

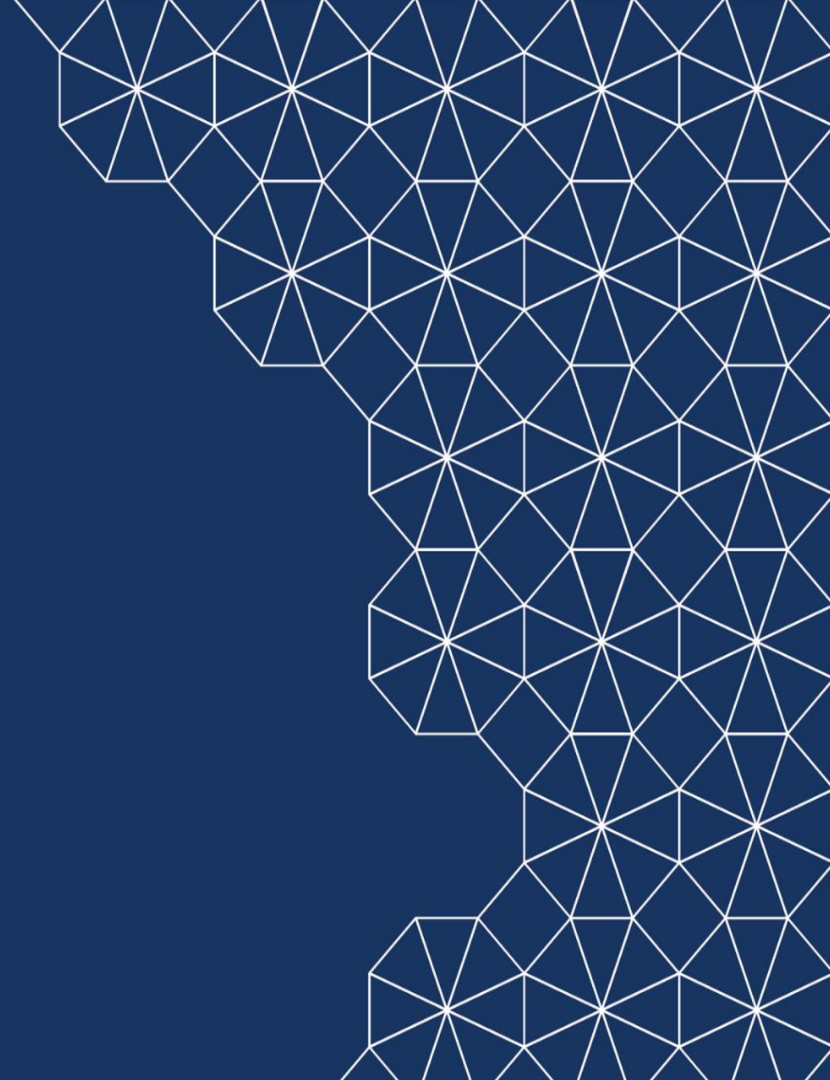
CS 182

Discussion 6

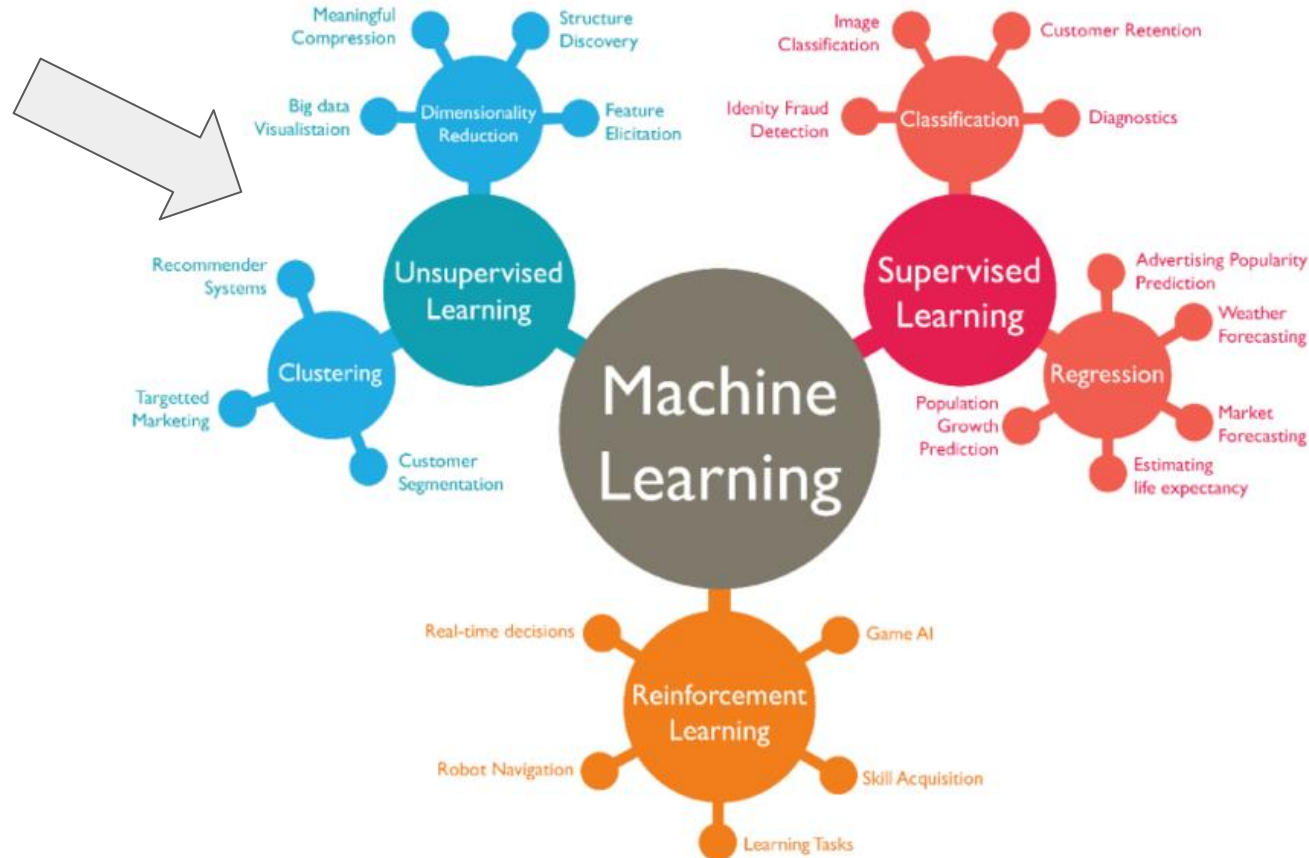
12/14/2023

Yuanming Shao

shaoym1@shanghaitech.edu.cn



Recall: Taxonomy of Machine Learning



Unsupervised Learning

The goal of **unsupervised learning** is to learn structure from **unlabeled** data.

Tasks:

- **Clustering** = partition data into similar points.
 - Find a k -partition of a set of points $X = X_1 \cup X_2 \cup \dots \cup X_k$ such that points within a subset (cluster) of points are “close” to each other.
- **Dimensionality reduction** = project high-dimensional points into a low-dimensional subspace.
 - Find a mapping $f: \mathbb{R}^n \rightarrow \mathbb{R}^d$ with $n > d$, such that if x_i is close to x_j , then $f(x_i)$ is close to $f(x_j)$.
- **Density estimation** = fit a continuous distribution to discrete data.
 - E.g. fitting a Gaussian to data via MLE.

- Let V be vector space over the same field F , $T: V \rightarrow V$ be a linear mapping. We call $v \in V$ an **eigenvector** of T if $v \neq \mathbf{0}$ and
$$T(v) = \lambda v$$

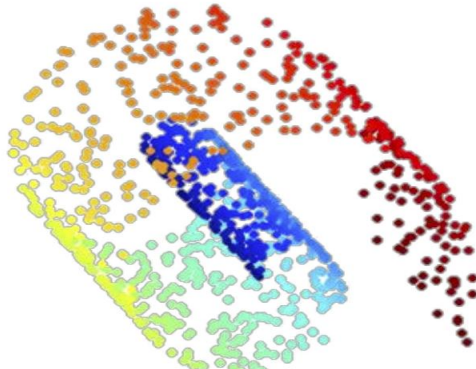
where $\lambda \in F$. We call λ the **eigenvalue** associated with v .

- Let F be a field (\mathbb{R} or \mathbb{C}). A matrix $A \in F^{n \times n}$ can be viewed as a linear mapping $T: F^n \rightarrow F^n$, $T(x) = Ax$. Hence we call λ the eigenvalue (of A) associated with eigenvector (of A) $v \in F^n$ if $v \neq \mathbf{0}$ and
$$Av = \lambda v$$

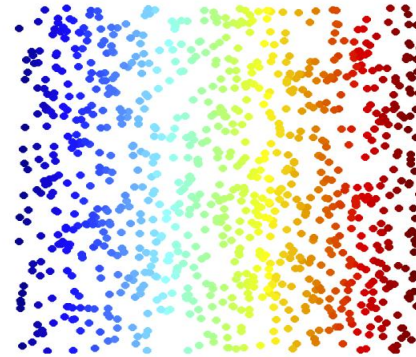
- Exercise: Eigenvector and eigenvalue of

$$A = \begin{bmatrix} -1 & 6 \\ 3 & 2 \end{bmatrix}$$

Dimension Reduction



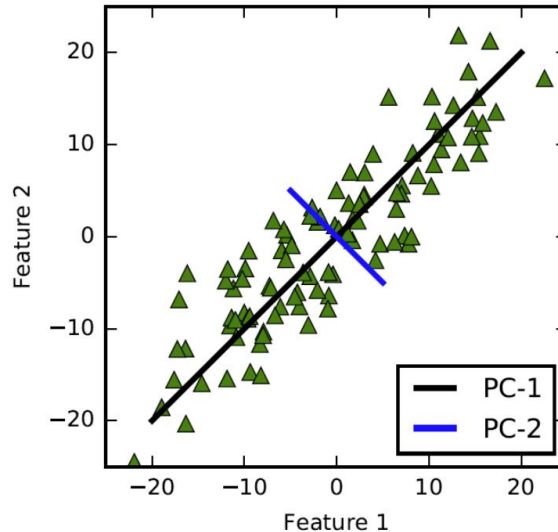
Looks like 3-D



Actually, 2-D

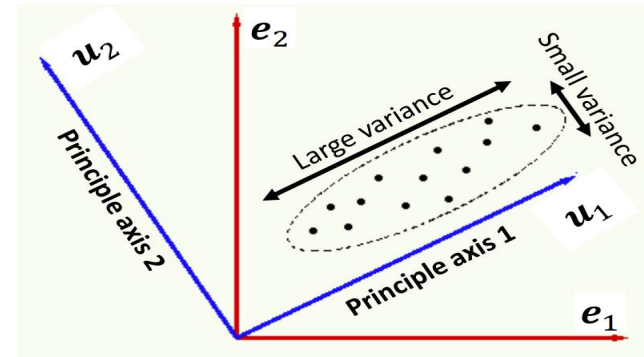
Principal Component Analysis (PCA)

- PCA's target: finding the best lower dimensional sub-space that conveys most of the variance in the original data
- Example: If we were to compress 2-D data to 1-D subspace, then PCA prefers projecting to the **black** line, since it preserves more variance comparing to **blue** line.



Principle Axes

- Objective of PCA: Given data in \mathbb{R}^M , want to *rigidly rotate* the axes to new positions (principle axes) with the following properties:
 - *Ordered such that principle axis 1 has the highest variance, axis 2 has the next highest variance, ..., and axis M has the lowest variance.*
 - *Covariance among each pair of the principal axes is zero.*
- The k 'th **principle component** is the projection to the k 'th principle axis.
- Keep the first $m < M$ principle components for dimensionality reduction.



- Given N data $\mathbf{x}_1, \dots, \mathbf{x}_N \in \mathbb{R}^M$, PCA first computes the covariance matrix for the data

$$\mathbf{\Sigma} = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^T = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T$$

where $\boldsymbol{\mu} \in \mathbb{R}^M$ is the data mean.

- Since $\mathbf{\Sigma}$ is symmetric, $\mathbf{\Sigma}$ can be written as $\mathbf{\Sigma} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T$, where $\mathbf{U} = [\mathbf{u}_1 \dots \mathbf{u}_M]$ is **orthogonal** matrix of eigenvectors (of $\mathbf{\Sigma}$), $\mathbf{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_M)$ is diagonal matrix of the associated eigenvalues arranged in non-ascending order $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_M \geq 0$. (Note that all eigenvalues are non-negative real scalars since $\mathbf{\Sigma}$ is **semi-positive definite**.)
- For data $\mathbf{x} \in \mathbb{R}^M$, compute its 1st principle component as $\mathbf{u}_1^T \mathbf{x}$, 2nd principle component as $\mathbf{u}_2^T \mathbf{x}, \dots$, M 'th principle component as $\mathbf{u}_M^T \mathbf{x}$

Given N data points $X_1, X_2, \dots, X_N \in \mathbb{R}^m$
Compute the Covariance matrix of data $\begin{bmatrix} \end{bmatrix}$

$$\Sigma = \frac{1}{N} \sum_{i=1}^N (X_i - \mu)(X_i - \mu)^T \in \mathbb{R}^{m \times m}$$

$$\mu = \frac{1}{N} \sum_{i=1}^N X_i \quad \begin{bmatrix} \end{bmatrix}^T \quad \end{bmatrix}$$

$\because \Sigma$ is symmetric, it can be diagonalized as:

$$\Sigma = U \Lambda U^T \quad \text{where } U = [u_1, u_2, \dots, u_m] \in \mathbb{R}^{m \times m}$$

$$\Lambda = \begin{bmatrix} \lambda_1 & & \\ & \lambda_m & \end{bmatrix} \in \mathbb{R}^{m \times m}$$

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_m \geq 0$$

Orthogonal matrix:

$U = [\mathbf{u}_1 \dots \mathbf{u}_M] \in \mathbb{R}^{M \times M}$ is an orthogonal matrix if $\mathbf{u}_1, \dots, \mathbf{u}_M$ are orthogonal and have unit length

$$\mathbf{u}_i^T \mathbf{u}_j = \begin{cases} 1 & \text{if } i = j \\ 0 & \text{if } i \neq j \end{cases}$$

That is, $U^T U = I$, namely, $U^{-1} = U^T$.

Positive definite:

$\Sigma \in \mathbb{R}^{M \times M}$ is positive semi-definite if $\mathbf{x}^T \Sigma \mathbf{x} \geq 0$ for all $\mathbf{x} \in \mathbb{R}^M$. If the equality holds only when $\mathbf{x} = \mathbf{0}$, then Σ is positive definite.

$$\bar{\Sigma} = U \Lambda U^T$$

v_i 是什么?

Σ ?

$$y^T \bar{\Sigma} y$$

$$= y^T \frac{1}{n} \sum_{i=1}^n (x_i - \mu)(x_i - \mu)^T y$$

$$= \frac{1}{n} \sum_{i=1}^n y^T (x_i - \mu)(x_i - \mu)^T y$$

$$= \frac{1}{n} \sum_{i=1}^n (y^T (x_i - \mu))^2 \geq 0.$$

$$\underline{\bar{\Sigma} v_i} = U \Lambda \underline{U^T v_i} = U \Lambda \underline{e_i} = U \underline{\lambda_i x}$$

$$= \lambda_i v_i$$

$$\begin{bmatrix} v_1^T \\ v_2^T \\ \vdots \\ v_m^T \end{bmatrix} v_i = \begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} = e_i$$

$$\lambda e_i = \lambda_i e_i$$

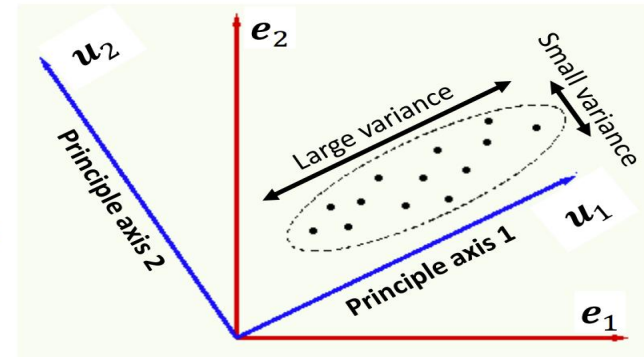
$$0 \leq v_i^T \bar{\Sigma} u_i = v_i^T \lambda_i u_i$$

$$= \lambda_i \|u_i\|^2$$

$$= \lambda_i$$

Principle Axes

- Objective of PCA: Given data in \mathbb{R}^M , want to *rigidly rotate* the axes to new positions (principle axes) with the following properties:
 - *Ordered such that principle axis 1 has the highest variance, axis 2 has the next highest variance, ..., and axis M has the lowest variance.*
 - *Covariance among each pair of the principal axes is zero.*
- The k 'th **principle component** is the projection to the k 'th principle axis.
- Keep the first $m < M$ principle components for dimensionality reduction.



$$\begin{aligned}\bar{\Sigma} &= V \Lambda V^T = [u_1 \quad \dots \quad u_m] \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_m \end{bmatrix} \begin{bmatrix} u_1^T \\ \vdots \\ u_m^T \end{bmatrix} \\ &= \sum_{i=1}^m \lambda_i u_i u_i^T\end{aligned}$$

from M to m

$$\Rightarrow T: x \in \mathbb{R}^m \mapsto \begin{bmatrix} u_1^T x \\ u_2^T x \\ \vdots \\ u_m^T x \end{bmatrix} \in \mathbb{R}^m$$

$$\hat{x}_i = T(x_i) \in \mathbb{R}^{\boxed{m}}$$

$$T(x_i) = \begin{bmatrix} u_1^T x_i \\ \vdots \\ u_m^T x_i \end{bmatrix}$$

$$\hat{X}_i \rightarrow \hat{X}_1, \hat{X}_2, \dots, \hat{X}_N$$

$$\hat{\mu} = \frac{1}{N} \sum_{i=1}^N \hat{X}_i = \frac{1}{N} \sum_{i=1}^N \begin{bmatrix} v_1^T \\ u_m^T \end{bmatrix} X_i = \begin{bmatrix} v_1^T \\ u_m^T \end{bmatrix} \mu.$$

$$(AB)^T = B^T A^T.$$

$$\hat{\Sigma} = \frac{1}{N} \sum_{i=1}^N (\hat{X}_i - \hat{\mu})(\hat{X}_i - \hat{\mu})^T$$

$$= \frac{1}{N} \sum_{i=1}^N \left(\begin{bmatrix} v_1^T \\ u_m^T \end{bmatrix} (X_i - \mu) \right) \left(\begin{bmatrix} v_1^T \\ u_m^T \end{bmatrix} (X_i - \mu) \right)^T.$$

$$= \frac{(X_i - \mu)^T \begin{bmatrix} u_1 & u_2 & \dots & u_m \end{bmatrix}}{N}$$

$$= \begin{bmatrix} v_1^T \\ u_m^T \end{bmatrix} \underbrace{\frac{1}{N} \sum_{i=1}^N (X_i - \mu)(X_i - \mu)^T}_{\Sigma} \begin{bmatrix} u_1 & \dots & u_m \end{bmatrix}$$

$$\Sigma = \begin{bmatrix} u_1^T \\ \vdots \\ u_m^T \end{bmatrix} \underbrace{\Sigma \begin{bmatrix} v_1 & \dots & v_m \end{bmatrix}}_{\substack{m \times m \quad m \times m}} \quad \boxed{m \times m}$$

$$= \begin{bmatrix} v_1^T \\ \vdots \\ v_m^T \end{bmatrix} \begin{bmatrix} \lambda_1 u_1 & \lambda_2 u_2 & \dots & \lambda_m u_m \end{bmatrix}$$

$$= \begin{bmatrix} \cancel{v_1^T \lambda_1 u_1} & v_1^T \lambda_2 u_2 \\ \vdots & \vdots \\ v_m^T \lambda_1 u_1 & v_m^T \lambda_2 u_2 \end{bmatrix} = \begin{bmatrix} \lambda_1 & & \\ & \lambda_2 & \\ & & \ddots \\ & & & \lambda_m \end{bmatrix}$$

Principle Components are Uncorrelated

- The covariance of the k 'th and ℓ 'th principle components of data $\mathbf{x}_1, \dots, \mathbf{x}_N$ is

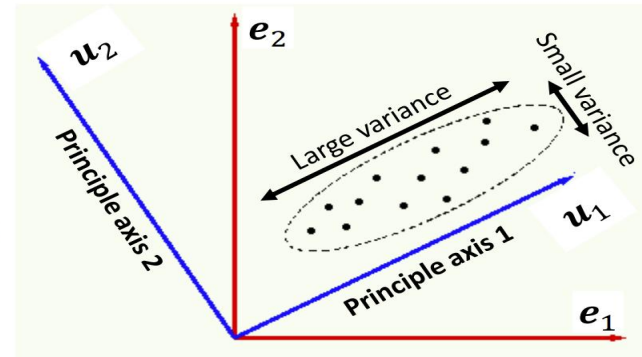
$$\begin{aligned}\frac{1}{N} \sum_{i=1}^N [\mathbf{u}_k^T (\mathbf{x}_i - \boldsymbol{\mu})][\mathbf{u}_\ell^T (\mathbf{x}_i - \boldsymbol{\mu})] &= \frac{1}{N} \sum_{i=1}^N \mathbf{u}_k^T (\mathbf{x}_i - \boldsymbol{\mu}) (\mathbf{x}_i - \boldsymbol{\mu})^T \mathbf{u}_\ell \\ &= \mathbf{u}_k^T \boldsymbol{\Sigma} \mathbf{u}_\ell = \mathbf{u}_k^T \mathbf{U} \boldsymbol{\Lambda} \mathbf{U}^T \mathbf{u}_\ell = \mathbf{e}_k^T \boldsymbol{\Lambda} \mathbf{e}_\ell = \begin{cases} \lambda_k & \text{if } k = \ell \\ 0 & \text{if } k \neq \ell \end{cases}\end{aligned}$$

Therefore

- The variance of the k 'th principle components is λ_k .
⇒ *principle axis 1 has the highest variance, axis 2 has the next highest variance, ..., and axis M has the lowest variance.*
- The covariance of different principle components is zero.
- ⇒ *Covariance among each pair of the principal axes is zero.*

Principle Axes

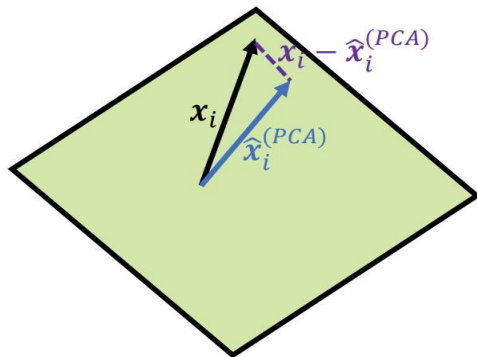
- Objective of PCA: Given data in \mathbb{R}^M , want to *rigidly rotate* the axes to new positions (principle axes) with the following properties:
 - *Ordered such that principle axis 1 has the highest variance, axis 2 has the next highest variance, ..., and axis M has the lowest variance.*
 - *Covariance among each pair of the principal axes is zero.*
- The k 'th **principle component** is the projection to the k 'th principle axis.
- Keep the first $m < M$ principle components for dimensionality reduction.



PCA and Reconstruction Error

WLOG assume zero mean $\frac{1}{N} \sum_{i=1}^N x_i = \mathbf{0}$

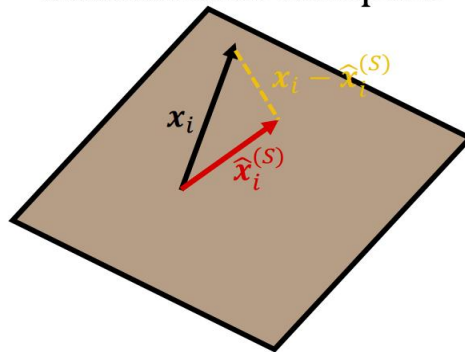
$$S_{PCA} = \text{Span}(\mathbf{u}_1, \dots, \mathbf{u}_m)$$



Variance after projection:

$$\sum_{i=1}^N \|\hat{x}_i^{(PCA)}\|^2 \geq \sum_{i=1}^N \|\hat{x}_i^{(S)}\|^2$$

S: Arbitrary m-dimensional subspace



Mean square error after projection:

$$\sum_{i=1}^N \|x_i - \hat{x}_i^{(PCA)}\|^2 \leq \sum_{i=1}^N \|x_i - \hat{x}_i^{(S)}\|^2$$

$$\hat{x}_i^{(PCA)} = (u_1^T x) u_1 + (u_2^T x) u_2 + \dots + (u_m^T x) u_m$$

Low Rank Approximation

Eckart-Young-Mirsky Theorem:

Let $\mathbf{X} \in \mathbb{R}^{M \times N}$ be a matrix with singular value decomposition $\mathbf{X} = \mathbf{U}\mathbf{D}\mathbf{V}^T$, where $\mathbf{U} \in \mathbb{R}^{M \times M}$, $\mathbf{V} \in \mathbb{R}^{N \times N}$ are orthogonal matrices of left- and right-eigenvectors (of \mathbf{X}), and $\mathbf{D} \in \mathbb{R}^{M \times N}$ is a diagonal matrix of singular values $\sigma_i = D_{ii}$, arranged by their magnitude

$$|\sigma_1| \geq |\sigma_2| \geq \cdots \geq |\sigma_{\min(M,N)}|$$

Let $m \leq \min(M, N)$, then both low rank approximation problems

$$\min_{\hat{\mathbf{X}}} \|\mathbf{X} - \hat{\mathbf{X}}\|_2 \quad \text{subject to } \text{rank}(\hat{\mathbf{X}}) \leq m$$

$$\min_{\hat{\mathbf{X}}} \|\mathbf{X} - \hat{\mathbf{X}}\|_F \quad \text{subject to } \text{rank}(\hat{\mathbf{X}}) \leq m$$

Has optimal solution $\hat{\mathbf{X}} = \sum_{i=1}^m \sigma_i \mathbf{u}_i \mathbf{v}_i^T$. Here \mathbf{u}_i and \mathbf{v}_i denotes the i 'th column in matrices \mathbf{U} , \mathbf{V} , respectively.

$$\text{Let } U = [U_1, U_2, \dots, U_m]$$

$$V = [V_1, V_2, \dots, V_N]$$

For each $i \leq \min\{m, N\}$, we have.

$$\underline{U_i^T X} = (U e_i)^T X = e_i^T U^T U D V^T = e_i^T D V^T = \sigma_i e_i^T V^T = \sigma_i (V e_i)^T$$

$$\underline{X V_i} = X V e_i = U D V^T V e_i = U D e_i = U \sigma_i e_i = \underline{\sigma_i U_i}$$

We call u_i the left eigenvector of X

v_i the right eigenvector of X

WLOG assume zero mean $\boldsymbol{\mu} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i = \mathbf{0}$

$$\mathbf{X} = [\mathbf{x}_1 \ \mathbf{x}_2 \ \dots \ \mathbf{x}_N] = \mathbf{U} \mathbf{D} \mathbf{V}^T$$

$$\boldsymbol{\Sigma} = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^T = \frac{1}{N} \mathbf{X} \mathbf{X}^T = \frac{1}{N} \mathbf{U} \mathbf{D} \mathbf{V}^T \mathbf{V} \mathbf{D}^T \mathbf{U}^T = \mathbf{U} \left(\frac{1}{N} \mathbf{D} \mathbf{D}^T \right) \mathbf{U}^T$$

$$\boldsymbol{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_M) = \frac{1}{N} \mathbf{D} \mathbf{D}^T = \text{diag} \left(\frac{\sigma_1^2}{N}, \dots, \frac{\sigma_{\min(M,N)}^2}{N}, 0, \dots, 0 \right)$$

$$|\sigma_1| \geq |\sigma_2| \geq \dots \text{ implies } \lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_M$$

Projection by PCA: $\hat{\mathbf{x}}_n^{(PCA)} = \sum_{i=1}^m \mathbf{u}_i \mathbf{u}_i^T \mathbf{x}_n$

$$\hat{\mathbf{X}}^{(PCA)} = [\hat{\mathbf{x}}_1^{(PCA)} \ \hat{\mathbf{x}}_2^{(PCA)} \ \dots \ \hat{\mathbf{x}}_N^{(PCA)}] = \sum_{i=1}^m \mathbf{u}_i \mathbf{u}_i^T \mathbf{X} = \sum_{i=1}^m \mathbf{u}_i \mathbf{u}_i^T \mathbf{U} \mathbf{D} \mathbf{V}^T = \sum_{i=1}^m \sigma_i \mathbf{u}_i \mathbf{v}_i^T$$

Projection to S: $\hat{\mathbf{x}}_n^{(S)} \in \mathcal{S}$

$$\hat{\mathbf{X}}^{(S)} = [\hat{\mathbf{x}}_1^{(S)} \ \hat{\mathbf{x}}_2^{(S)} \ \dots \ \hat{\mathbf{x}}_N^{(S)}] \Rightarrow \text{rank}(\hat{\mathbf{X}}^{(S)}) \leq \dim(\mathcal{S}) = m$$

Hence by Eckart-Young-Mirsky Theorem,

$$\|\mathbf{X} - \hat{\mathbf{X}}^{(PCA)}\|_F \leq \|\mathbf{X} - \hat{\mathbf{X}}^{(S)}\|_F, \text{ for all } m\text{-dimensional subspace } \mathcal{S}$$

That is,

$$\sum_{i=1}^N \|\mathbf{x}_i - \hat{\mathbf{x}}_i^{(PCA)}\|^2 \leq \sum_{i=1}^N \|\mathbf{x}_i - \hat{\mathbf{x}}_i^{(S)}\|^2, \text{ for all } m\text{-dimensional subspace } \mathcal{S}$$

Projection by PCA

$$\hat{X}_n^{(PCA)} = \sum_{i=1}^m u_i \underbrace{u_i^T X_n}_{z^{(n)}} \in \mathbb{R}^m$$

$X \in \mathbb{R}^m \rightarrow z \in \mathbb{R}^m$
reconstruct x by z linearly.
say $\hat{X} = Wz$ $W \in \mathbb{R}^{m \times m}$

$$\hat{X}^{(PCA)} = [\hat{X}_1^{(PCA)} \quad \dots \quad \hat{X}_N^{(PCA)}] \in \mathbb{R}^{m \times N}$$

$$= \left[\sum_{i=1}^m u_i u_i^T X_1 \quad \sum_{i=1}^m u_i u_i^T X_2 \quad \dots \quad \sum_{i=1}^m u_i u_i^T X_m \right]$$

$$= \sum_{i=1}^m u_i u_i^T [X_1 \quad X_2 \quad \dots \quad X_N]$$

$$= \sum_{i=1}^m u_i u_i^T X = \sum_{i=1}^m u_i \underbrace{[u_i^T U]_{e_i^T}}_{= \gamma_i} P V^T = \sum_{i=1}^m u_i \sigma_i e_i^T V^T = \sum_{i=1}^m \sigma_i u_i v_i^T$$

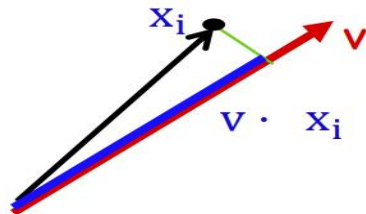
Two Interpretations

So far: **Maximum Variance Subspace**. PCA finds vectors \mathbf{v} such that projections on to the vectors capture maximum variance in the data

$$\frac{1}{n} \sum_{i=1}^n (\mathbf{v}^T \mathbf{x}_i)^2 = \mathbf{v}^T \mathbf{X} \mathbf{X}^T \mathbf{v}$$

Alternative viewpoint: **Minimum Reconstruction Error**. PCA finds vectors \mathbf{v} such that projection on to the vectors yields minimum MSE reconstruction

$$\frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i - (\mathbf{v}^T \mathbf{x}_i) \mathbf{v}\|^2$$



Why? Pythagorean Theorem

E.g., for the first component.

Maximum Variance Direction: 1st PC a vector \mathbf{v} such that projection on to this vector capture maximum variance in the data (out of all possible one dimensional projections)

$$\frac{1}{n} \sum_{i=1}^n (\mathbf{v}^T \mathbf{x}_i)^2 = \mathbf{v}^T \mathbf{X} \mathbf{X}^T \mathbf{v}$$

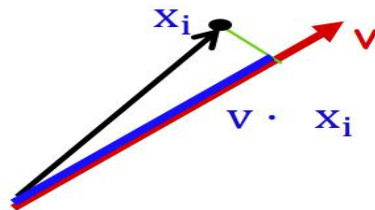
$$\frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i - (\mathbf{v}^T \mathbf{x}_i) \mathbf{v}\|^2$$

Minimum Reconstruction Error: 1st PC a vector \mathbf{v} such that projection on to this vector yields minimum MSE reconstruction

$$\text{blue}^2 + \text{green}^2 = \text{black}^2$$

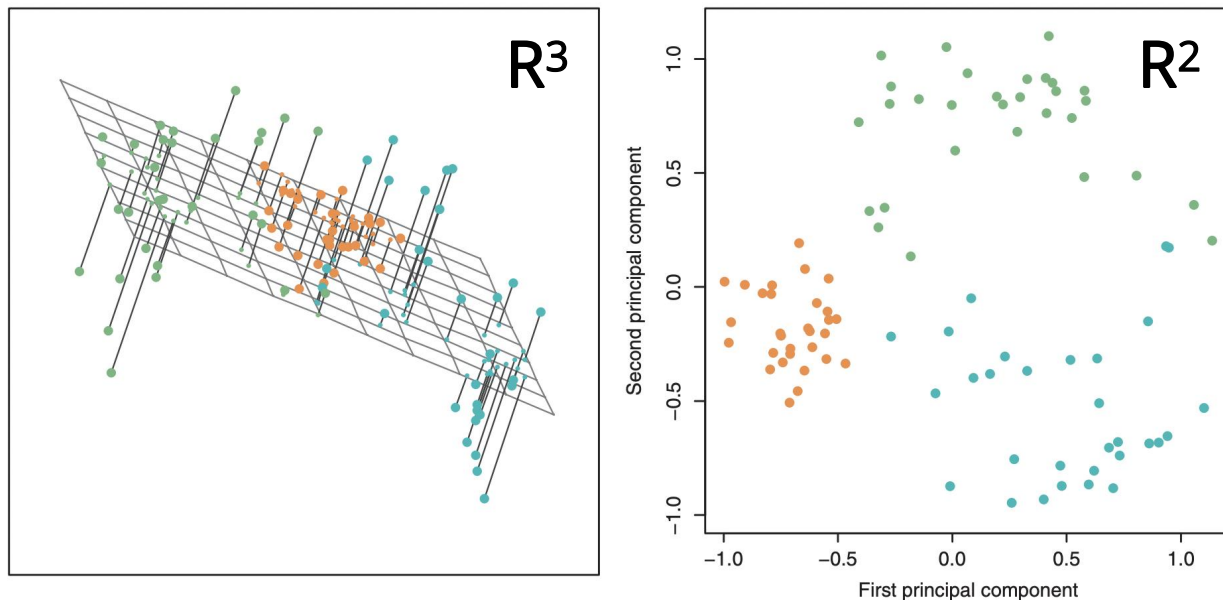
black² is fixed (it's just the data)

So, maximizing blue² is equivalent to minimizing green²



Principal Component Analysis (PCA)

Goal: Given points in \mathbb{R}^d , find the k directions of most variation, effectively performing dimensionality reduction.



Why is PCA used?

- **Lower latency inference** = embed high-dim points into a low-dim space, to be used for many downstream tasks
- **Reduce overfitting** = reducing irrelevant dimensions helps learning algorithms to generalize better
 - More powerful than just feature selection, because we're computing linear combinations of features!
- **Rich representation** = super high-dim points can be represented with a few dimensions, which helps to represent variation among complex things

Concept Check

1. What decomposition allows us to calculate the principal components?
 - a. Eigendecomposition $\Sigma = V\Lambda V^T$.
2. Why are we always guaranteed a solution by eigendecomposition for any full-rank data matrix X ?
 - a. Because we perform eigendecomposition on the **covariance matrix**, which is $\Sigma = (1/n) X^T X$.
 - b. This operation generates a square symmetric matrix, which is the criterion for performing an eigendecomposition.
3. Intuitively, why do we want to maximize sample variance in PCA?
 - a. The goal of PCA is to find as information-rich of a representation as possible. When projecting down, we want to lose as little information as possible.
 - b. When points are close together (low sample variance), we lose information.

Thanks for listening