \phi(x)$ and $X' = [x_1' \quad \dots \quad x_n']^T$
 2. Center X' so that, $\forall i, \sum_i x_i' = 0$
 3. Find covariance matrix, $M' = \frac{1}{n} \sum_i x_i'^T x_i'$
 4. Eigendecompose M'
- But the dimension of x' , $\dim(x') \gg n$ (or even ∞)
 - $\Rightarrow M'$ has large (or even ∞) dimensions
 - \Rightarrow Eigendecomposition leads to **large number of eigenvectors**, each of **large (or infinite) dimensions**

Kernel PCA idea foiled

Problem 1: Large number of eigenvectors

- How many eigenvectors are there actually
 - $\text{rank}(M')$, bounded by the number of datapoints
 - Recall that eigenvectors can be expressed as a linear combination of the datapoints by solving the equations $x'_i = \sum_j \langle x'_i, u_j \rangle u_j$
 - j is bounded by $\text{rank}(M') \Rightarrow$ may be manageable
 - However, working with the system of equations is hard because x'_i and u_j are of...

Problem 2: Large (or ∞) dimensions

Kernel method

- Do not compute $\phi(x_1), \dots, \phi(x_n)$ or eigenvectors of M'
- Limit operations in the mapped space to those that involves only the inner products $\langle x'_i, x'_j \rangle$
 - Sufficient for finding a relation between the eigenvectors of M' and $\phi(x_1), \dots, \phi(x_n)$
 - Given point x , sufficient for finding the projection of $\phi(x)$ on the eigenvectors of M'
- Use a function $K(x_i, x_j) = \langle x'_i, x'_j \rangle$ (called a kernel function) that does not require computing $\phi(x_i)$ or $\phi(x_j)$
 - More on these functions in later slides

Relating eigenvector to x'_1, \dots, x'_n

- Relate eigenvectors of M' with x'_1, \dots, x'_n using a computation that involves only $\langle x'_i, x'_j \rangle$
- Start with the definition of $M' = \frac{1}{n} (\sum_{i=1}^n x_i'^T x_i')$
 - Solving $M'u = \lambda u$ means $(\sum_i x_i'^T x_i')u = n\lambda u$
 - This implies $u = \frac{1}{n\lambda} \sum_i x_i'^T x_i' u$. Since
Proof later
 $x^T x u = x u x^T$, $u = \frac{1}{n\lambda} \sum_i \overbrace{x_i'^T u}^{\text{scalar}} x_i'$

Hence can let $u = \sum_{i=1}^n \alpha_i x_i'^T$ for $\alpha_i \in \mathbb{R}$

- $\alpha_1, \dots, \alpha_n$ **relate eigenvector u to x'_1, \dots, x'_n**
- Place $u^{(r)} = \sum_i \alpha_i^{(r)} x_i'^T$ back in $(\sum_i x_i'^T x_i')u = n\lambda u$
 - Use superscript r to associate α with its corresponding u and λ

Solving $\alpha_1, \dots, \alpha_n$

(Terms in bold cannot be reordered)

System of $\dim(u)$ equations

$$\left(\sum_{i=1}^n \mathbf{x}_i'^T \mathbf{x}_i' \right) \mathbf{u}^{(r)} = n\lambda^{(r)} \mathbf{u}^{(r)}$$

Replace $\mathbf{u}^{(r)}$ with $\sum_j \alpha_j^{(r)} \mathbf{x}_j'^T$

$$\left(\sum_{i=1}^n \mathbf{x}_i'^T \mathbf{x}_i' \right) \sum_{j=1}^n \alpha_j^{(r)} \mathbf{x}_j'^T = n\lambda^{(r)} \sum_{k=1}^n \alpha_k^{(r)} \mathbf{x}_k'^T$$

Reorder

$$\left(\sum_i \mathbf{x}_i'^T \right) \sum_j \overbrace{\mathbf{x}_i' \mathbf{x}_j'^T}^{\text{scalar}} \alpha_j^{(r)} = n\lambda^{(r)} \sum_k \mathbf{x}_k'^T \alpha_k^{(r)}$$

Multiply from the left with \mathbf{x}_l' (equation holds for each l)

$$\left(\sum_i \overbrace{\mathbf{x}_l' \mathbf{x}_i'^T}^{\text{scalar}} \right) \sum_j \mathbf{x}_i' \mathbf{x}_j'^T \alpha_j^{(r)} = n\lambda^{(r)} \sum_k \underbrace{\mathbf{x}_l' \mathbf{x}_k'^T}_{\text{scalar}} \alpha_k^{(r)}$$

System of one equation

Replace $\mathbf{x}_i' \mathbf{x}_j'^T$ with the kernel function

$$\sum_i K(x_l, x_i) \sum_j K(x_i, x_j) \alpha_j^{(r)} = n\lambda^{(r)} \sum_k K(x_l, x_k) \alpha_k^{(r)}$$

Reorder

$$\sum_i \sum_j K(x_l, x_i) K(x_i, x_j) \alpha_j^{(r)} = n\lambda^{(r)} \sum_k K(x_l, x_k) \alpha_k^{(r)}$$

Solving $\alpha_1, \dots, \alpha_n$

$$\sum_i \sum_j K(x_l, x_i) K(x_i, x_j) \alpha_j^{(r)} = n \lambda^{(r)} \sum_k K(x_l, x_k) \alpha_k^{(r)}$$

- **Replace $K(x_i, x_j)$ with a matrix K where $k_{ij} = K(x_i, x_j)$ (K is called a kernel matrix)**

$$\sum_i \sum_j k_{li} k_{ij} \alpha_j^{(r)} = n \lambda^{(r)} \sum_k k_{lk} \alpha_k^{(r)}$$

- For each l this gives one single equation with a linear combination of the variables $\alpha_1^{(r)}, \dots, \alpha_n^{(r)}$

■ e.g. $r = 2$

$$K_r \rightarrow \begin{bmatrix} k_{11} & k_{12} & \dots \\ k_{21} & k_{22} & \dots \\ \vdots & \vdots & \ddots \end{bmatrix} \begin{bmatrix} \overset{K_1^T}{\downarrow} k_{11} \\ k_{21} \\ \vdots \end{bmatrix} \begin{bmatrix} \overset{K_2^T}{\downarrow} k_{12} \\ k_{22} \\ \vdots \end{bmatrix} \begin{bmatrix} \alpha_1^{(r)} \\ \alpha_2^{(r)} \\ \vdots \end{bmatrix} = n \lambda^{(r)} \begin{bmatrix} k_{11} & k_{12} & \dots \\ k_{21} & k_{22} & \dots \\ \vdots & \vdots & \ddots \end{bmatrix} \begin{bmatrix} \alpha_1^{(r)} \\ \alpha_2^{(r)} \\ \vdots \end{bmatrix}$$

$$(k_{21}k_{11} + k_{22}k_{21} + \dots)\alpha_1^{(r)} + (k_{21}k_{12} + k_{22}k_{22} + \dots)\alpha_2^{(r)} + \dots$$

$$= n \lambda^{(r)} (k_{21}\alpha_1^{(r)} + k_{21}\alpha_2^{(r)} + \dots)$$

Solving $\alpha_1, \dots, \alpha_n$

$$K_r \rightarrow \begin{bmatrix} k_{11} & k_{12} & \dots \\ k_{21} & k_{22} & \dots \\ \vdots & \vdots & \ddots \end{bmatrix} \begin{bmatrix} \overset{K_1^T}{\downarrow} k_{11} \\ k_{21} \\ \vdots \end{bmatrix} \begin{bmatrix} \overset{K_2^T}{\downarrow} k_{12} \\ k_{22} \\ \vdots \end{bmatrix} \dots \begin{bmatrix} \alpha_1^{(r)} \\ \alpha_2^{(r)} \\ \vdots \end{bmatrix} = n\lambda^{(r)} \begin{bmatrix} k_{11} & k_{12} & \dots \\ k_{21} & k_{22} & \dots \\ \vdots & \vdots & \ddots \end{bmatrix} \begin{bmatrix} \alpha_1^{(r)} \\ \alpha_2^{(r)} \\ \vdots \end{bmatrix}$$

System of one equation

- Repeat r for 1 to n

$$\begin{bmatrix} k_{11} & k_{12} & \dots \\ k_{21} & k_{22} & \dots \\ \vdots & \vdots & \ddots \end{bmatrix} \begin{bmatrix} k_{11} \\ k_{21} \\ \vdots \end{bmatrix} \begin{bmatrix} k_{12} \\ k_{22} \\ \vdots \end{bmatrix} \dots \begin{bmatrix} \alpha_1^{(r)} \\ \alpha_2^{(r)} \\ \vdots \end{bmatrix} = n\lambda^{(r)} \begin{bmatrix} k_{11} & k_{12} & \dots \\ k_{21} & k_{22} & \dots \\ \vdots & \vdots & \ddots \end{bmatrix} \begin{bmatrix} \alpha_1^{(r)} \\ \alpha_2^{(r)} \\ \vdots \end{bmatrix}$$

System of n equations

- This in matrix notation is

$$\mathbf{K}^2 \boldsymbol{\alpha}^{(r)} = n\lambda^{(r)} \mathbf{K} \boldsymbol{\alpha}^{(r)}$$

- Each $\boldsymbol{\alpha}^{(r)}$ that fulfills the equation relates a eigenvector $\mathbf{u}^{(r)}$ of the covariance matrix M' in terms of the data \mathbf{x}'_i

Solving $\alpha_1, \dots, \alpha_n$

- Removing K from both sides will only affect the $\alpha^{(r)}$ with zero $\lambda^{(r)}$ (proof omitted), leaving the final form of the eigenvalue system

$$K\alpha^{(r)} = n\lambda^{(r)}\alpha^{(r)}$$

- Since $\|\mathbf{u}\| = 1$, we require $n\lambda\alpha^T\alpha = 1 \Rightarrow \|\alpha\|^2 = 1/n\lambda \Rightarrow \|\alpha\| = \sqrt{1/n\lambda}$ Proof later

- However, α^* from the eigendecomposition of K have unit length and eigenvalue $\lambda^* = n\lambda^{(r)}$

To correct for this, $\alpha^{(r)} = \frac{\alpha^*}{\sqrt{n\lambda^{(r)}}} = \frac{\alpha^*}{\sqrt{n\lambda^*/n}} = \frac{\alpha^*}{\sqrt{\lambda^*}}$

- Since $\lambda^{(r)} = \lambda^*/n$, the relative importance of the eigenvectors can be determined from λ^*

Proof for $\|\mathbf{u}\| = 1 \Rightarrow n\lambda\boldsymbol{\alpha}^T\boldsymbol{\alpha} = 1$

□ Since $\|\mathbf{u}\| = 1$

$$\mathbf{u}^T\mathbf{u} = 1$$

$$\left(\sum_i \alpha_i \mathbf{x}_i'^T\right)^T \left(\sum_j \alpha_j \mathbf{x}_j'^T\right) = 1$$

$$\sum_i \sum_j \alpha_i \alpha_j \mathbf{x}_i' \mathbf{x}_j'^T = 1$$

$$\sum_i \sum_j \alpha_i K_{ij} \alpha_j = 1$$

□ Multiply α_i to $\sum_j K_{ij} \alpha_j = n\lambda \sum_k \alpha_k$ gives

$$n\lambda \sum_i \sum_k \alpha_i \alpha_k = \sum_i \sum_j \alpha_i K_{ij} \alpha_j$$

$$n\lambda \sum_i \sum_k \alpha_i \alpha_k = 1$$

$$n\lambda\boldsymbol{\alpha}^T\boldsymbol{\alpha} = 1$$

Proof for $x^T x u = x u x^T$

$$\begin{aligned}(v^T v)u &= \begin{pmatrix} v_1 v_1 & \cdots & v_1 v_n \\ \vdots & \ddots & \vdots \\ v_n v_1 & \cdots & v_n v_n \end{pmatrix} \begin{pmatrix} u_1 \\ \vdots \\ u_n \end{pmatrix} \\&= \begin{pmatrix} v_1 v_1 u_1 + \cdots + v_1 v_n u_n \\ \vdots \\ v_n v_1 u_1 + \cdots + v_n v_n u_n \end{pmatrix} \\&= \begin{pmatrix} (v_1 u_1 + \cdots + v_n u_n) v_1 \\ \vdots \\ (v_1 u_1 + \cdots + v_n u_n) v_n \end{pmatrix} \\&= (v_1 u_1 + \cdots + v_n u_n) \begin{pmatrix} v_1 \\ \vdots \\ v_n \end{pmatrix}\end{aligned}$$

Projection of $\phi(x)$ on u

- Given a point y , the projection of $\phi(y)$ on the eigenvector $u^{(r)}$ of M' can be computed using $\alpha^{(r)}$ as

$$\begin{aligned}\phi(y)u^{(r)} &= \sum_{i=1}^n \alpha_i^{(r)} \phi(y)^T x'_i \\ &= \sum_i \alpha_i^{(r)} K(y, x_i)\end{aligned}$$

- This allows the principal components to be used for clustering the existing datapoints as well as classifying new datapoints into the clusters

Normalizing M'

- X' has been assumed to be normalized so far
- To normalize a matrix X' , subtract every column with the mean of the column:

$$x^* = x' - \frac{1}{n} \sum_{i=1}^n x'_i$$

- The corresponding kernel,

$$\begin{aligned} K^*(x_i, x_j) &= x_i^* x_j^* = \left(x' - \frac{1}{n} \sum_{i=1}^n x'_i \right) \left(x' - \frac{1}{n} \sum_{i=1}^n x'_i \right) \\ &= K(x_i, x_j) - \frac{1}{n} \sum_{k=1}^n K(x_i, x_k) \\ &\quad - \frac{1}{n} \sum_{k=1}^n K(x_j, x_k) + \frac{1}{n^2} \sum_{l,k=1}^n K(x_l, x_k) \end{aligned}$$

Or in matrix notation

$$\mathbf{K}^* = \mathbf{K} - 2\mathbf{1}_{1/n}\mathbf{K} + \mathbf{1}_{1/n}\mathbf{K}\mathbf{1}_{1/n}$$