# Lecture 13: Principal components analysis (PCA)

Lecturer: Jie Fu

# High-Dimensional Data

- High-Dimensions = Lot of Features

*Surveys Netflix*

| | movie 1 | movie 2 | movie 3 | movie 4 | movie 5 | movie 6 |
|---|---|---|---|---|---|---|
| Tom | 5 | ? | ? | 1 | 3 | ? |
| George | ? | ? | 3 | 1 | 2 | 5 |
| Susan | 4 | 3 | 1 | ? | 5 | 1 |
| Beth | 4 | 3 | ? | 2 | 4 | 2 |

*Food preference*

| | kale | taco bell | sashimi | pop tarts |
|---|---|---|---|---|
| Alice | 10 | 1 | 2 | 7 |
| Bob | 7 | 2 | 1 | 10 |
| Carolyn | 2 | 9 | 7 | 3 |
| Dave | 3 | 6 | 10 | 2 |

- PCA: Unsupervised learning techniques to extract hidden dimensional structure from high dimensional dataset

  - Visualization
  - Efficient use of resources.
  - Statistical: lower dimension --> better generalization.
  - Further processing for other machine learning algorithm.

# Motivating problem

- Friends' preferences of four different food choice.

- Dimension of data points: 4

- Number of data points: 4

Can we visualize the data in less than 4 dimension?

|         | kale | taco bell | sashimi | pop tarts |
|---------|------|-----------|---------|-----------|
| Alice   | 10   | 1         | 2       | 7         |
| Bob     | 7    | 2         | 1       | 10        |
| Carolyn | 2    | 9         | 7       | 3         |
| Dave    | 3    | 6         | 10      | 2         |

Table 1: Your friends' ratings of four different foods.

- Each row of the data can be expressed approximately:

| | kale | taco bell | sashimi | pop tarts |
|---|---|---|---|---|
| Alice | 10 | 1 | 2 | 7 |
| Bob | 7 | 2 | 1 | 10 |
| Carolyn | 2 | 9 | 7 | 3 |
| Dave | 3 | 6 | 10 | 2 |

Table 1: Your friends' ratings of four different foods.

| Name | $(a_1, a_2)$ |
|---|---|
| Alice | $(1, \ 1)$ |
| Bob | $(1, \ -1)$ |
| Carolyn | $(-1, \ -1)$ |
| Dave | $(-1, \ 1)$ |

Table 1: Values of $(a_1, a_2)$ for each person

$$\bar{\mathbf{x}} + a_1 \mathbf{v}_1 + a_2 \mathbf{v}_2,$$

where

$$\bar{\mathbf{x}} = (5.5, 4.5, 5, 5.5)$$

is the average of the data points,

$$\mathbf{v}_1 = (3, -3, -3, 3),$$

$$\mathbf{v}_2 = (1, -1, 1, -1),$$

- Reduce the dimensionality of data points (eg. 4 to 2):
- Given a list of m n-dimensional vectors (data points),
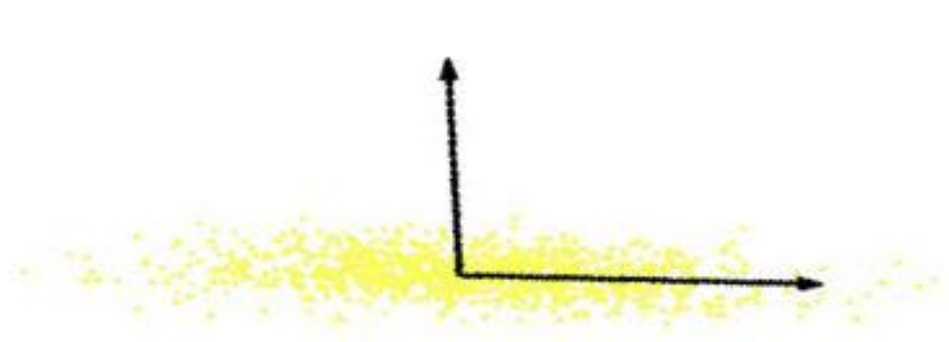
$$x_1, x_2, \ldots, x_m \in R^n$$

For each vector $x_i$, express it as linear combinations of k n-dimensional vectors $v_1, \ldots, v_n \in R^n$ such that
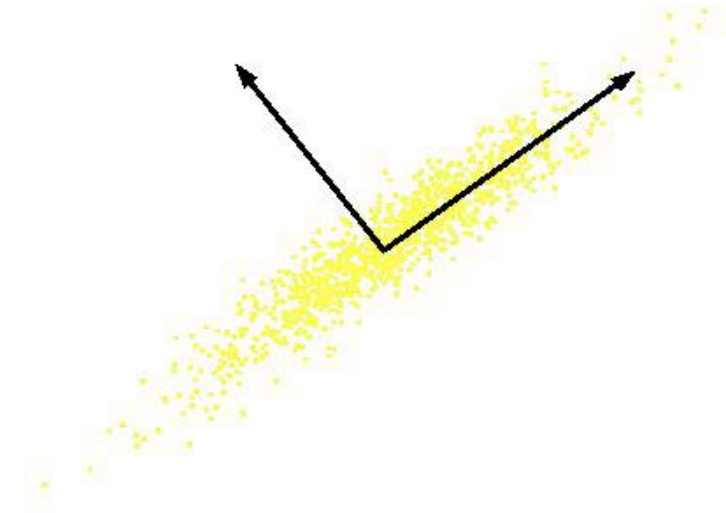
$$x_i \approx \sum_{j \in 1}^{k} a_{ij} v_j$$

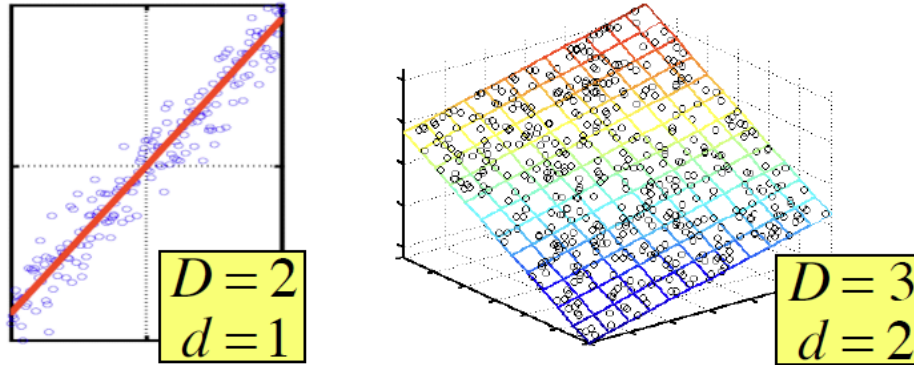Dimension reduction: n → k, which is smaller than n.

# PCA

- PCA is an orthogonal projection or transformation of the data into a possible lower dimensional subspace so that **the variance of the projected data is maximized.**



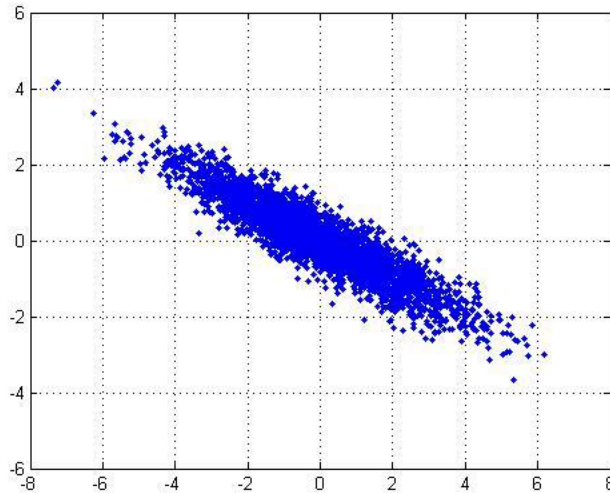*Only one relevant feature*



*Both features are relevant, but*

# PCA



$D = 2$
$d = 1$

$D = 3$
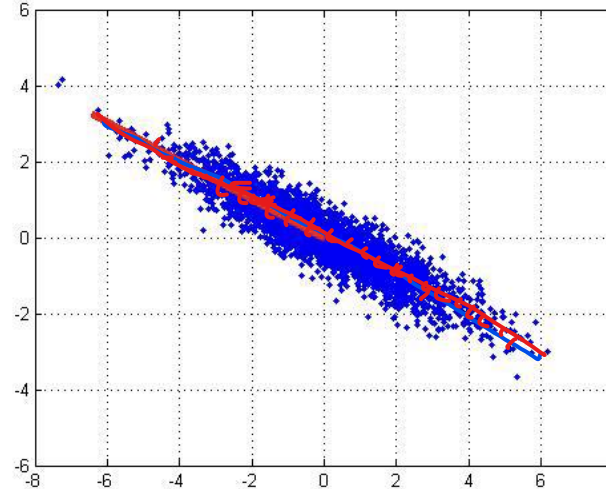$d = 2$

**Does the data mostly lie in a subspace?**
**If so, what is its dimensionality?**

- The goal is to identify the axes or subspace the high-dimensional data should be projected into.
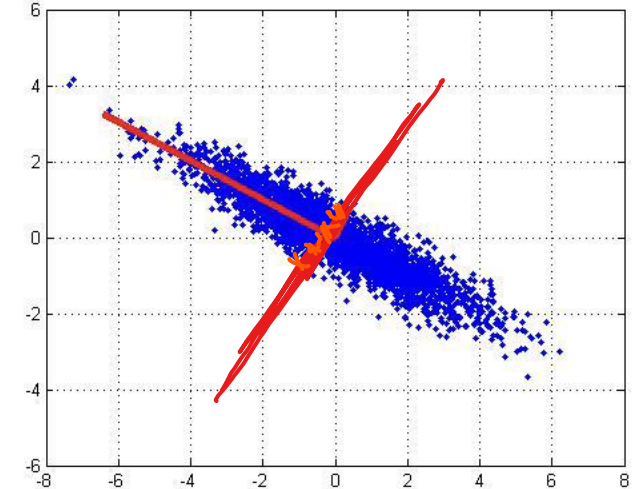
# Maximize the variance

- Why maximize the variance of the projected data?



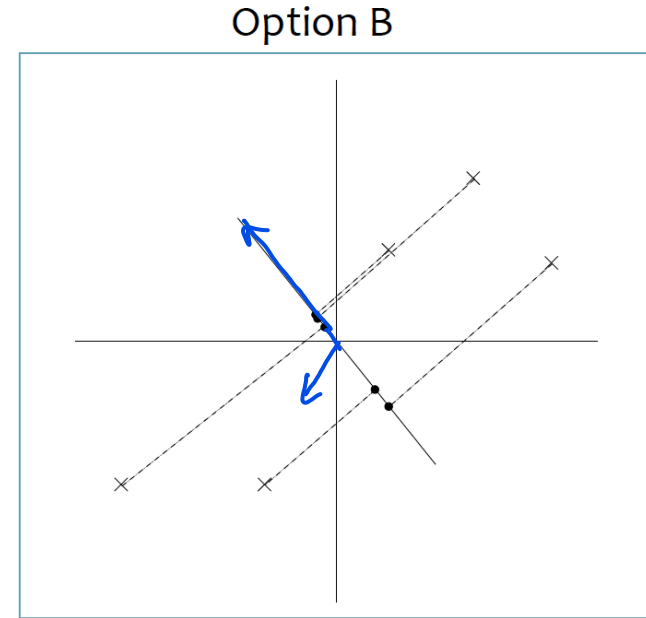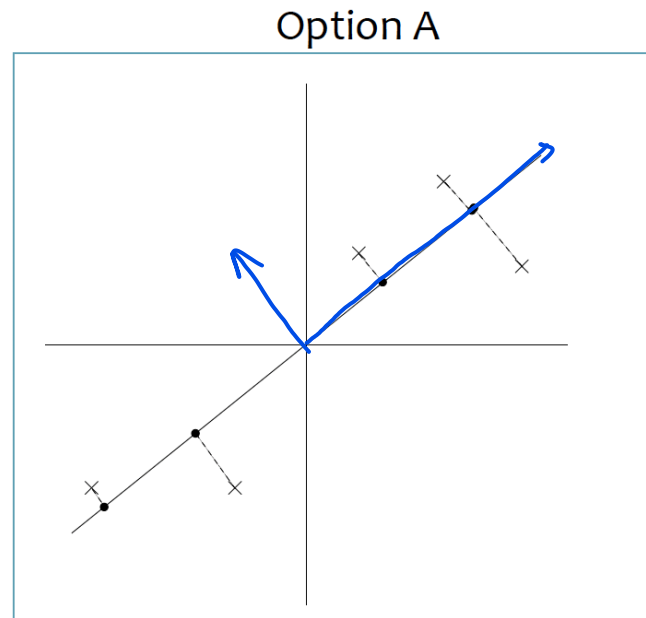*2D Gaussian dataset*  *1st principle component*  *2nd principle component*

Variance tells us how much information or "spread" a dataset has. In PCA, we assume directions with higher variance are more informative.

# Maximize the variance

- Which of the two projections maximize the variance?


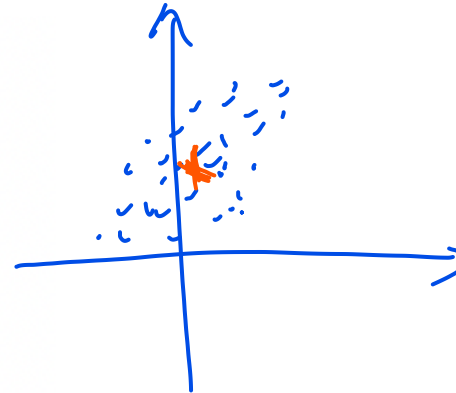
*Figures from Andrew Ng*
*(CS229 Lecture Notes)*

# Maximize the variance

We want to find new axes (directions) to project our data such that:

• The projected data has **maximum variance**.

• The new features (called principal components) are **uncorrelated**.

|        | kale | taco bell | sashimi | pop tarts |
|--------|------|-----------|---------|-----------|
| Alice  | 10   | 1         | 2       | 7         |
| Bob    | 7    | 2         | 1       | 10        |
| Carolyn| 2    | 9         | 7       | 3         |
| Dave   | 3    | 6         | 10      | 2         |

Table 1: Your friends' ratings of four different foods.



Step 1: center the data matrix
Step 2: compute the covariance matrix of the centered data
Step 3: select top k principal components/features

- Center the data
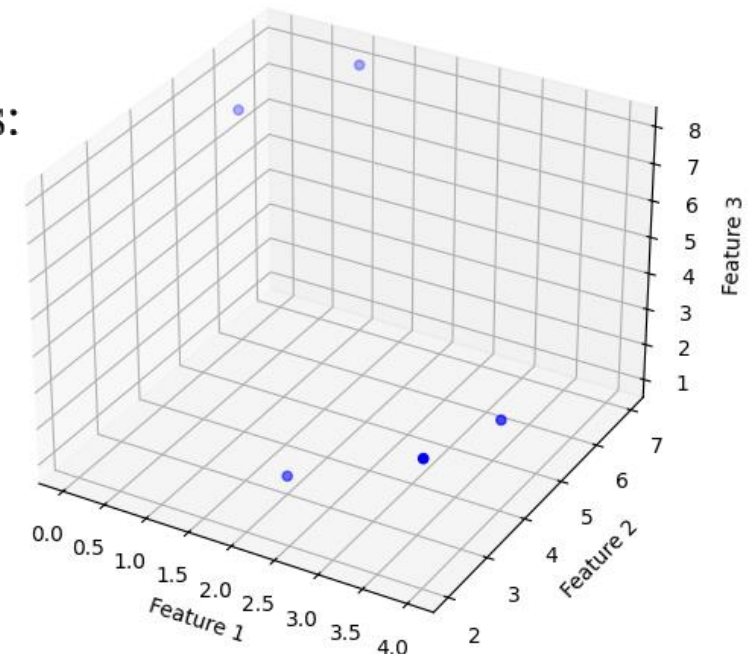
$$X_c = X - \bar{X}$$

- Example:

Consider the following dataset with 5 samples and 3 features:

feature 1  2  3

$$X = \begin{bmatrix} 2 & 3 & 1 \\ 4 & 2 & 4 \\ 4 & 4 & 3 \\ 0 & 6 & 7 \\ 1 & 7 & 8 \end{bmatrix}$$

Mean of each feature

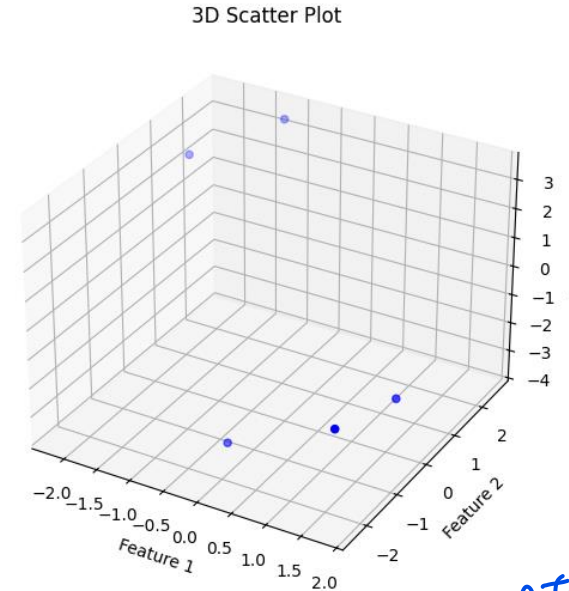mean 1   2   3
2.2   4.4   4.6

3D Scatter Plot

- Centered data matrix

$$X_c = \begin{bmatrix} -0.2 & -1.4 & -3.6 \\ 1.8 & -2.4 & -0.6 \\ 1.8 & -0.4 & -1.6 \\ -2.2 & 1.6 & 2.4 \\ -1.2 & 2.6 & 3.4 \end{bmatrix}$$

3D Scatter Plot



$$K = U \Lambda U^{-1}$$

covariance matrix $\leftarrow$ after rotation

$$\Lambda \begin{bmatrix} \lambda_1 & & 0 \\ & \lambda_2 & \\ 0 & & \ddots \lambda_n \end{bmatrix}$$

Step 2: compute the covariance matrix of the centered data (use the transposed.)

np.cov(M_c.T)

The covariance matrix $K$ is given by:

$$K = \frac{1}{n-1} X_c^{\top} X_c$$

$$\text{Var}(X_i) \qquad \text{Cov}(X_1, X_2)$$

$$\text{Cov}(X_2, X_1) \begin{bmatrix} 3.2 & -2.85 & -3.15 \\ -2.85 & 4.3 & 4.95 \\ -3.15 & 4.95 & 8.3 \end{bmatrix}$$

UF
UNIVERSITY of
FLORIDA

Let $A$ be a $n \times n$ matrix.

$$A = \begin{bmatrix} 5 & 0 \\ 0 & 10 \end{bmatrix} \qquad ? \qquad A\vec{x} = \lambda \cdot \vec{x}$$

$$\begin{bmatrix} 5 & 0 \\ 0 & 10 \end{bmatrix} \cdot \begin{bmatrix} c \\ 0 \end{bmatrix} = 5 \cdot \begin{bmatrix} c \\ 0 \end{bmatrix}$$

- $\vec{x} \neq 0$ is an *eigenvector* of $A$ if there is a scalar $\lambda$ such that

$$\begin{cases} 5x_1 + 0 \cdot x_2 = \lambda x_1 \\ \overline{0 \cdot x_1 + 10 \cdot x_2 = \lambda x_2} \end{cases}$$

$$A\vec{x} = \lambda\vec{x}$$

$$\lambda = 5: \quad 10 x_2 = 5 x_2$$
$$\Downarrow$$
$$x_2 = 0$$

- the corresponding $\lambda$ is called the *eigenvalue*.

- Example: find the eigenvalue and eigenvector of A.

$$\lambda = 10: \quad x_1 = 0$$

$$\begin{bmatrix} 0 \\ c_2 \end{bmatrix}$$

$$A = \begin{bmatrix} 4 & 2 \\ 1 & 3 \end{bmatrix}$$

$$A X = \lambda X \implies (A - \lambda I) X = 0$$

$\lambda$ be selected such that

$$\det \left( A - \lambda I \right) = 0$$

A $n \times n$ matrix with $n$ linearly independent eigenvectors is said to be **diagonalizable**.

$$A\, u_1 = \lambda_1\, u_1,$$
$$A\, u_2 = \lambda_2\, u_2,$$
$$\dots$$
$$A\, u_n = \lambda_n\, u_n,$$

In matrix form:

$$A\, (u_1 \quad \dots \quad u_n) = (\lambda_1 u_1 \quad \dots \quad \lambda_n u_n) = (u_1 \quad \dots \quad u_n) \begin{pmatrix} \lambda_1 & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \lambda_n \end{pmatrix}$$

This corresponds to a similarity transformation

$$AU = UD \iff A = UDU^{-1}$$

# PCA and eigen-decomposition of covariance matrix.

- Covariance matrix:

$$K[i,j] = COV(X_i - \bar{X}_i, X_j - \bar{X}_j)$$

mean $i$     mean $j$
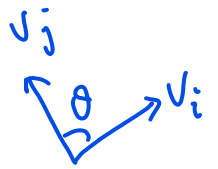
diagonal matrix.

$$K = U \wedge U^{\dagger}$$

Property of covariance matrix:

1. It is symmetric → for symmetric matrix, eigenvectors for distinct eigenvalues are orthogonal.

   $$U = [v_1 \; v_2 \; \cdots \; v_n]$$

   any $i, j$     $v_i, v_j$ are orthogonal.

   $v_j$   $\theta$   $v_i$

   $$\theta = \pm \frac{\pi}{2} \; : \; v_i \cdot v_j = v_i^T v_j = 0$$

2. It is real: -> All eigenvalues of a real symmetric matrix are real.

   linalg · eig(A)     orthonormal vectors     eg.   $v_1 = \begin{bmatrix} 10 \\ 0 \end{bmatrix}$   $v_2 = \begin{bmatrix} 0 \\ 5 \end{bmatrix}$

   $v_i, v_j$     $\|v_i\| = 1$

   $\frac{v}{\|v\|}$

   $v_1' = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$   $v_2' = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$

- Eigen-decomposition of covariance matrix

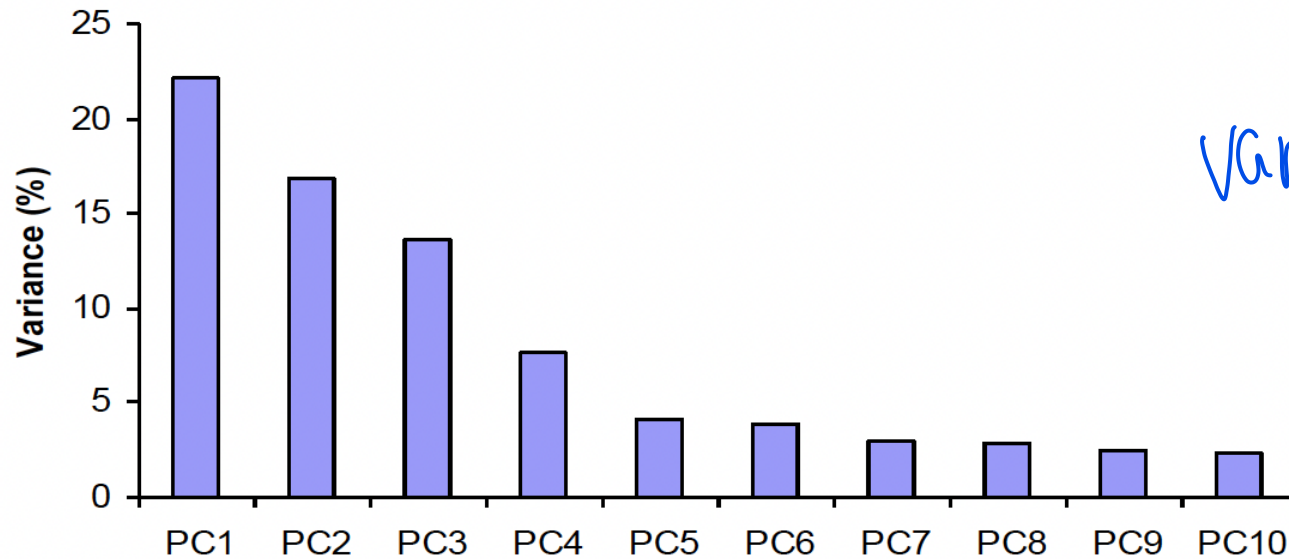$$K = U\Lambda U^{-1}$$

$$K = U\Lambda U^{-1}$$

↳ diagonal

cov after coordinate transformate

$$\begin{bmatrix} \lambda_1 & 0 \\ 0 & \ddots & \lambda_n \end{bmatrix}$$

- Columns of U are *eigenvectors* of $K$.

- Diagonal matrix $\Lambda$ are eigenvalues of $K$, ordered in the order of eigenvectors.

$$\Lambda = diag(\lambda_1, \lambda_2, \ldots, \lambda_n)$$

- We order these eigenvectors in an order of the values of eigenvalues and called these: 1st principal component, 2nd principal component, etc.

- Where does dimensionality reduction come from?
  - Can ignore the components of lesser significance.

$\lambda_1 \cdots \lambda_n$

$\downarrow$

variance of features
after transform

The covariance matrix $K$ is given by:

$X_c = X - \overline{X} \qquad \rightarrow \overline{\mathbb{E}(X_c)} = 0$

$$K = \frac{1}{n-1} X_c^\top X_c$$

$X_c: \qquad t_1 \quad t_2 \quad t_3$

① ② ③ ④ ⑤

$$\begin{bmatrix} 3.2 & -2.85 & -3.15 \\ -2.85 & 4.3 & 4.95 \\ -3.15 & 4.95 & 8.3 \end{bmatrix}$$

eigenvalues, eigenvectors = LA.eig(K)

$$\begin{bmatrix} 13.38070762 & 1.82004592 & 0.59924646 \end{bmatrix}$$

$$\begin{bmatrix} -0.38263617 & 0.77297413 & -0.50606379 \\ 0.53188845 & -0.26357343 & -0.80475072 \\ 0.75543646 & 0.57709622 & 0.31028329 \end{bmatrix}$$

- Project the Data onto the Principal Components:

$$\text{proj}_v \, X = \frac{X \cdot V}{\| V \|} \cdot \vec{v} = (X \cdot V) \cdot \frac{\vec{v}}{\| v \|}$$

normalization

- If we want 2D dimension, project each centered data point into the first two pc:

$$\text{proj}_{V_1} V_2 = \frac{V_2 \cdot V_1}{\| V_1 \|} \cdot \vec{V_1}$$

$$X_c = \begin{bmatrix} -0.2 & -1.4 & -3.6 \\ 1.8 & -2.4 & -0.6 \\ 1.8 & -0.4 & -1.6 \\ -2.2 & 1.6 & 2.4 \\ -1.2 & 2.6 & 3.4 \end{bmatrix}$$

①②③④⑤

$$\begin{bmatrix} -0.38263617 & 0.77297413 & -0.50606379 \\ 0.53188845 & -0.26357343 & -0.80475072 \\ 0.75543646 & 0.57709622 & 0.31028329 \end{bmatrix}$$

$V_1 \qquad V_2 \qquad U_3$

$$X_1 = a_1 V_1 + a_2 V_2 + a_3 V_3$$

$$\text{proj}_{V_1} \vec{X_1} = \text{proj}_{V_1} (a_1 V_1 + a_2 V_2 + a_3 V_3) = a_1 \, \text{proj}_{V_1}^{V_1} + a_2 \, \text{proj}_{V_1} V_2 + a_3 \, \text{proj}_{V_1} V_3$$

$$= a_1 \cdot (V_1 \cdot V_1) \cdot \frac{\vec{V_1}}{\| V_1 \|} + a_2 \cdot 0 + a_3 \cdot 0 = a_1 \vec{V_1}$$

$$a_i \vec{V_i} = \text{proj}_{V_i} X_k \qquad V_i : \text{eigenvector (normalized)}$$
$$\text{orthogonal.}$$

$$X = \sum_{i=1}^{n} a_i \vec{V_i} \qquad n - \text{features}$$

$$\Downarrow$$

$$\tilde{X} = \sum_{i=1}^{K} a_i \vec{V_i} \qquad k < n$$

$$X - \tilde{X} = \sum_{i=k+1}^{n} a_i \vec{V_i}$$