

Machine Learning & Pattern Recognition

Xin Xin(辛鑫)

xinxin@sdu.edu.cn

<https://xinxin-me.github.io/>

Unsupervised Feature Extraction

- **Principal Component Analysis (PCA)**

What is feature extraction?

- Feature extraction (dimensionality reduction/feature reduction) refers to the mapping of the original **high-dimensional** data into a **low-dimensional** space.
- Criterion for feature reduction can be different based on different problem setting.
 - ✓ Unsupervised setting: minimize the information loss
 - ✓ Supervised setting: maximize the class discrimination

Feature Extraction VS. Feature Selection

- Feature extraction
 - All original features are used.
 - Transformed features are linear combinations of the original features
- Feature selection
 - Only a subset of the original features are used.

Why Feature Extraction?

- **Visualization:** projection of high-dimensional data onto 2D or 3D
- **Data compression:** efficient storage and retrieval
- **Noise removal:** positive effect on accuracy

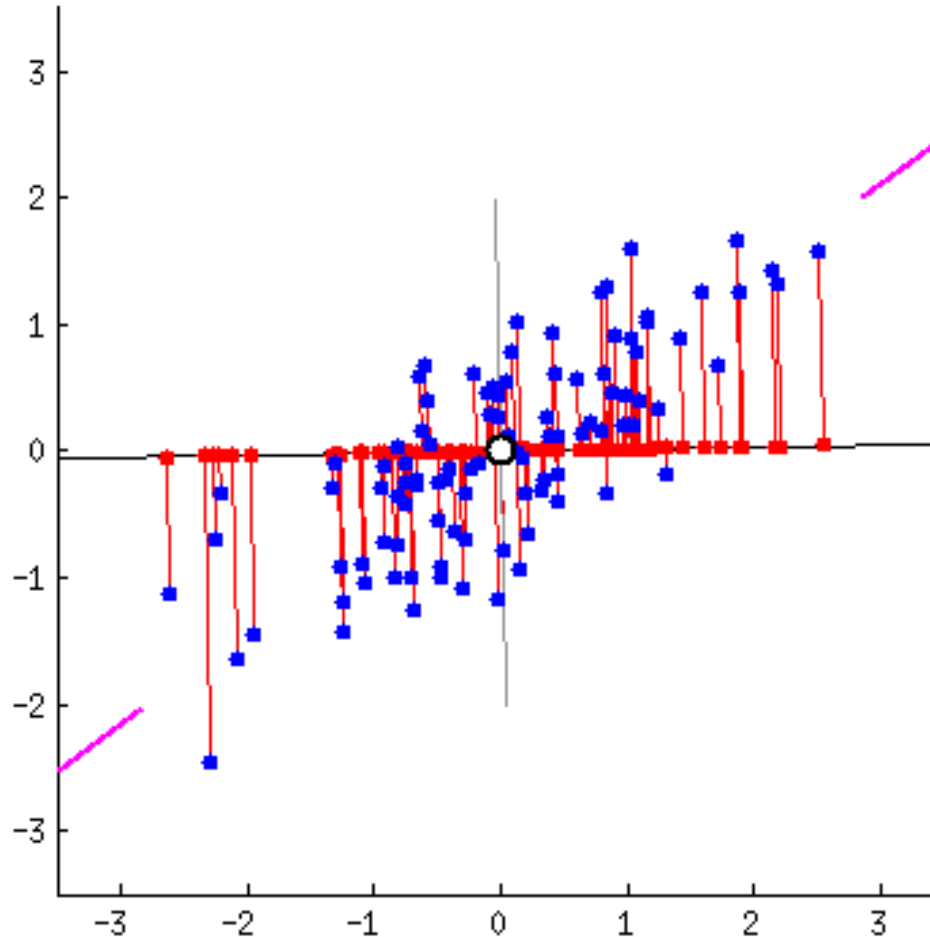
Feature Extraction Algorithms

- **Unsupervised**
 - Principal Component Analysis (**PCA**)
 - Nonnegative Matrix Factorization (NMF)
 - Independent Component Analysis (ICA) [Reading]
- **Supervised**
 - Linear Discriminant Analysis (**LDA**)
 - General Graph Embedding (GE) [Reading]
 - Canonical Correlation Analysis (CCA) [Reading, encouraged]
- **Semi-supervised**
 - Research topic [Further study, encouraged]

Principal Component Analysis (PCA)

PCA is a technique that is widely used for applications such as dimensionality reduction, lossy data compression, feature extraction, and data visualization (Jolliffe, 2002).

Map the original high-dimensional data into a low-dimensional space.



Principal Component Analysis (PCA)

Principal component analysis, or PCA, is a technique that is widely used for applications such as dimensionality reduction, lossy data compression, feature extraction, and data visualization (Jolliffe, 2002).

- Two commonly used definitions of PCA
 - **Maximum variance formulation**
 - The variance of the projected data is maximized.
 - **Minimum-error formulation**
 - Minimizes the average projection cost

Maximum Variance Formulation

- The central idea of PCA is to reduce the dimensionality of a data set consisting of a large number of **interrelated variables**, while retaining as much as possible of the **variation** present in the data set.
- This is achieved by transforming to a new set of variables, the **principal components (PCs)**, which are **uncorrelated**, and are ordered by the fraction of the total information each retains, so that the **first few** retain most of the **variation** present in all of the original variables.

Algebraic Derivation of PCs

Given a sample set of m observations on a vector of d variables

$$\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m\} \in \mathbb{R}^d$$

Define the first PC of the samples by the linear projection $\mathbf{w}_1 \in \mathbb{R}^d$

$$z_{1i} = \mathbf{w}_1^T \mathbf{x}_i = \sum_{k=1}^d w_{1k} x_{ik}, i = 1, \dots, m$$

where $\mathbf{w}_1 = (w_{11}, w_{12}, \dots, w_{1d})^T$ $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{id})^T$

$$\mathbf{z}_1 = \{z_{11}, z_{12}, \dots, z_{1m}\}$$

\mathbf{w}_1 is chosen such that $var[\mathbf{z}_1]$ is maximum.

Algebraic Derivation of PCs

To find \mathbf{w}_1 , first note that

$$\begin{aligned} \text{var}[z_1] &= E[(z_1 - \bar{z}_1)^2] = \frac{1}{m} \sum_{i=1}^m (\mathbf{w}_1^T \mathbf{x}_i - \mathbf{w}_1^T \bar{\mathbf{x}})^2 \\ &= \frac{1}{m} \sum_{i=1}^m \mathbf{w}_1^T (\mathbf{x}_i - \bar{\mathbf{x}}) (\mathbf{x}_i - \bar{\mathbf{x}})^T \mathbf{w}_1 = \mathbf{w}_1^T \mathbf{S} \mathbf{w}_1 \end{aligned}$$

where $\mathbf{S} = \frac{1}{m} \sum_{i=1}^m (\mathbf{x}_i - \bar{\mathbf{x}}) (\mathbf{x}_i - \bar{\mathbf{x}})^T$ is the **covariance** matrix.

$\bar{\mathbf{x}} = \frac{1}{m} \sum_{i=1}^m \mathbf{x}_i$ is the mean.

Covariance and Correlation Coefficient

- For two random variables X and Y ,

Covariance $COV[X, Y] = E[\{X - E[X]\}\{Y - E[Y]\}] = E[XY] - E[X]E[Y]$

The extent to which X and Y vary together.

$$|COV[X, Y]| \leq \sqrt{VAR[X]VAR[Y]}$$

Cauchy–Schwarz inequality.

柯西-施瓦茨不等式

- Correlation coefficient ρ (normalized covariance)

$$\rho(X, Y) = \frac{COV[X, Y]}{\sqrt{VAR[X]VAR[Y]}}$$

Interpretation of The Correlation Coefficient ρ

- Correlation coefficient ρ (normalized covariance)

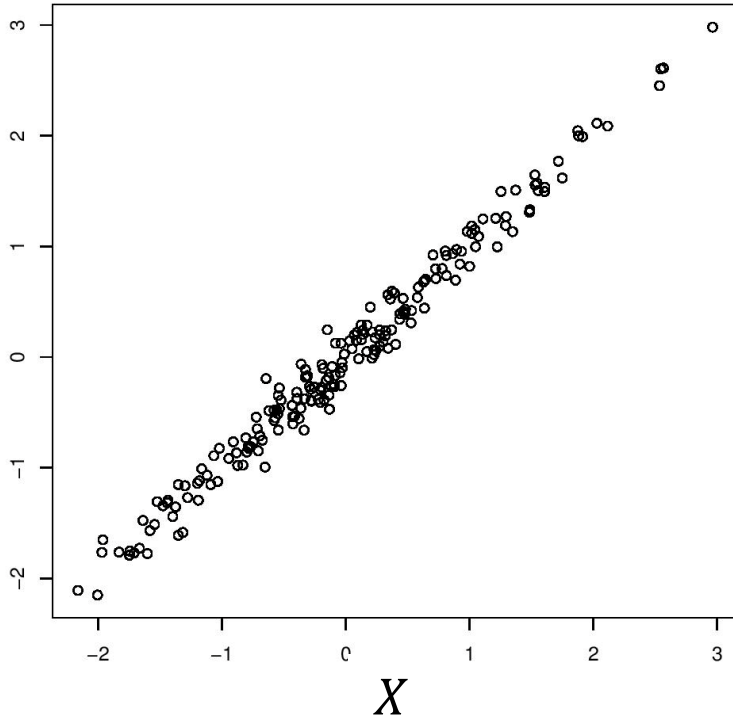
$$\rho(X, Y) = \frac{COV[X, Y]}{\sqrt{VAR[X]VAR[Y]}}$$

- $\rho(X, Y)$ measures the strength and direction of the **linear relationship** between X and Y .
- If X and Y have non-zero variance, then $\rho(X, Y) \in [-1, 1]$.
- Y is a linearly **increasing** function of X if and only if $\rho(X, Y) = 1$
- Y is a linearly **decreasing** function of X if and only if $\rho(X, Y) = -1$
- X and Y are **uncorrelated**, if and only if $\rho(X, Y) = 0$

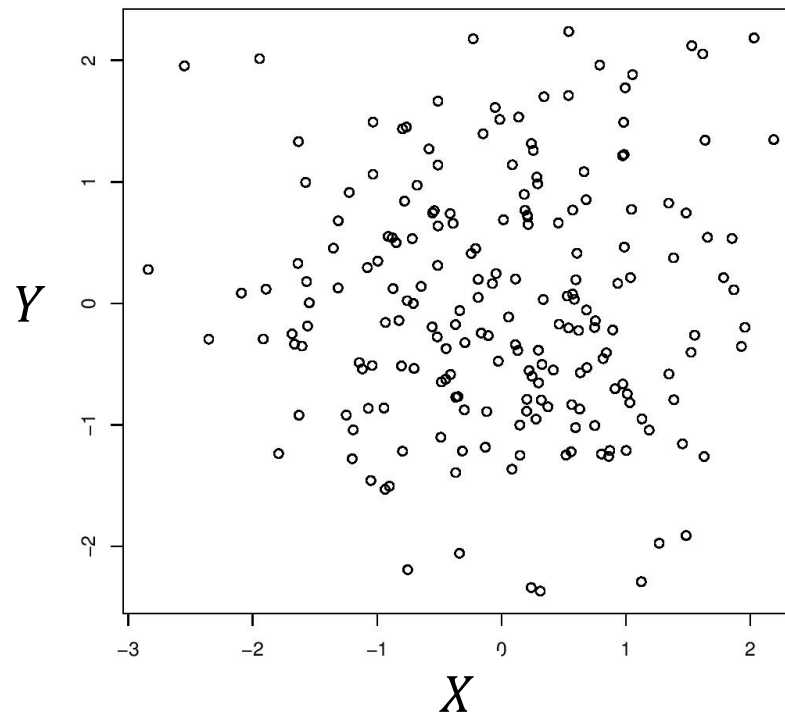
Interpretation of The Correlation Coefficient ρ

- Y is a linearly **increasing** function of X if and only if $\rho(X, Y) = 1$
- Y is a linearly **decreasing** function of X if and only if $\rho(X, Y) = -1$
- X and Y are **uncorrelated**, if and only if $\rho(X, Y) = 0$

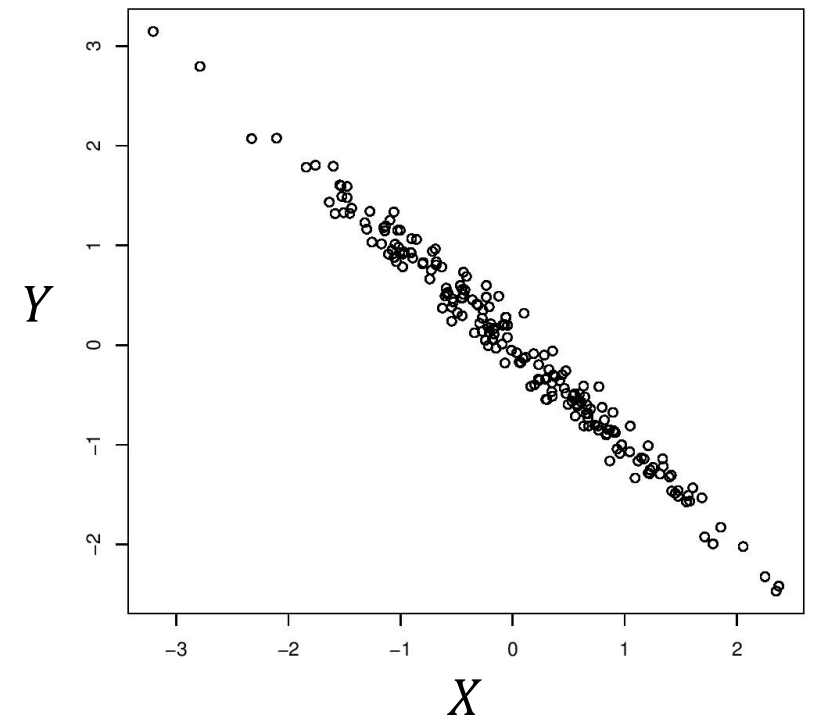
$\rho(X, Y) = 0.99$



$\rho(X, Y) = 0$



$\rho(X, Y) = -0.99$

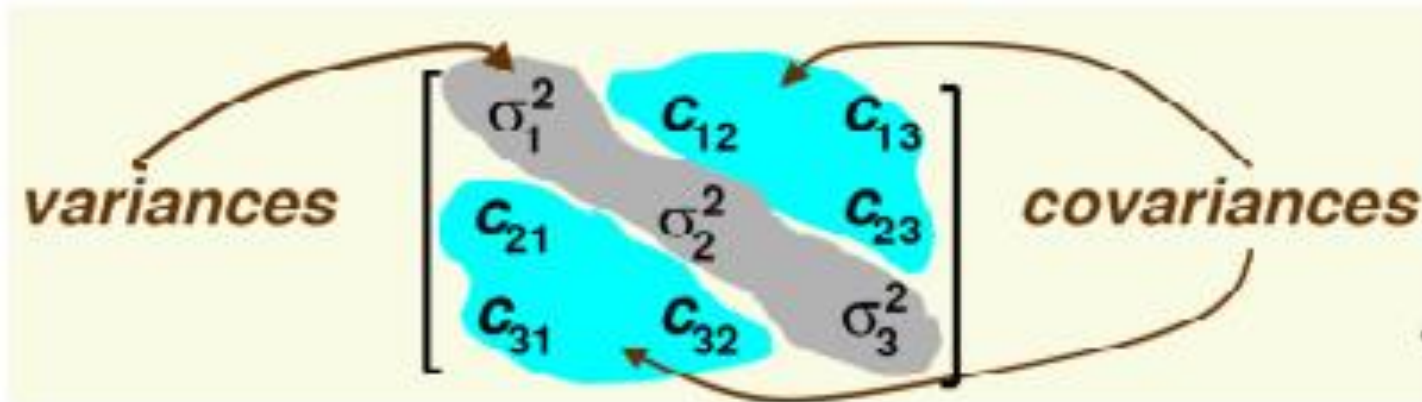


Covariance Matrix

- Given random vector, $\vec{X} = [x_1, x_2, \dots, x_N]^T$, we define,

Mean vector $E[X] = [E[X_1], E[X_2], \dots, E[X_N]]^T = [\mu_1 \mu_2 \dots \mu_N] = \mu$

Covariance matrix
$$COV[X] = \Sigma = E[(X - \mu)(X - \mu)^T]$$
$$= \begin{bmatrix} E[(X_1 - \mu_1)(X_1 - \mu_1)] & \dots & E[(X_1 - \mu_1)(X_N - \mu_N)] \\ & \ddots & \\ E[(X_N - \mu_N)(X_1 - \mu_1)] & \dots & E[(X_N - \mu_N)(X_N - \mu_N)] \end{bmatrix} = \begin{bmatrix} \sigma_1^2 & \dots & c_{1N} \\ & \dots & \\ c_{N1} & \dots & \sigma_N^2 \end{bmatrix}$$



Covariance Matrix

- The covariance matrix indicates the tendency of each pair of dimensions (features) in a random vector to vary together, i.e., to co-vary.

■ Important Properties

- If x_i and x_k tend to increase together, then $c_{ik} > 0$
- If x_i tends to decrease when x_k increases, then $c_{ik} < 0$
- If x_i and x_k are uncorrelated, then $c_{ik} = 0$
- $|c_{ik}| \leq \sigma_i \sigma_k$, where σ_i is the standard deviation of x_i
- $c_{ii} = \sigma_i^2 = \text{VAR}(x_i)$
- **Symmetric:** $c_{ji} = c_{ij}$
- **Positive semi-definite:**
 - Eigenvalues are nonnegative
 - Determinant is nonnegative, $|C| \geq 0$

Covariance Matrix

- You are given the **heights** and **weights** of a certain set of individuals in unknown units. Which one of the following four matrices is the most likely to be the sampled covariance matrix?

(a) $\begin{bmatrix} 1.232 & 0.867 \\ -0.867 & 2.791 \end{bmatrix}$

(b) $\begin{bmatrix} 1.232 & -0.867 \\ -0.867 & 2.791 \end{bmatrix}$

(c) $\begin{bmatrix} 1.232 & 0.867 \\ 0.867 & 2.791 \end{bmatrix}$

(d) $\begin{bmatrix} 1.232 & 3.307 \\ 3.307 & 2.791 \end{bmatrix}$

Algebraic Derivation of PCs

To find \mathbf{w}_1 , first note that

$$\begin{aligned} \text{var}[z_1] &= E[(z_1 - \bar{z}_1)^2] = \frac{1}{m} \sum_{i=1}^m (\mathbf{w}_1^T \mathbf{x}_i - \mathbf{w}_1^T \bar{\mathbf{x}})^2 \\ &= \frac{1}{m} \sum_{i=1}^m \mathbf{w}_1^T (\mathbf{x}_i - \bar{\mathbf{x}}) (\mathbf{x}_i - \bar{\mathbf{x}})^T \mathbf{w}_1 = \mathbf{w}_1^T \mathbf{S} \mathbf{w}_1 \end{aligned}$$

where $\mathbf{S} = \frac{1}{m} \sum_{i=1}^m (\mathbf{x}_i - \bar{\mathbf{x}}) (\mathbf{x}_i - \bar{\mathbf{x}})^T$ is the covariance matrix.

$\bar{\mathbf{x}} = \frac{1}{m} \sum_{i=1}^m \mathbf{x}_i$ is the mean.

The covariance matrix \mathbf{S} is **symmetric**.

- The eigenvectors must be **orthogonal** (正交) to one another.
- The eigenvalues of \mathbf{S} must all be ≥ 0

Algebraic Derivation of PCs

To find \mathbf{w}_1 that maximizes $\text{var}[z_1]$ subject to $\mathbf{w}_1^T \mathbf{w}_1 = 1$

We use the **Lagrange multiplier**, we maximize the following function:

$$L = \mathbf{w}_1^T \mathbf{S} \mathbf{w}_1 + \lambda(1 - \mathbf{w}_1^T \mathbf{w}_1)$$

$$\Rightarrow \frac{\partial L}{\partial \mathbf{w}_1} = \mathbf{S} \mathbf{w}_1 - \lambda \mathbf{w}_1 = 0 \quad \Rightarrow \quad (\mathbf{S} - \lambda \mathbf{I}_d) \mathbf{w}_1 = 0$$

Eigenvectors and Eigenvalues

- **Definition:** \boldsymbol{v} is an eigenvector of matrix $\boldsymbol{A} \in \mathbb{R}^{m \times m}$ if there exists a scalar λ , such that:

$$\boldsymbol{A}\boldsymbol{v} = \lambda\boldsymbol{v} \quad \left\{ \begin{array}{l} \boldsymbol{v}: \text{an eigenvector (nonzero vector)} \\ \lambda: \text{the corresponding eigenvalue} \end{array} \right.$$

- **Computation**

$$\boldsymbol{A}\boldsymbol{v} = \lambda\boldsymbol{v} \quad (\boldsymbol{A} - \lambda\boldsymbol{I})\boldsymbol{v} = \mathbf{0}$$

$$\boldsymbol{v} \neq \mathbf{0} \quad \Rightarrow \quad |\boldsymbol{A} - \lambda\boldsymbol{I}| = 0$$

Eigenvectors and Eigenvalues

- **Definition:** \mathbf{v} is an eigenvector of matrix $\mathbf{A} \in \mathbb{R}^{m \times m}$ if there exists a scalar λ , such that:

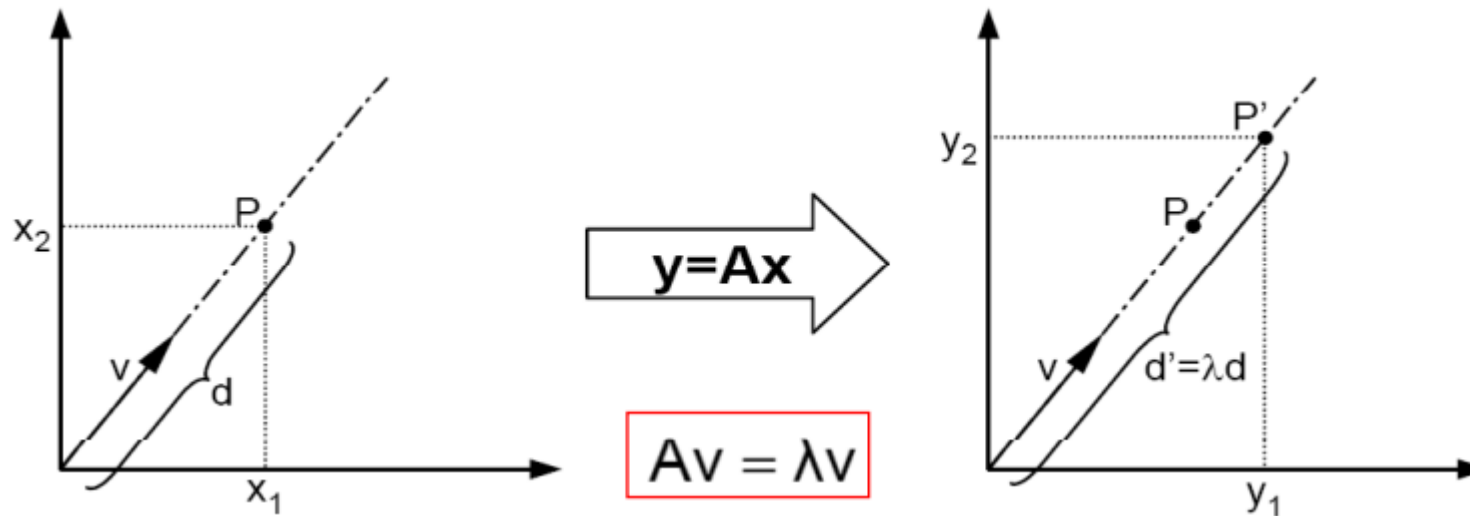
$$\mathbf{A}\mathbf{v} = \lambda\mathbf{v} \quad \left\{ \begin{array}{l} \mathbf{v}: \text{an eigenvector} \\ \lambda: \text{the corresponding eigenvalue} \end{array} \right.$$

- **Note**

- $\text{tr}(\mathbf{A}) = \sum_i \lambda_i$
- $|\mathbf{A}| = \prod_i \lambda_i$
- If λ is an eigenvalue of the matrix \mathbf{A} , then λ^2 is an eigenvalue of \mathbf{A}^2 .
($\mathbf{A}^2 = \mathbf{A}\mathbf{A}$)
- If λ is an eigenvalue of the matrix \mathbf{A} , then λ is an eigenvalue of \mathbf{A}^T .

Eigenvectors and Eigenvalues

- **Intepretation:** an eigenvector represents an **invariant** direction in the vector space.
- Any point lying on the direction defined by v remains on that direction.
- Its magnitude is multiplies by the corresponding eigenvalue λ



Algebraic Derivation of PCs

To find \mathbf{w}_1 that maximizes $\text{var}[z_1]$ subject to $\mathbf{w}_1^T \mathbf{w}_1 = 1$

We use the **Lagrange multiplier**, we maximize the following function:

$$L = \mathbf{w}_1^T \mathbf{S} \mathbf{w}_1 + \lambda(1 - \mathbf{w}_1^T \mathbf{w}_1)$$

$$\Rightarrow \frac{\partial L}{\partial \mathbf{w}_1} = \mathbf{S} \mathbf{w}_1 - \lambda \mathbf{w}_1 = 0 \quad \Rightarrow \quad (\mathbf{S} - \lambda \mathbf{I}_d) \mathbf{w}_1 = 0$$

Therefore \mathbf{w}_1 is an eigenvector of \mathbf{S} , λ is the associated eigenvalue.

Which eigenvector should we choose?

Algebraic Derivation of PCs

If we recognize that the quantity to be maximized

$$L = \mathbf{w}_1^T \mathbf{S} \mathbf{w}_1 = \lambda$$

then we should choose λ to be as big as possible.

So $\lambda = \lambda_1$ is the largest eigenvalue, \mathbf{w}_1 is the corresponding eigenvector.

We have got the 1st PC. Now, let's go to the 2nd PC.

Algebraic Derivation of PCs

The 2nd PC, $\mathbf{w}_2^T \mathbf{x}$ maximize $\mathbf{w}_2^T \mathbf{S} \mathbf{w}_2$ subject to being **uncorrelated** with $\mathbf{w}_1^T \mathbf{x}$

The uncorrelation constraint can be expressed as

$$COV(\mathbf{w}_1^T \mathbf{x}, \mathbf{w}_2^T \mathbf{x}) = \mathbf{w}_1^T \mathbf{S} \mathbf{w}_2 = \mathbf{w}_1^T \lambda_1 \mathbf{w}_2 = \lambda_1 \mathbf{w}_1^T \mathbf{w}_2 = 0$$

Using Lagrange multiplier, we will maximize

$$\mathbf{w}_2^T \mathbf{S} \mathbf{w}_2 + \lambda_2 (1 - \mathbf{w}_2^T \mathbf{w}_2) - \varphi \mathbf{w}_1^T \mathbf{w}_2$$

Algebraic Derivation of PCs

Differentiation w.r.t. \mathbf{w}_2 yields,

$$\frac{d(\mathbf{w}_2^T \mathbf{S} \mathbf{w}_2 + \lambda_2 (1 - \mathbf{w}_2^T \mathbf{w}_2) - \varphi \mathbf{w}_1^T \mathbf{w}_2)}{d\mathbf{w}_2} = \mathbf{0}$$

$$2\mathbf{S}\mathbf{w}_2 - 2\lambda_2\mathbf{w}_2 - \varphi\mathbf{w}_1 = \mathbf{0}$$

If we left multiply \mathbf{w}_1^T ,

$$\mathbf{w}_1^T \mathbf{S} \mathbf{w}_2 - \lambda_2 \mathbf{w}_1^T \mathbf{w}_2 - \frac{\varphi}{2} \mathbf{w}_1^T \mathbf{w}_1 = 0$$

what does it imply?

Algebraic Derivation of PCs

$$\mathbf{w}_1^T \mathbf{S} \mathbf{w}_2 - \lambda_2 \mathbf{w}_1^T \mathbf{w}_2 - \frac{\varphi}{2} \mathbf{w}_1^T \mathbf{w}_1 = 0$$

$$0 - 0 - \frac{\varphi}{2} 1 = 0$$

φ must be **zero** and then we re-check the derivative,

$$\mathbf{S} \mathbf{w}_2 - \lambda_2 \mathbf{w}_2 - \varphi \mathbf{w}_1 = \mathbf{0} \Rightarrow \mathbf{S} \mathbf{w}_2 - \lambda_2 \mathbf{w}_2 = \mathbf{0}$$

The same strategy of choosing \mathbf{w}_2 to be the **eigenvector** associated with the **second largest eigenvalue** λ_2 yields the second PC.

Algebraic Derivation of PCs

This process can be repeated for $k = 1, \dots, p$ yielding up to p different eigenvectors of \mathbf{S} along with the corresponding eigenvalues $\lambda_1, \dots, \lambda_p$

The variance of each of the PCs are given by

$$\text{var}[z_k] = \text{var}[\mathbf{w}_k^T \mathbf{x}] = \lambda_k, k = 1, \dots, p$$

Algebraic Derivation of PCs

- Main steps for computing PCs:
 - Form the covariance matrix \mathbf{S} .
 - Compute its eigenvectors: $\{\mathbf{w}_i\}_{i=1}^d$
 - The first p eigenvectors $\{\mathbf{w}_i\}_{i=1}^p$ form the p PCs
 - The transformation \mathbf{G} consists of the p PCs

$$\mathbf{G} \leftarrow [\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_p] \in \mathbb{R}^{d \times p}$$

$$\mathbf{y} = \mathbf{G}^T \mathbf{x} \in \mathbb{R}^p$$

Principal Component Analysis (PCA)

Principal component analysis, or PCA, is a technique that is widely used for applications such as dimensionality reduction, lossy data compression, feature extraction, and data visualization (Jolliffe, 2002).

- Two commonly used definitions of PCA
 - Maximum variance formulation
 - The variance of the projected data is maximized.
 - Minimum-error formulation
 - Minimizes the average projection cost

Minimum-error Formulation

■ Basis

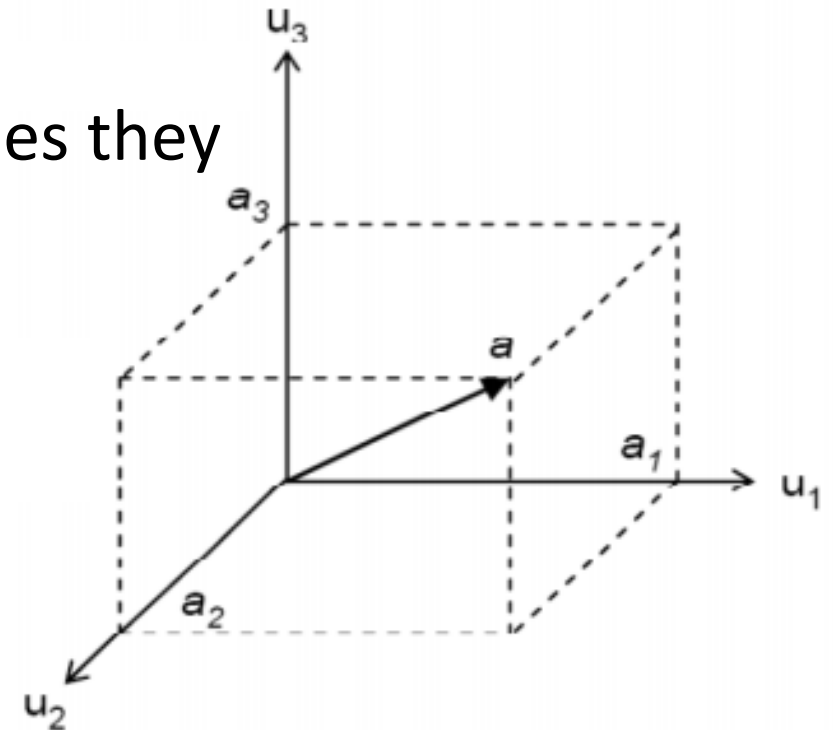
- A set of vectors $\{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_n\}$ are called a **basis** for a vector space \mathbb{R}^n if any vector $\mathbf{x} \in \mathbb{R}^n$ can be written as a linear combination of $\{\mathbf{u}_i\}$

$$\mathbf{x} = a_1 \mathbf{u}_1 + a_2 \mathbf{u}_2 + \dots + a_n \mathbf{u}_n$$

- $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_n$ are linearly **independent** implies they form a basis and vice versa.

$$a_1 \mathbf{u}_1 + a_2 \mathbf{u}_2 + \dots + a_n \mathbf{u}_n = \mathbf{0} \Rightarrow a_k = 0$$

- A basis $\{\mathbf{u}_i\}$ is **orthonormal** if
 - Basis vectors are pairwise orthogonal
 - Have unit length, i.e., $|\mathbf{u}_i| = 1$.



Minimum-error Formulation

Here, we first introduce a complete **orthonormal** set of basis vectors $\{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_d\} \in \mathbb{R}^d$, where

$$\mathbf{u}_i^T \mathbf{u}_j = \delta_{ij} \text{ (0 or 1)}$$

Each data point can be represented by a linear combination of the basis vectors:

$$\mathbf{x}_n = \sum_{i=1}^d \alpha_{ni} \mathbf{u}_i$$

$$\alpha_{ni} = \mathbf{x}_n^T \mathbf{u}_i$$

$$\mathbf{x}_n = \sum_{i=1}^d (\mathbf{x}_n^T \mathbf{u}_i) \mathbf{u}_i$$

Minimum-error Formulation

- Our goal is to approximate $\mathbf{x}_n \in \mathbb{R}^d$ with $\tilde{\mathbf{x}}_n \in \mathbb{R}^p$, $p < d$ (projection onto a lower-dimensional subspace).
- The p -D subspace can be represented, without loss of generality, by the first p of the basis vectors, and we approximate each sample \mathbf{x}_n by

$$\tilde{\mathbf{x}}_n = \sum_{i=1}^p z_{ni} \mathbf{u}_i + \sum_{i=p+1}^d b_i \mathbf{u}_i$$

- $\{z_{ni}\}$ depend on the particular data point.
- $\{b_i\}$ are constants that are the same for all data points.

We aim to minimize the loss introduced by the dimensionality reduction, i.e., the squared distance between the original \mathbf{x}_n and approximation $\tilde{\mathbf{x}}_n$

$$J = \frac{1}{m} \sum_{n=1}^m \|\mathbf{x}_n - \tilde{\mathbf{x}}_n\|^2$$

Minimum-error Formulation

To minimize

$$J = \frac{1}{m} \sum_{n=1}^m \|\mathbf{x}_n - \tilde{\mathbf{x}}_n\|^2 \quad \tilde{\mathbf{x}}_n = \sum_{i=1}^p z_{ni} \mathbf{u}_i + \sum_{i=p+1}^d b_i \mathbf{u}_i$$

- Taking the derivative w.r.t. z_{nj} and setting to zero, we have

$$z_{nj} = \mathbf{x}_n^T \mathbf{u}_j, \quad j = 1, \dots, p$$

- Taking the derivative w.r.t. b_j and setting to zero, we have

$$b_j = \bar{\mathbf{x}}^T \mathbf{u}_j, \quad j = p + 1, \dots, d$$

Leaving for
Your
homework.

- Substituting z_{nj} and b_j , and make use of $\mathbf{x}_n = \sum_{i=1}^d (\mathbf{x}_n^T \mathbf{u}_i) \mathbf{u}_i$, we have

$$\mathbf{x}_n - \tilde{\mathbf{x}}_n = \sum_{i=p+1}^d \{(\mathbf{x}_n - \bar{\mathbf{x}})^T \mathbf{u}_i\} \mathbf{u}_i$$

$$x_n - \hat{x}_n = x_n - \sum_{i=1}^p z_{ni} u_i - \sum_{i=p+1}^d b_i u_i.$$

$$= x_n - \sum_{i=1}^p x_n u_i u_i - \sum_{i=p+1}^d \bar{x}^T u_i u_i$$

$$= \sum_{i=1}^d x_n^T u_i u_i - \sum_{i=1}^p x_n u_i u_i - \sum_{i=p+1}^d \bar{x}^T u_i u_i$$

$$= \underbrace{\sum_{i=1}^p x_n^T u_i u_i + \sum_{i=p+1}^d x_n^T u_i u_i - \sum_{i=1}^p x_n u_i u_i - \sum_{i=p+1}^d \bar{x}^T u_i u_i.}_{\quad}$$

$$= \sum_{i=p+1}^d (x_n^T - \bar{x}^T) u_i u_i = \sum_{i=p+1}^d (x_n - \bar{x})^T u_i u_i.$$

Minimum-error Formulation $\tilde{\mathbf{x}}_n = \sum_{i=1}^p z_{ni} \mathbf{u}_i + \sum_{i=p+1}^d b_i \mathbf{u}_i$

$$\mathbf{x}_n - \tilde{\mathbf{x}}_n = \sum_{i=p+1}^d \{(\mathbf{x}_n - \bar{\mathbf{x}})^T \mathbf{u}_i\} \mathbf{u}_i$$

- $\mathbf{x}_n - \tilde{\mathbf{x}}_n$ lies in the space orthogonal to the principal subspace, **why?**
- Then

$$\min J = \frac{1}{m} \sum_{n=1}^m \|\mathbf{x}_n - \tilde{\mathbf{x}}_n\|^2 = \frac{1}{m} \sum_{n=1}^m \sum_{i=p+1}^d (\mathbf{x}_n^T \mathbf{u}_i - \bar{\mathbf{x}}^T \mathbf{u}_i)^2 = \sum_{i=p+1}^d \mathbf{u}_i^T \mathbf{S} \mathbf{u}_i$$

- Let us try to obtain some intuition about the result by considering the case of $d = 2$ and $p = 1$.

Minimum-error Formulation

$$J = \sum_{i=p+1}^d \mathbf{u}_i^T \mathbf{S} \mathbf{u}_i, d = 2, p = 1$$

- We choose \mathbf{u}_2 to minimize $J = \mathbf{u}_2^T \mathbf{S} \mathbf{u}_2$.
 - Using Lagrange multiplier, we know \mathbf{u}_2 should be the eigenvector corresponding to the **smaller** of the two eigenvalues.
 - Thus we should choose the principle subspace to be aligned with the eigenvector having the **larger** eigenvalue.
-
- For the case when the eigenvalues are equal, any choice of principal direction will give rise to the same value of J .

Minimum-error Formulation

- Generally, the solution to the minimization of J is obtained by choosing the $\{\mathbf{u}_i\}$ to be the eigenvectors of the covariance matrix given by $\mathbf{S}\mathbf{u}_i = \lambda_i \mathbf{u}_i, i = 1, \dots, d$.
- The corresponding value of the loss measure is then given by

$$J = \sum_{i=p+1}^d \lambda_i$$

- Eigenvectors defining the **principal** subspace are those corresponding to the p **largest** eigenvalues.

Eigenvalue Decomposition

- Given a **square** matrix with m linearly **independent** eigenvectors $A \in \mathbb{R}^{m \times m}$, we have an eigenvalue decomposition

$$AV = V\Lambda \quad A = V\Lambda V^{-1}$$

- Note

- Columns of V are eigenvectors of A
- Diagonal elements of Λ are eigenvalues of A

$$\Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_m), \lambda_i \geq \lambda_{i+1}$$

Eigenvalue Decomposition

■ Note

- If A is non-singular
 - All eigenvalues are non-zero
- If A is real and **symmetric**
 - All eigenvalues are **real**.

$$\text{If } |A - \lambda I| = 0 \text{ and } A = A^T \Rightarrow \lambda \in \mathbb{R}$$

- The eigenvectors for distinct eigenvalues are **orthogonal**.

$$A\mathbf{v}_{\{1,2\}} = \lambda_{\{1,2\}}\mathbf{v}_{\{1,2\}} \text{ and } \lambda_1 \neq \lambda_2 \Rightarrow \mathbf{v}_1 \cdot \mathbf{v}_2 = 0$$

- If A is positive definite, then all eigenvalues are **positive**.

$$\forall \mathbf{w} \in \mathbb{R}^n, \mathbf{w}^T A \mathbf{w} > 0, \text{ then if } A\mathbf{v} = \lambda \mathbf{v} \Rightarrow \lambda > 0$$

Singular Value Decomposition (SVD)

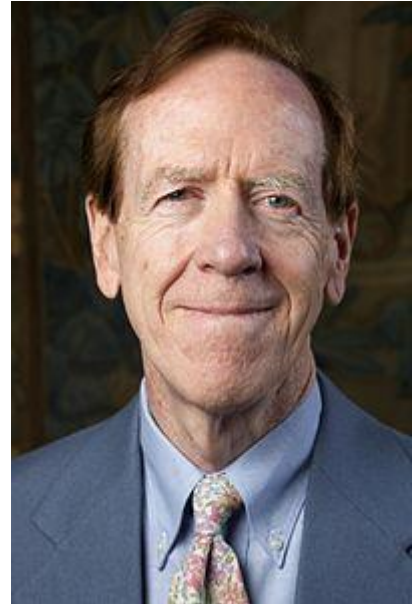
- In practice, we compute the PCs via singular value decomposition (**SVD**) on the **centered** data matrix.

Singular Value Decomposition (SVD)

■ “The highpoint of linear algebra”--- *Gillert Strang*

Awards and honors [\[edit \]](#)

- Rhodes Scholar (1955)
- Alfred P. Sloan Fellow (1966–1967)
- Chauvenet Prize, Mathematical Association of America (1976)
- Honorary Professor, Xian Jiaotong University, China (1980)
- American Academy of Arts and Sciences (1985)
- Honorary Fellow, Balliol College, Oxford University (1999)
- Honorary Member, Irish Mathematical Society (2002)
- Award for Distinguished Service to the Profession, Society for Industrial and Applied Mathematics (2003)
- Lester R. Ford Award (2005)^[3]
- Von Neumann Medal, US Association for Computational Mechanics (2005)
- Haimo Prize, Mathematical Association of America (2007)^[4]
- Su Buchin Prize, International Congress (ICIAM, 2007)
- Henrici Prize (2007)
- National Academy of Sciences (2009)
- Doctor Honoris Causa, University of Toulouse (2010)
- Fellow of the American Mathematical Society (2012)^[5]
- Doctor Honoris Causa, Aalborg University (2013)
- Fellow of the Society for Industrial and Applied Mathematics (2009) ^[6]



1934~

SVD

rank: the number of linearly independent rows (or columns)
 $0 \leq \text{rank}(\mathbf{A}) \leq \min(m, n)$

■ Any $m \times n$ matrix \mathbf{A} of rank r can be decomposed into: $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$

➤ For $m > n$

$$\mathbf{A} = \begin{matrix} & \mathbf{U}_{m \times m} & & \mathbf{\Sigma}_{m \times n} & & \mathbf{V}_{n \times n} \\ & & & & & \\ \begin{bmatrix} \vdots & & \vdots & \vdots & & \vdots \\ \mathbf{u}_1 & \dots & \mathbf{u}_r & \mathbf{u}_{r+1} & \dots & \mathbf{u}_m \\ \vdots & & \vdots & \vdots & & \vdots \end{bmatrix} & & \begin{bmatrix} \sigma_1 & & & & & \\ & \ddots & & & & \\ & & \sigma_r & & & \\ & & & 0 & & \\ & & & & \ddots & \\ 0 & \dots & \dots & \dots & \dots & 0 \\ \vdots & & & & & \vdots \\ 0 & \dots & \dots & \dots & \dots & 0 \end{bmatrix} & & \begin{bmatrix} \dots & \mathbf{v}_1^\top & \dots \\ \vdots & & \\ \dots & \mathbf{v}_r^\top & \dots \\ \dots & \mathbf{v}_{r+1}^\top & \dots \\ \vdots & & \\ \dots & \mathbf{v}_n^\top & \dots \end{bmatrix} \end{matrix}$$

■ Orthogonal matrix \mathbf{U} : $\mathbf{U}^T \mathbf{U} = \mathbf{I}_{m \times m}$, and $\mathbf{U} \mathbf{U}^T = \mathbf{I}_{m \times m}$.

■ The columns of \mathbf{U} are defined as *left singular vectors*.

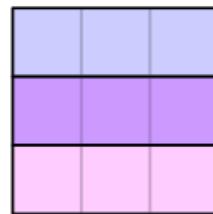
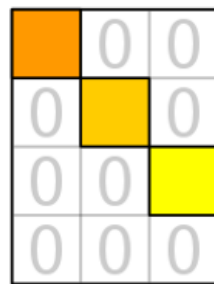
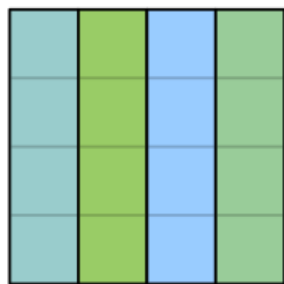
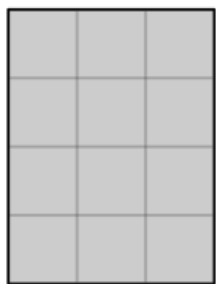
■ Orthogonal matrix \mathbf{V} : $\mathbf{V}^T \mathbf{V} = \mathbf{I}_{n \times n}$, and $\mathbf{V} \mathbf{V}^T = \mathbf{I}_{n \times n}$.

■ The columns of \mathbf{V} are defined as *right singular vectors*.

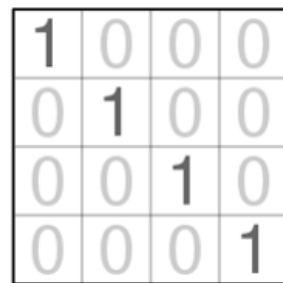
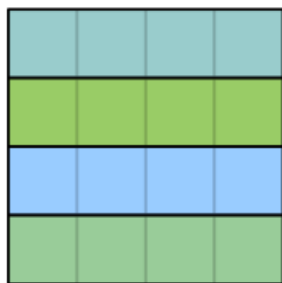
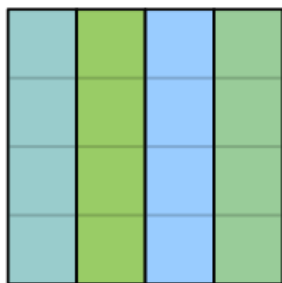
■ Rectangular diagonal matrix $\mathbf{\Sigma}$: we arrange the *singular values* as $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0$.

An orthogonal matrix is a square matrix whose columns and rows are orthogonal unit vectors (orthonormal vectors).

SVD



$$\begin{matrix} \mathbf{A} \\ m \times n \end{matrix} = \begin{matrix} \mathbf{U} \\ m \times m \end{matrix} \begin{matrix} \mathbf{\Sigma} \\ m \times n \end{matrix} \begin{matrix} \mathbf{V}^* \\ n \times n \end{matrix}$$

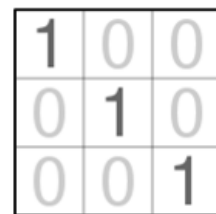
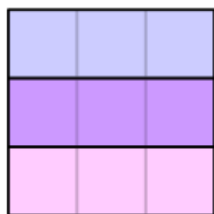
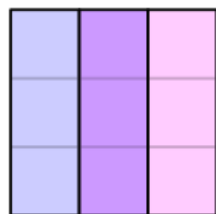


\mathbf{U}

\mathbf{U}^*

$=$

\mathbf{I}_m



\mathbf{V}

\mathbf{V}^*

$=$

\mathbf{I}_n

Singular Value Decomposition (SVD)

- U and V are orthogonal matrices. The columns of U and V are basis vectors in \mathbb{R}^m and \mathbb{R}^n , respectively.

The linear transformation A can be interpreted as a composition of three geometrical transformations: a **rotation or reflection** (V^T), followed by a coordinate-by-coordinate **scaling** (Σ), followed by another **rotation or reflection** (U).

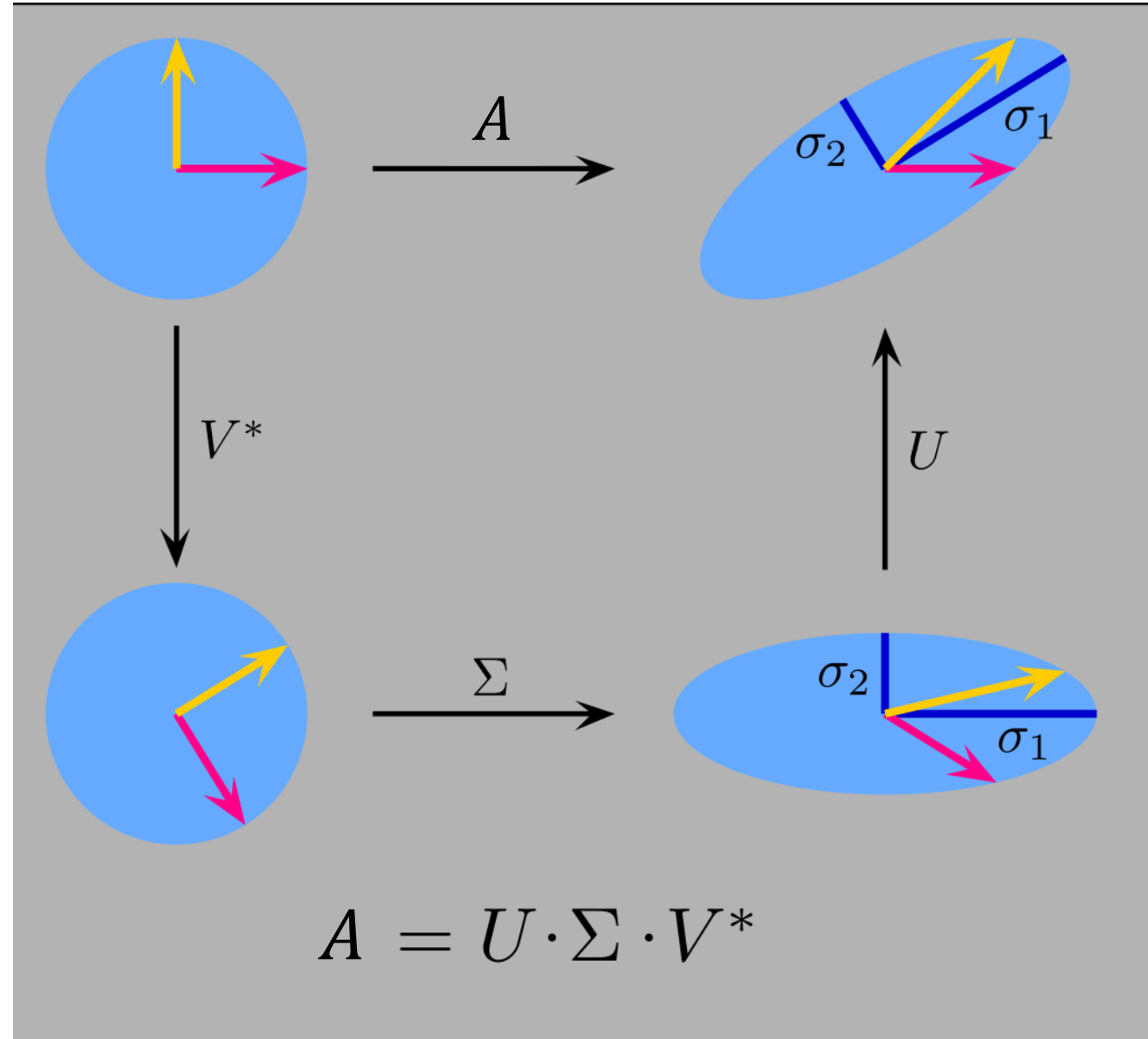
Singular Value Decomposition (SVD)

Illustration of the SVD of a real 2×2 matrix $A = U\Sigma V^T$.

- **Top:** The action of A .
- **Left:** The action of V^T , a rotation.
- **Bottom:** The action of Σ , a scaling by the singular values σ_1 horizontally and σ_2 vertically.
- **Right:** The action of U , another rotation.

$$\begin{bmatrix} \textcircled{3} \\ -2 \end{bmatrix} \begin{bmatrix} \textcircled{2} \\ 1 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix}$$

Where \hat{i} lands Where \hat{j} lands
变换后的 i 变换后的 j



Singular Value Decomposition (SVD)

Given the matrix $A = \begin{bmatrix} 3 & 1 \\ 2 & 1 \end{bmatrix}$

$$U = \begin{bmatrix} 0.8174 & -0.5760 \\ 0.5760 & 0.8174 \end{bmatrix} \quad \Sigma = \begin{bmatrix} 3.8643 & 0 \\ 0 & 0.2588 \end{bmatrix} \quad V^T = \begin{bmatrix} 0.9327 & 0.3606 \\ -0.3606 & 0.9327 \end{bmatrix}$$

Let us focus on the transformation of the original basis vectors by A

$$e_1 = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \quad e_2 = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

First, by rotation V^T , we have

$$V^T e_1 = \begin{bmatrix} 0.9327 \\ -0.3606 \end{bmatrix} \quad V^T e_2 = \begin{bmatrix} 0.3606 \\ 0.9327 \end{bmatrix}$$

Singular Value Decomposition (SVD)

Given the matrix $A = \begin{bmatrix} 3 & 1 \\ 2 & 1 \end{bmatrix}$

$$U = \begin{bmatrix} 0.8174 & -0.5760 \\ 0.5760 & 0.8174 \end{bmatrix} \quad \Sigma = \begin{bmatrix} 3.8643 & 0 \\ 0 & 0.2588 \end{bmatrix} \quad V^T = \begin{bmatrix} 0.9327 & 0.3606 \\ -0.3606 & 0.9327 \end{bmatrix}$$

Second, by Σ , we scale $V^T e_1$ and $V^T e_2$ along with the original coordinate system (e_1, e_2)

$$\Sigma V^T e_1 = \begin{bmatrix} 3.6042 \\ -0.0933 \end{bmatrix} \quad \Sigma V^T e_2 = \begin{bmatrix} 1.3935 \\ 0.2414 \end{bmatrix}$$

Last, by rotation U , we have

$$U \Sigma V^T e_1 = \begin{bmatrix} 3 \\ 2 \end{bmatrix} \quad U \Sigma V^T e_2 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

Singular Value Decomposition (SVD)

- Any $m \times n$ matrix A of rank r can be decomposed into: $A = U\Sigma V^T$

► For $m > n$

$$\mathbf{A} = \begin{bmatrix} \vdots & & \vdots & \vdots & & \vdots \\ \mathbf{u}_1 & \dots & \mathbf{u}_r & \mathbf{u}_{r+1} & \dots & \mathbf{u}_m \\ \vdots & & \vdots & \vdots & & \vdots \end{bmatrix} \begin{bmatrix} \sigma_1 & & & & & \\ & \ddots & & & & \\ & & \sigma_r & & & \\ & & & 0 & & \\ & & & & \ddots & \\ 0 & \dots & \dots & \dots & \dots & 0 \\ \vdots & & & & & \vdots \\ 0 & \dots & \dots & \dots & \dots & 0 \end{bmatrix} \begin{bmatrix} \dots & \mathbf{v}_1^\top & \dots \\ & \vdots & \\ \dots & \mathbf{v}_r^\top & \dots \\ \dots & \mathbf{v}_{r+1}^\top & \dots \\ & \vdots & \\ \dots & \mathbf{v}_n^\top & \dots \end{bmatrix}$$

■ Special Properties:

- The columns of \mathbf{U} (i.e., left singular vectors) are **eigenvectors** of $\mathbf{A}\mathbf{A}^T$.
- The columns of \mathbf{V} (i.e., right singular vectors) are **eigenvectors** of $\mathbf{A}^T\mathbf{A}$.
- Eigenvalues $\lambda_1, \dots, \lambda_r$ of $\mathbf{A}\mathbf{A}^T$ are the eigenvalues of $\mathbf{A}^T\mathbf{A}$.
- Singular value $\sigma_i = \sqrt{\lambda_i}$.

Singular Value Decomposition (SVD)

- Prove that the columns of V (i.e., right singular vectors) are eigenvectors of $A^T A$.

$A^T A = (U\Sigma V^T)^T U\Sigma V^T = V\Sigma^T \Sigma V^T$ This is the **eigen-decomposition** of $A^T A$.

V is the eigenvector matrix of $A^T A$, and $\Sigma^T \Sigma$ is the eigenvalue matrix of $A^T A$, i.e., singular values are positive square roots of eigenvalues.

Theorem1: Every **symmetric** matrix M is **orthogonally** diagonalizable, i.e., there exists an **orthogonal** matrix Q (i.e., $Q^T = Q^{-1}$) such that $Q^T M Q = D$ (i.e., $M = Q D Q^T$) is a diagonal matrix. (https://en.wikipedia.org/wiki/Diagonalizable_matrix)

Compact SVD

- Only the $r = \text{rank}(A)$ column vectors of U and r row vectors of V^T corresponding to the non-zero singular values Σ_r are calculated.

$$A = \begin{bmatrix} \vdots & \vdots \\ \mathbf{u}_1 & \dots & \mathbf{u}_r & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \vdots \end{bmatrix} \begin{bmatrix} \sigma_1 & & & & \\ & \ddots & & & \\ & & \sigma_r & & \\ & & & 0 & \\ & & & & \ddots \\ 0 & \dots & \dots & \dots & 0 \\ \vdots & & & & \vdots \\ 0 & \dots & \dots & \dots & 0 \end{bmatrix} \begin{bmatrix} \dots & \mathbf{v}_1^\top & \dots \\ \vdots & & \\ \dots & \mathbf{v}_r^\top & \dots \\ \vdots & & \\ \dots & \mathbf{v}_{r+1}^\top & \dots \\ \vdots & & \\ \dots & \mathbf{v}_n^\top & \dots \end{bmatrix} \quad m > n$$

- Economy** version $A = \underbrace{U_r}_{m \times r} \underbrace{\Sigma_r}_{r \times r} \underbrace{V_r^T}_{r \times n} \quad \Sigma_r = \text{diag}(\sigma_1, \dots, \sigma_r)$

Any information loss?

Truncated SVD

- Only k column vectors of \mathbf{U} and k row vectors of \mathbf{V}^T corresponding to the non-zero singular values Σ_k are calculated, $0 < k < r, r = \text{rank}(\mathbf{A})$.

$$\mathbf{A} = \begin{bmatrix} \vdots & & & & \\ \mathbf{u}_1 & \dots & \mathbf{u}_r & \mathbf{u}_{r+1} & \dots & \mathbf{u}_m \\ \vdots & & & & & \end{bmatrix} \begin{bmatrix} \boxed{\begin{matrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_r \end{matrix}} & & & & \\ & & 0 & & \\ & & & \ddots & \\ & & & & 0 \\ 0 & \dots & \dots & \dots & \dots & 0 \\ \vdots & & & & & \vdots \\ 0 & \dots & \dots & \dots & \dots & 0 \end{bmatrix} \begin{bmatrix} \dots & \mathbf{v}_1^\top & \dots \\ \vdots & & \\ \dots & \mathbf{v}_r^\top & \dots \\ \dots & \mathbf{v}_{r+1}^\top & \dots \\ \vdots & & \\ \dots & \mathbf{v}_n^\top & \dots \end{bmatrix} \quad m > n$$

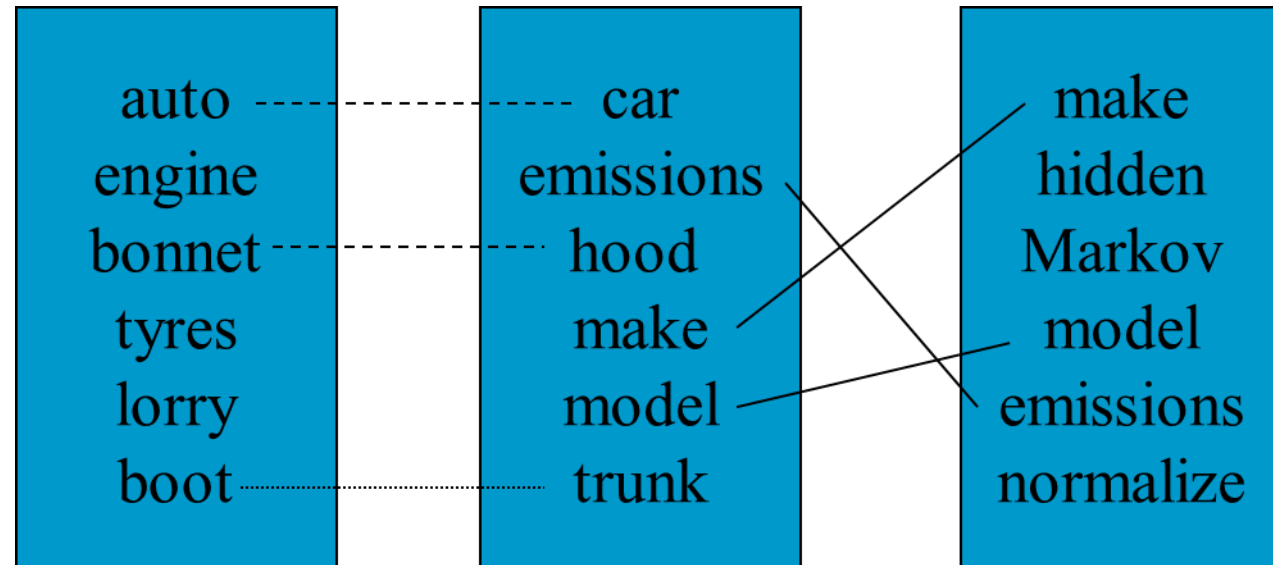
- More Economical** $\mathbf{A} = \underbrace{\mathbf{U}_k}_{m \times k} \underbrace{\Sigma_k}_{k \times k} \underbrace{\mathbf{V}_k^T}_{k \times n} \quad \Sigma_k = \text{diag}(\sigma_1, \dots, \sigma_k)$

Truncated SVD is no longer an exact decomposition of the original matrix.

SVD Application1-Latent Semantic Indexing (LSI)

- LSI was proposed to address two problems with the vector space model
 - **synonymy**: many ways to refer to the same object, e.g. *car* and *automobile*
 - leads to poor recall (*recall*: portion of the target items that the system selected)
 - **polysemy**: most words have more than one distinct meaning, e.g. *model*, *python*, *chip*
 - leads to poor precision (*precision*: portion of selected items that the system got right)

Index Words	Titles								
	T1	T2	T3	T4	T5	T6	T7	T8	T9
book			1	1					
dads						1			1
dummies		1						1	
estate							1		1
guide	1					1			
investing	1	1	1	1	1	1	1	1	1
market	1		1						
real							1		1
rich						2			1
stock	1		1					1	
value				1	1				



Synonymy

Will have small cosine

but are related

Polysemy

Will have large cosine

but not truly related

SVD Application1-Latent Semantic Indexing (LSI)

- LSI is a technique that projects queries and documents into a space with “latent” semantic dimensions.
- In the latent semantic space, a query and a document can have high cosine similarity even if they don’t share any terms, as long as their terms are semantically similar in a sense.

Index Words	Titles								
	T1	T2	T3	T4	T5	T6	T7	T8	T9
book			1	1					
dads						1			1
dummies		1						1	
estate							1		1
guide	1					1			
investing	1	1	1	1	1	1	1	1	1
market	1		1						
real							1		1
rich						2			1
stock	1		1					1	
value				1	1				

A



book	0.15	-0.27	0.04
dads	0.24	0.38	-0.09
dummies	0.13	-0.17	0.07
estate	0.18	0.19	0.45
guide	0.22	0.09	-0.46
investing	0.74	-0.21	0.21
market	0.18	-0.30	-0.28
real	0.18	0.19	0.45
rich	0.36	0.59	-0.34
stock	0.25	-0.42	-0.28
value	0.12	-0.14	0.23

U_k

$k = 3$

3.91	0	0
0	2.61	0
0	0	2.00

Σ_k

T1	T2	T3	T4	T5	T6	T7	T8	T9
0.35	0.22	0.34	0.26	0.22	0.49	0.28	0.29	0.44
-0.32	-0.15	-0.46	-0.24	-0.14	0.55	0.07	-0.31	0.44
-0.41	0.14	-0.16	0.25	0.22	-0.51	0.55	0.00	0.34

V_k^T

SVD Application2-Pseudoinverse

- SVD can be used for computing the pseudoinverse of a matrix.

$$A = U\Sigma V^T$$

$$A = \begin{bmatrix} \vdots & & \vdots & \vdots & & \vdots \\ \mathbf{u}_1 & \dots & \mathbf{u}_r & \mathbf{u}_{r+1} & \dots & \mathbf{u}_m \\ \vdots & & \vdots & \vdots & & \vdots \end{bmatrix} \begin{bmatrix} \sigma_1 & & & & & \\ & \ddots & & & & \\ & & \sigma_r & & & \\ & & & 0 & & \\ & & & & \ddots & \\ 0 & \dots & \dots & \dots & \dots & 0 \\ \vdots & & & & & \vdots \\ 0 & \dots & \dots & \dots & \dots & 0 \end{bmatrix} \begin{bmatrix} \dots & \mathbf{v}_1^\top & \dots \\ \vdots & & \\ \dots & \mathbf{v}_r^\top & \dots \\ \dots & \mathbf{v}_{r+1}^\top & \dots \\ \vdots & & \\ \dots & \mathbf{v}_n^\top & \dots \end{bmatrix}$$

$$A^\dagger = V\Sigma^\dagger U^T$$

Σ^\dagger : pseudoinverse of Σ by replacing every non-zero diagonal entry by its reciprocal and transposing the matrix.

$$\begin{bmatrix} \frac{1}{\sigma_1} & & & 0 \dots 0 \\ & \ddots & & \vdots \\ & & \frac{1}{\sigma_r} & \vdots \\ & & & 0 \\ & & & & \ddots & \vdots \\ & & & & & 0 \dots 0 \end{bmatrix}$$

$$AA^\dagger A = U\Sigma V^T V\Sigma^\dagger U^T U\Sigma V^T = U\Sigma\Sigma^\dagger\Sigma V^T = U\Sigma V^T = A$$

SVD Application3-PCA

- In practice, we compute the PCs via singular value decomposition (**SVD**) on the **centered** data matrix.

- Form the **centered** data matrix:

$$\mathbf{X} = [(\mathbf{x}_1 - \bar{\mathbf{x}}); \dots; (\mathbf{x}_m - \bar{\mathbf{x}})] \in \mathbb{R}^{d \times m}$$

- Compute its SVD:

$$\mathbf{X} = \mathbf{U}_{d \times d} \mathbf{D}_{d \times m} (\mathbf{V}_{m \times m})^T$$

where \mathbf{U} and \mathbf{V} are **orthogonal** matrices, \mathbf{D} is a diagonal matrix.

SVD Application3-PCA

- Note that the scatter/covariance matrix can be written as

$$\mathbf{S} = \mathbf{X}\mathbf{X}^T = \mathbf{U}\mathbf{D}^2\mathbf{U}^T \quad \mathbf{X} = \mathbf{U}_{d \times d} \mathbf{D}_{d \times m} (\mathbf{V}_{m \times m})^T$$

- The eigenvectors of \mathbf{S} are the columns of \mathbf{U} and the eigenvalues are the diagonal elements of \mathbf{D}^2 .
- Take **only a few significant** eigenvalue-eigenvector pairs $p \ll d$. The new reconstructed sample from low-dim space is:

$$\hat{\mathbf{x}}_i = \bar{\mathbf{x}} + \mathbf{U}_{d \times p} (\mathbf{U}_{d \times p})^T (\mathbf{x}_i - \bar{\mathbf{x}})$$

SVD Application4-Reconstruction

Original



Reconstruction using 50 PCs



$$\hat{\mathbf{x}}_i = \bar{\mathbf{x}} + \mathbf{U}_{d \times p} (\mathbf{U}_{d \times p})^T (\mathbf{x}_i - \bar{\mathbf{x}})$$

Advantages of Using SVD for PCA

- No need to compute the covariance matrix $S = \mathbf{X}\mathbf{X}^T$
- Numerically more accurate, since the formation of $\mathbf{X}\mathbf{X}^T$ can cause loss of precision.
 - For example, the Läuchli matrix:

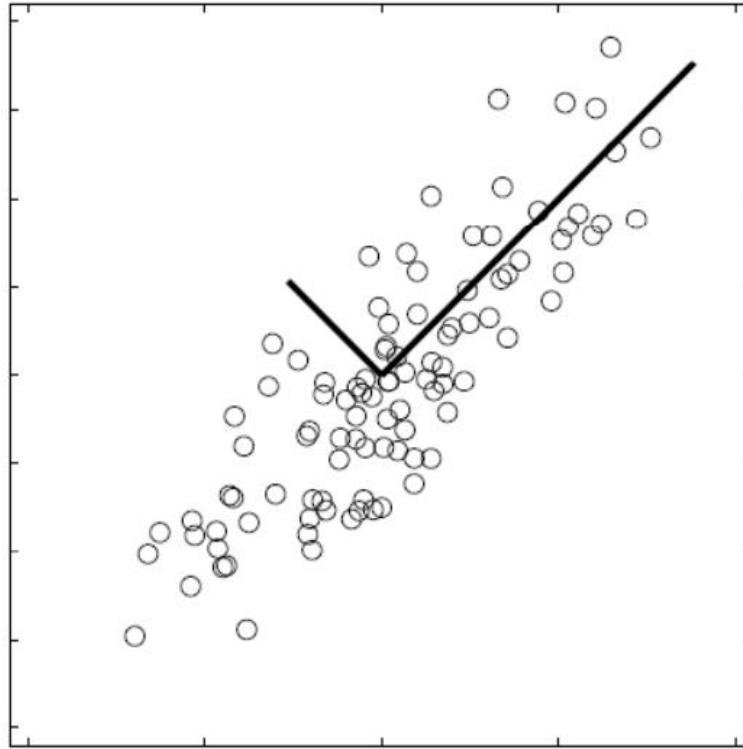
$$\begin{pmatrix} 1 & 1 & 1 \\ \epsilon & 0 & 0 \\ 0 & \epsilon & 0 \\ 0 & 0 & \epsilon \end{pmatrix}^T$$

where ϵ is a tiny number.

Visualize PCs

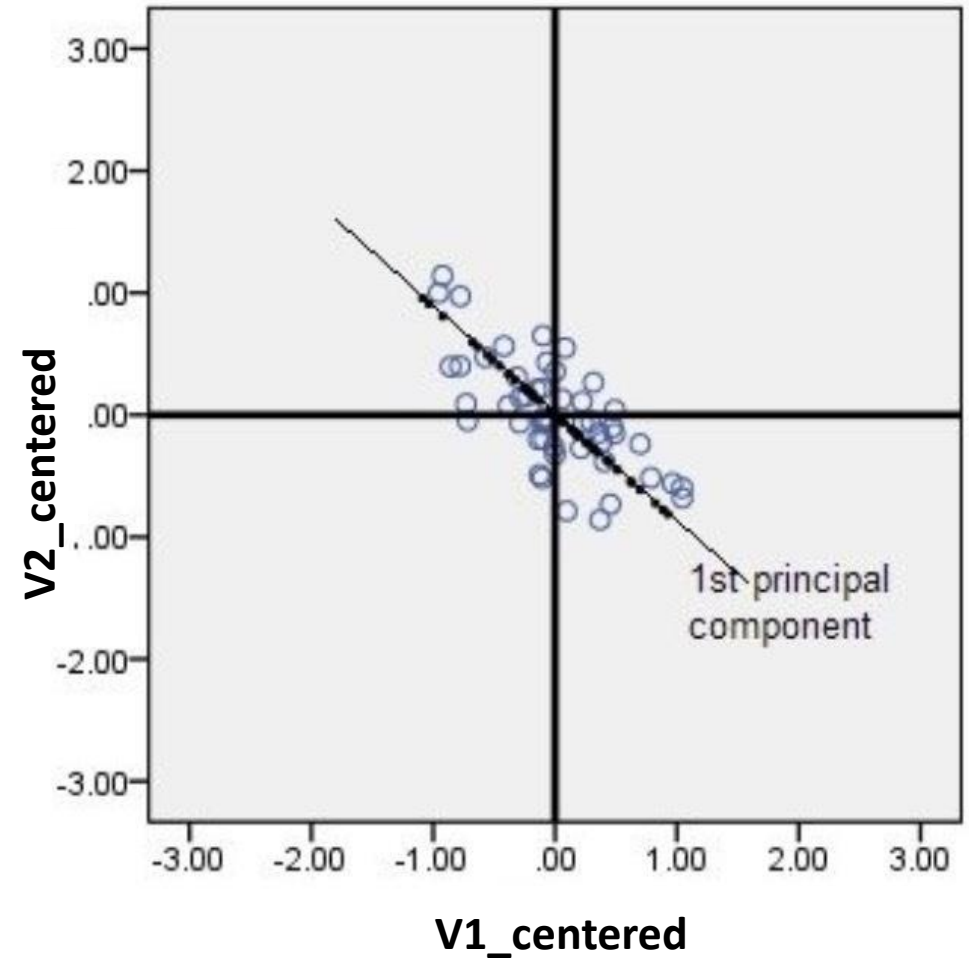
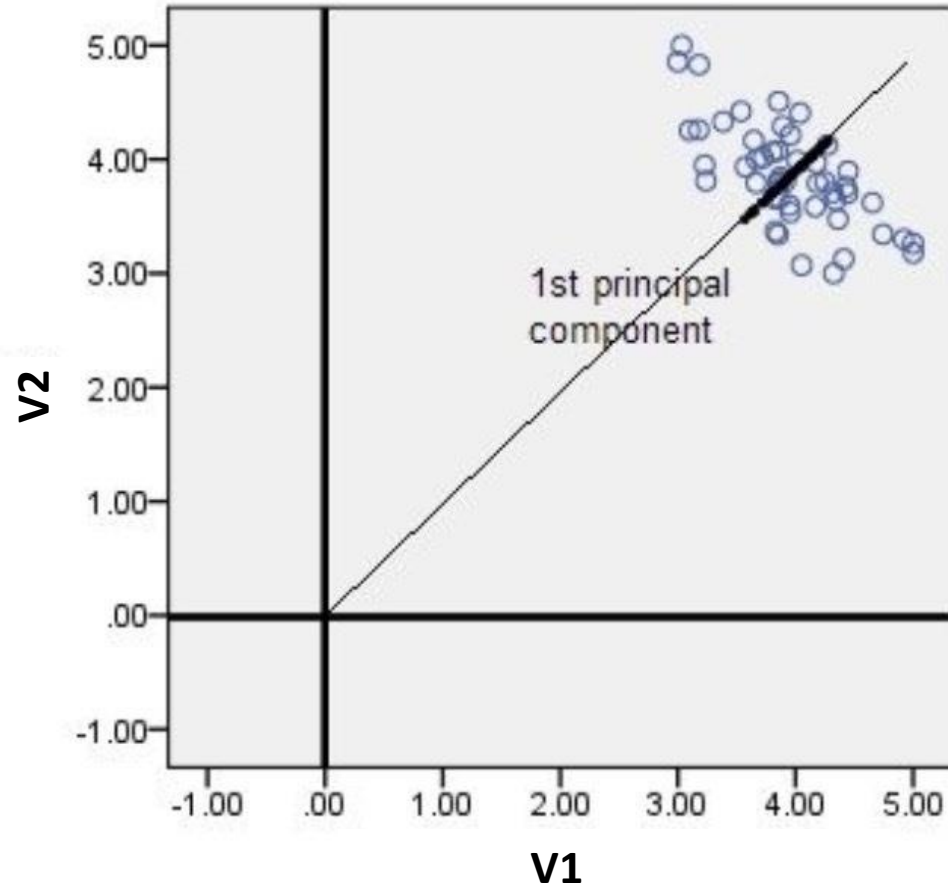
The columns of \mathbf{U} are eigenvectors of $\mathbf{S} = \mathbf{A}\mathbf{A}^T$.

$$\mathbf{y} = \mathbf{U}^T \mathbf{x}$$



Data points are represented in a **rotated** orthogonal coordinate system: the origin is the mean of the data points and the axes are provided by the eigenvectors.

The Necessity of Centralization



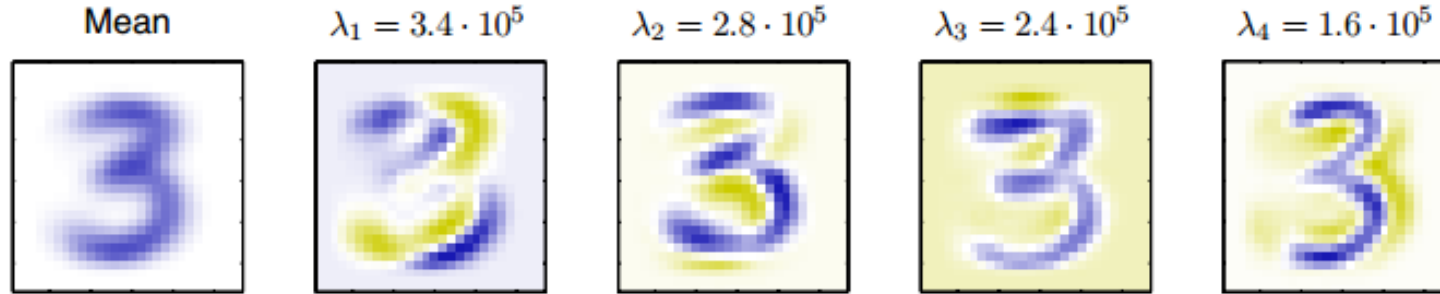
How Many PCs to Keep?

To choose p based on percentage of energy to retain, we can use the following criterion (smallest p):

$$\frac{\sum_{i=1}^p \lambda_i}{\sum_{i=1}^d \lambda_i} \geq \textit{Threshold} \quad (e.g., 0.95)$$

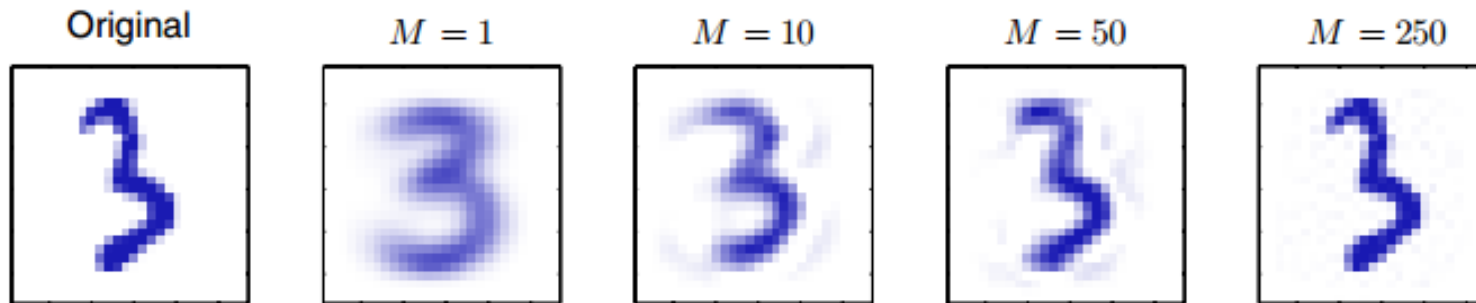
PCA-Applications

Data Compression



We represent the eigenvectors as images of the same size as the data points.

The mean vector \bar{x} along with the first four PCA eigenvectors u_1, \dots, u_4 for the off-line digits data set, together with the corresponding eigenvalues.



An original example from the off-line digits data set together with its PCA reconstructions obtained by retaining M principal components for various values of M . As M increases the reconstruction becomes more accurate and would become perfect when $M = D = 28 \times 28 = 784$.

PCA-Applications

Data Preprocessing

- The goal is **not** dimensionality reduction but rather the transformation of a data set in order to **standardizing** the data.
- Important in allowing subsequent pattern recognition algorithms to be applied successfully to the data set.
- Typically, it is done when the original variables are measured in different order of magnitudes or have significantly different variability.

PCA-Applications

Data Preprocessing

- The goal is **not** dimensionality reduction but rather the transformation of a data set in order to **standardizing** the data.

Traditionally, we can make a linear re-scaling of the individual variables such that each variable had zero **mean** and unit **variance**.

$$\frac{x_{ni} - \bar{x}_i}{\sigma_i}$$

However, using PCA we can make a more substantial normalization of the data to give it **zero mean** and **unit covariance**, so that variables become **decorrelated**.

PCA-Applications

Data Preprocessing

- The goal is **not** dimensionality reduction but rather the transformation of a data set in order to **standardizing** the data.

We first write the eigenvalue decomposition $SU = U\Lambda$

Then for each data point, we define $\mathbf{y}_n = \Lambda^{-\frac{1}{2}}U^T(\mathbf{x}_n - \bar{\mathbf{x}})$

Clearly, $\{\mathbf{y}_n\}$ have zero mean, and we thus have the covariance matrix

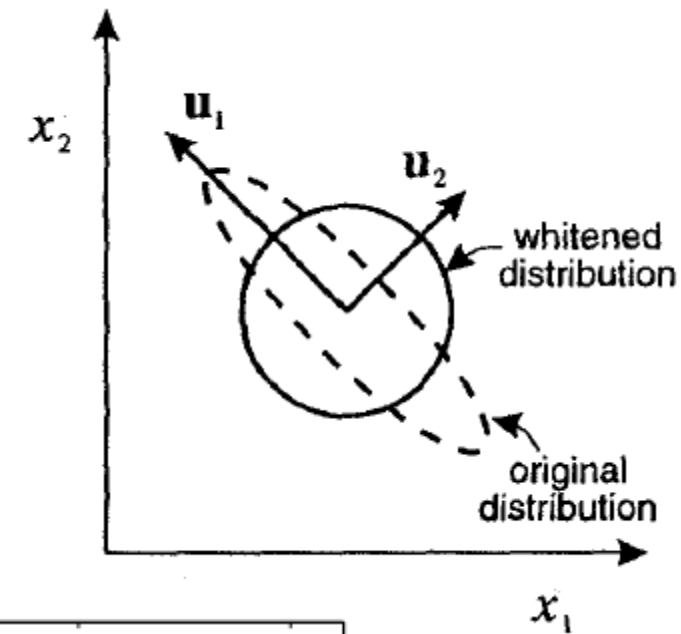
$$\frac{1}{m} \sum_{n=1}^m \mathbf{y}_n \mathbf{y}_n^T = \frac{1}{m} \sum_{n=1}^m \Lambda^{-\frac{1}{2}} U^T (\mathbf{x}_n - \bar{\mathbf{x}}) (\mathbf{x}_n - \bar{\mathbf{x}})^T U \Lambda^{-\frac{1}{2}} = \Lambda^{-\frac{1}{2}} U^T S U \Lambda^{-\frac{1}{2}} = \Lambda^{-\frac{1}{2}} \Lambda \Lambda^{-\frac{1}{2}} = I$$

This operation is known as *whitening* or *sphereing* the data

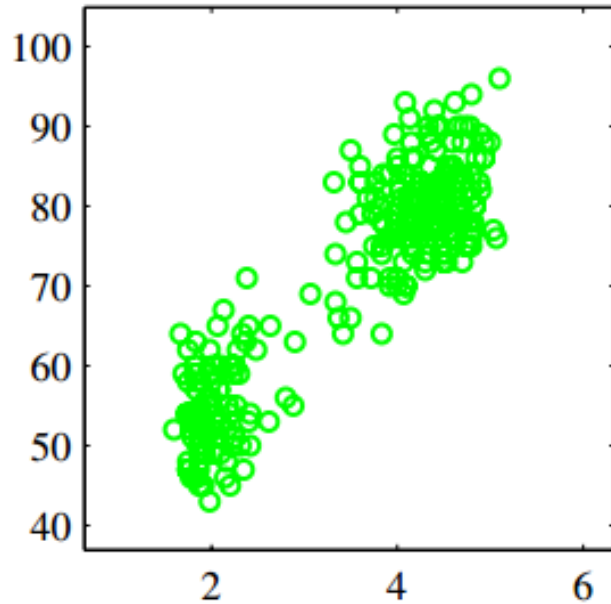
PCA-Applications

Data Preprocessing

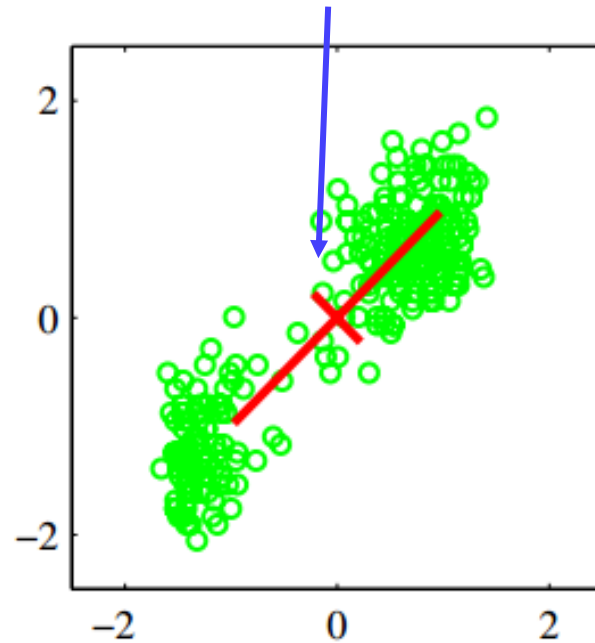
- The goal is **not** dimensionality reduction but rather the transform order to **standardizing** the data.



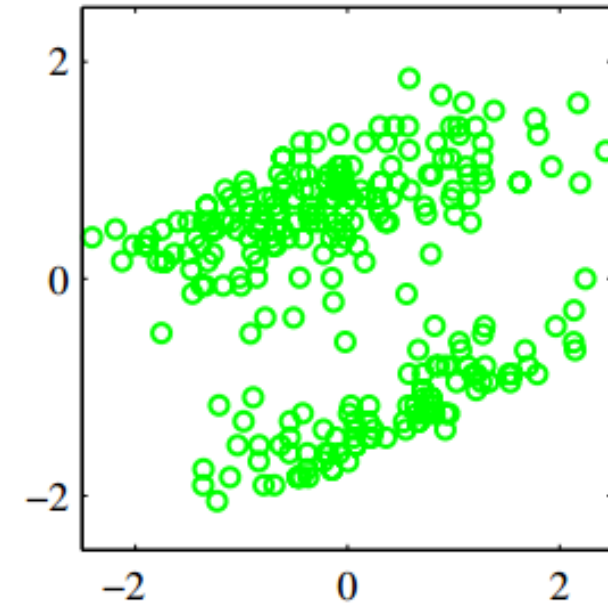
Principal axes of this normalized data set.



Original data



$$\frac{x_{ni} - \bar{x}_i}{\sigma_i}$$

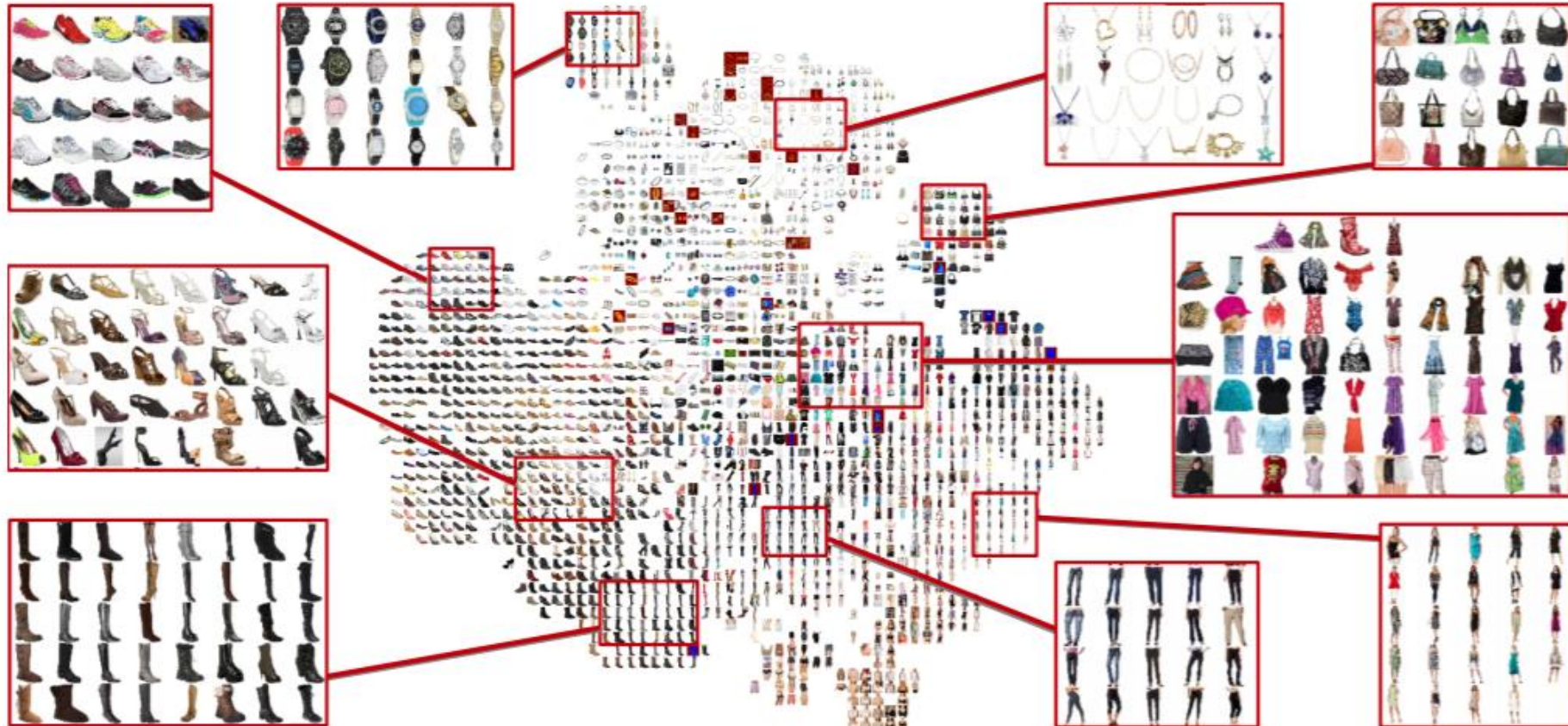


$$y_n = \Lambda^{-\frac{1}{2}} U^T (x_n - \bar{x})$$

PCA-Applications

Data Visualization

- Each data point is projected onto a two-dimensional principal subspace.



PCA and Classification

- Classification with PCA
 - Project both training and testing data into the PCs space
 - For each testing sample, use Nearest Neighbor for classification
 - **Issue:** accuracy is sensitive to the number of PCs
- PCA may not be always an optimal feature extraction technique for classification.
 - Suppose there are C classes in the training data
 - PCA is based on the sample covariance which characterizes the scatter of the **entire** data set, **irrespective of class-membership.**
 - The projection axes chosen by PCA might not provide good discrimination power.

Summary

- Two commonly used definitions of PCA
 - Maximum variance formulation
 - Minimum-error formulation
- Covariance Matrix
 - Symmetric, Positive semi-definite
- Lagrange Multiplier
 - Constrained system \rightarrow Unconstrained system
 - Critical points
- Data Point Reconstruction
- SVD
 - U, V orthogonal matrix
- PCA Applications