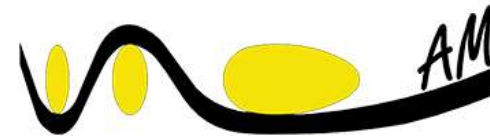# LECTURE 6.  PCA

## (Principal Component Analysis)

### MANU 465
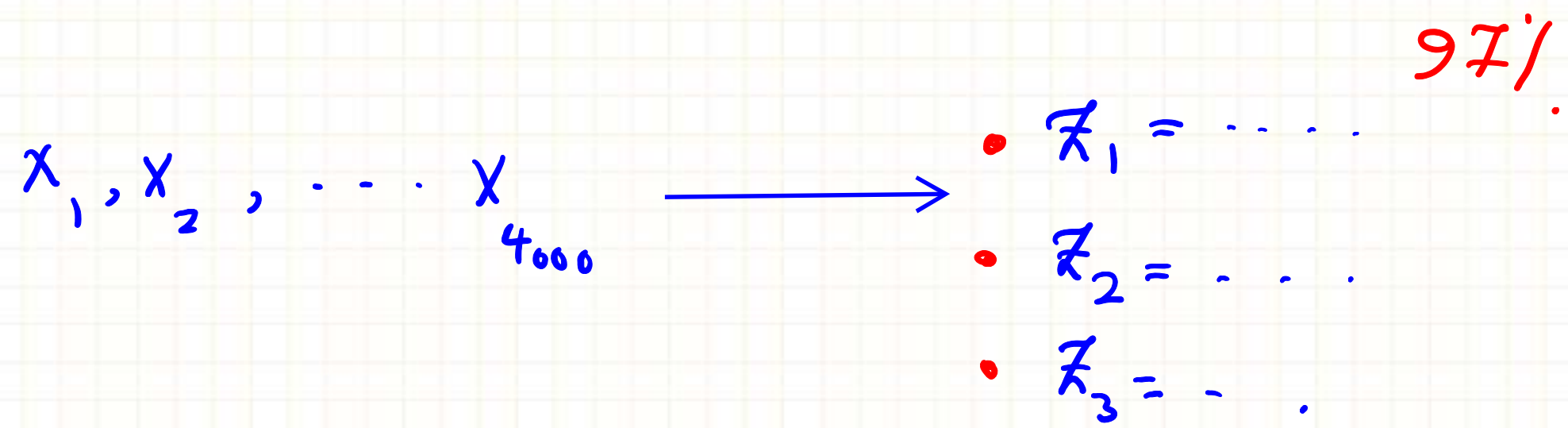
Ahmad Mohammadpanah

PhD, PEng



AIntelligentManufacturing.com

The Prostate, Lung, Colorectal, and Ovarian (PLCO) Cancer Screening Trial is a large, randomized trial designed and sponsored by the National Cancer Institute (NCI) in 2006 to determine the chance of getting PLCO cancer. Participants are being followed and additional data will be collected through 2015, for 216 patients, labeled as Cancer and Healthy, with 4000 features (info related to gene, blood, etc.)

| | $X_1$ | $X_2$ | $X_3$ | $X_4$ | ... | $X_{4000}$ | Result |
|---|---|---|---|---|---|---|---|
| Person 1 | 0.063915364 | 0.025408624 | 0.025536250 | 0.012817321 | | 0.036122788 | Cancer |
| Person 2 | 0.033241734 | 0.051084790 | 0.036122788 | 0.029651841 | | 0.079289645 | Cancer |
| Person 3 | 0.018484138 | 0.056304950 | 0.054195240 | 0.079289645 | | 0.039348744 | Healthy |
| Person 4 | 0.0086176926 | 0.021738490 | 0.0097349817 | 0.050676957 | | 0.039736748 | Cancer |
| Person 5 | 0.035628796 | 0.027409980 | 0.027520513 | 0.039736744 | | 0.039736744 | Healthy |
| Person 6 | 0.037925478 | 0.014913797 | 0.052254751 | 0.057712860 | | 0.0147349818 | Healthy |
| . | ... | | | | | | ... |
| Person 216 | 0.035628796 | 0.027409980 | 0.027520513 | 0.039736744 | | 0.039736475 | Cancer |

$$X_1, X_2, \cdots X_{4000} \longrightarrow$$

$$97\%$$

- $Z_1 = \cdots$
- $Z_2 = \cdots$
- $Z_3 = \cdots$
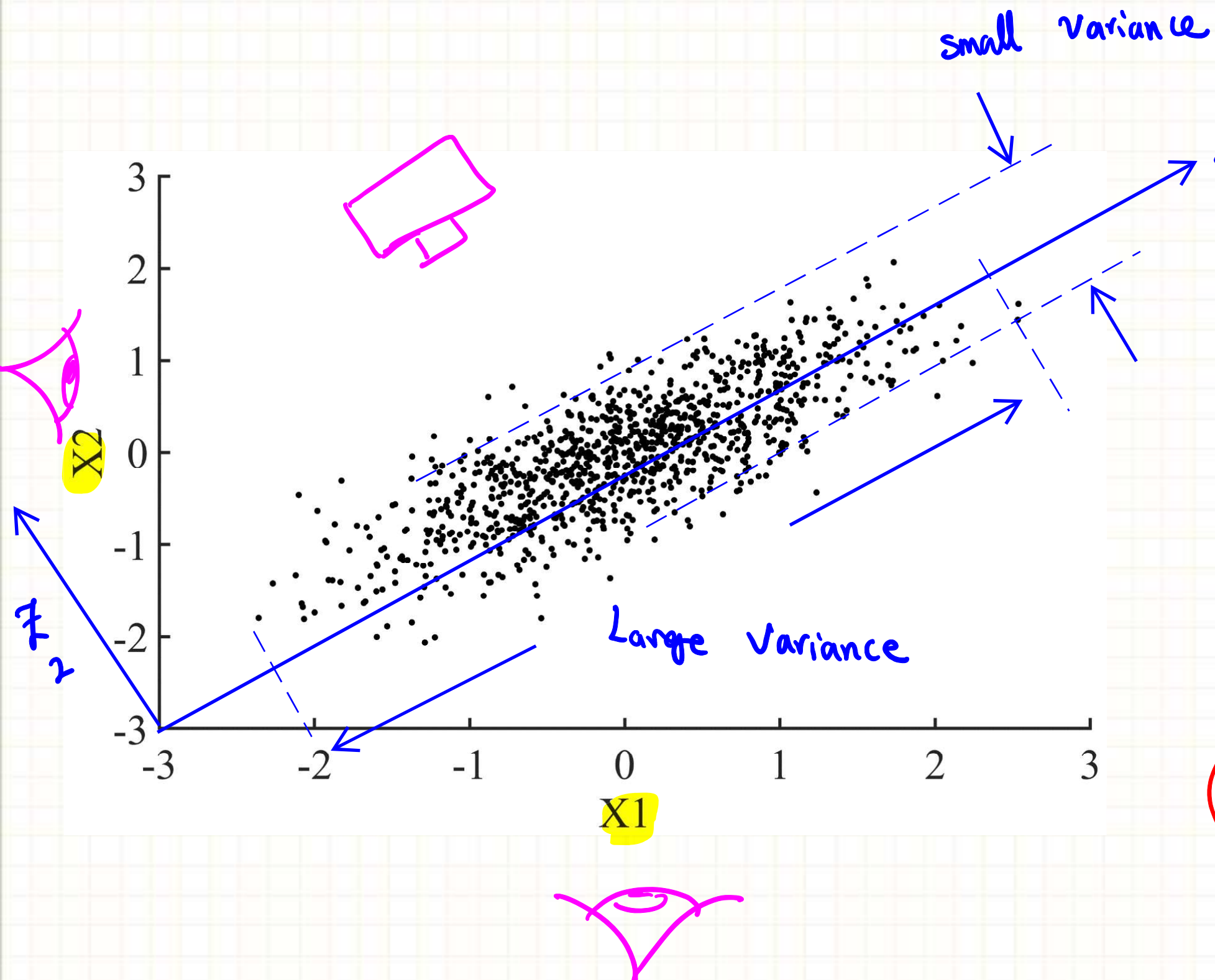
# PCA  (Principal Component Analysis)

**PCA is a powerful tool which:**

- Reduces the dimension of the data.

- Takes 4 or more variables and makes a 2D plot.

- Finds the dominant combinations of variables that describes as much of the data as possible.

## Objective

In this lecture, you will learn <u>what PCA does</u>, <u>how it does it</u>, and <u>how to apply it</u> to get deeper insight into your data.

# What does PCA do?

small variance

$Z_1$

Large Variance

$X2$

$X1$

$Z_2$

PCA will find the direction of max variance.

$\lambda_1 = Var(Z_1)$

$\lambda_2 = Var(Z_2)$

if you only use $Z_1$, then

$$\frac{\lambda_1}{\lambda_1 + \lambda_2} \times 100\%$$ of information
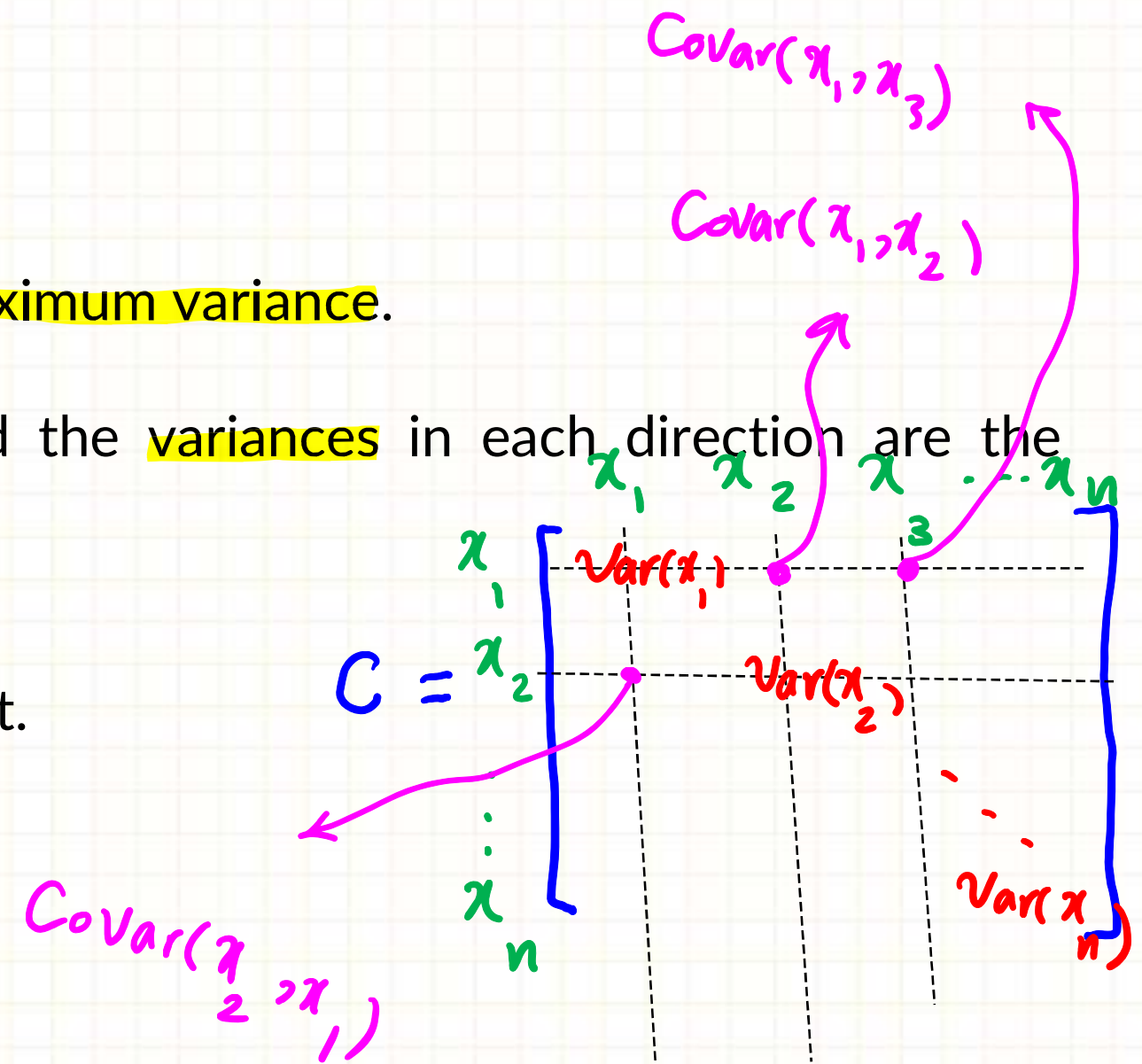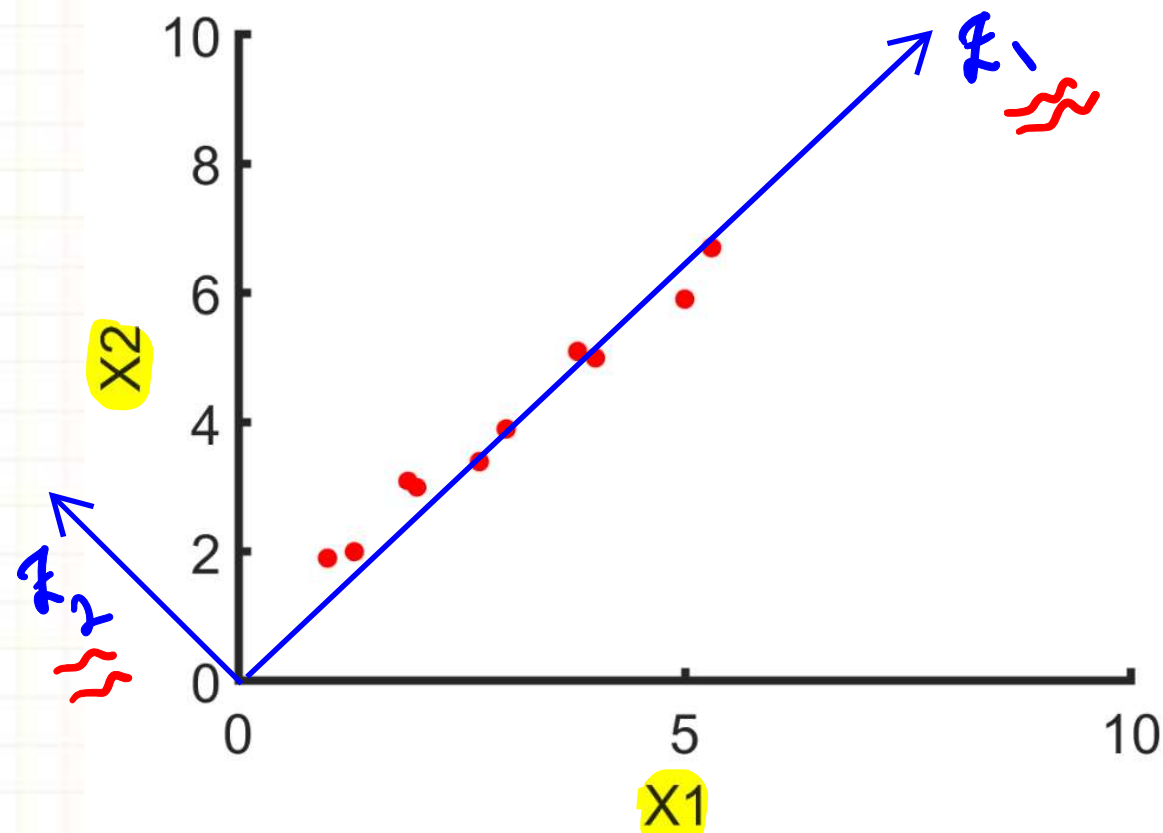
preserves.

# How does it work?

o Fundamentally, PCA transform the data into directions of maximum variance.

o Mathematically, these directions are the Eigenvectors; and the variances in each direction are the corresponding Eigenvalues of the covariance matrix of data.

**Example**) Find the direction of maximum variance for this dataset.

| X1 | X2 |
|-----|-----|
| 3.8 | 5.1 |
| 2 | 3 |
| 1.9 | 3.1 |
| 5 | 5.9 |
| 1.3 | 2 |
| 3 | 3.9 |
| 5.3 | 6.7 |
| 4 | 5 |
| 1 | 1.9 |
| 2.7 | 3.4 |



$$Covar(x_1, x_3)$$

$$Covar(x_1, x_2)$$

$$C = \begin{array}{c} x_1 \\ x_2 \\ \vdots \\ x_n \end{array} \begin{bmatrix} Var(x_1) & & & \\ & Var(x_2) & & \\ & & \ddots & \\ & & & Var(x_n) \end{bmatrix}$$

$$x_1 \quad x_2 \quad x_3 \cdots x_n$$

$$Covar(x_2, x_1)$$

$$Var(x_1) = \mathbf{2.23}$$

$$Var(x_2) = \mathbf{2.63}$$

$$Cov(x_1, x_2) = \mathbf{2.40} = Cov(x_2, x_1)$$

$$Var(x) = \frac{\sum_{i=1}^{n}(x_i - \bar{X})^2}{n-1} \quad , \quad Covar(x_1, x_2) = \frac{\sum_{i=1}^{n}(x_i^1 - \bar{X}_1)(x_i^2 - \bar{X}_2)}{n-1}$$

$$C = \begin{bmatrix} 2.23 & 2.4 \\ 2.4 & 2.63 \end{bmatrix} ,$$

finding Eigenvalues & Eigenvectors :

$$|C - \lambda I| = 0 \longrightarrow \boxed{\begin{array}{l} \lambda_1 = 4.84 \\ \lambda_2 = 0.01 \end{array}}$$

if we only use $z_1$ information,

$$\frac{\lambda_1}{\lambda_1 + \lambda_2} \times 100\% = 99\%$$

of information still are available.

$$\begin{cases} [C - \lambda I]\begin{bmatrix} u \\ v \end{bmatrix} = 0 \\ u^2 + v^2 = 1 \end{cases} \longrightarrow$$

for $\lambda = \lambda_1 \longrightarrow V_1 = \begin{bmatrix} 0.67 \\ 0.73 \end{bmatrix}$

for $\lambda = \lambda_2 \longrightarrow V_2 = \begin{bmatrix} -0.73 \\ 0.67 \end{bmatrix}$

eigenvectors

the new variables:

$$Z_1 = [X] V_1 = \begin{array}{cc} X_1 & X_2 \\ \begin{bmatrix} 3.8 & 5.1 \\ 2 & 3 \\ 1.9 & 3.1 \\ 5 & 5.9 \\ 1.3 & 2 \\ 3 & 3.9 \\ 5.3 & 6.7 \\ 4 & 5 \\ 1 & 1.9 \\ 2.7 & 3.4 \end{bmatrix} \end{array} \begin{bmatrix} 0.67 \\ 0.73 \end{bmatrix} = \begin{bmatrix} 6.3 \\ 3.5 \\ 3.5 \\ 7.7 \\ 2.3 \\ 4.9 \\ 8.5 \\ 6.3 \\ 2 \\ 4.3 \end{bmatrix} ,$$

PC1

our input matrix

$$Z_2 = [X] V_2 = \begin{bmatrix} 0.65 \\ 0.56 \\ 0.7 \\ 0.3 \\ 0.39 \\ 0.43 \\ 0.63 \\ 0.31 \\ 0.51 \\ 0.44 \end{bmatrix}$$
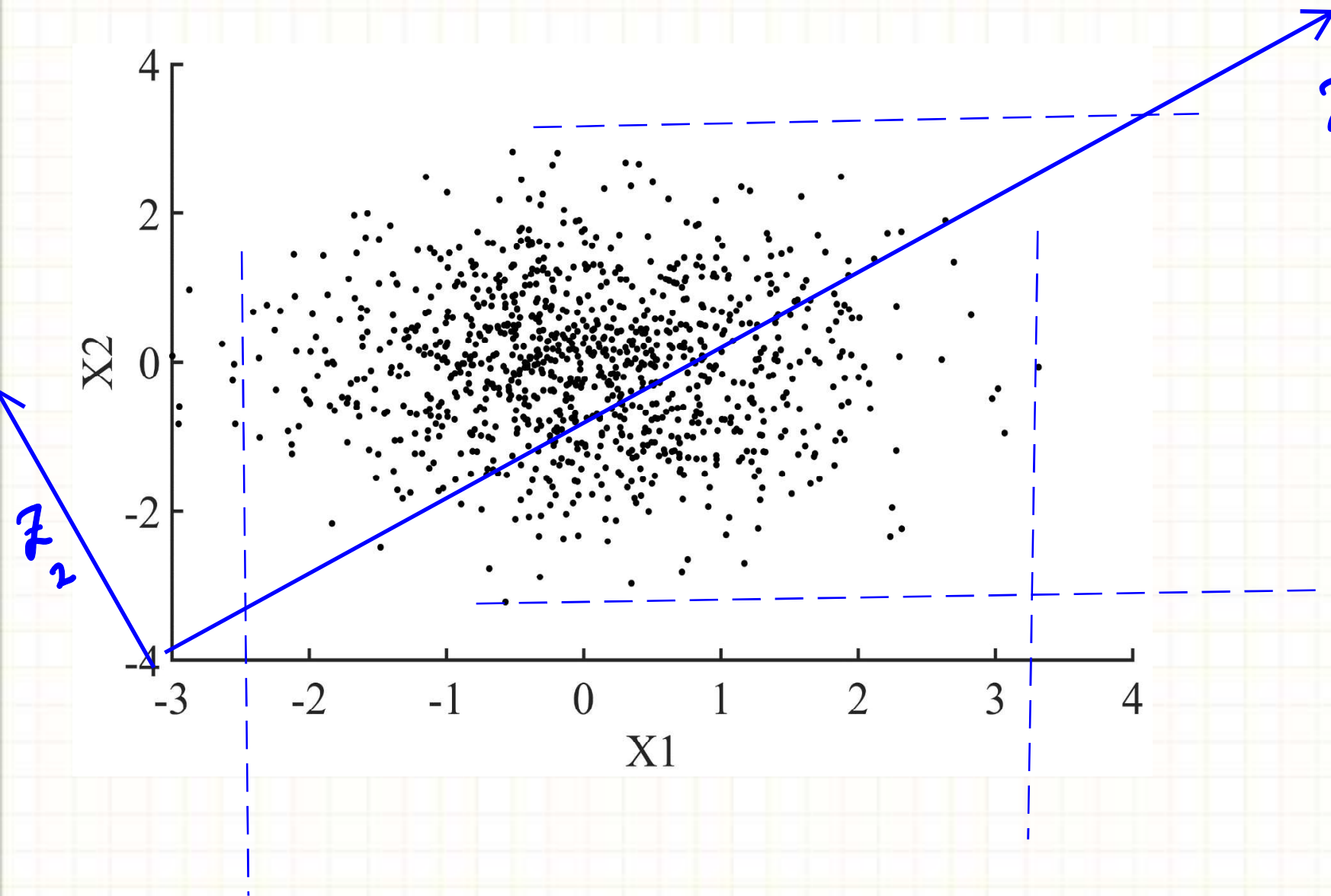
PC2

Variance of the new variables are the same as Eigenvalues.

$Var(Z_1) = 4.84$

$Var(Z_2) = 0.01$

**Does PCA always work?**

$$Var(x_1) + Var(x_2) = Var(z_1) + Var(z_2)$$



$$z_1 = c_1 x_1 + c_2 x_2$$
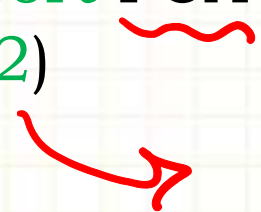
$$\lambda_1$$

$$\lambda_2$$

$$\vdots$$

$z_2$

# How to apply it?

1. Compute the Covariance Matrix C for X (the input data).

2. Compute the Eigenvectors ($V_i$) and Eigenvalues ($\lambda_i$) of C.

3. Sort the eigenvalues from the maximum to minimum ($\lambda_1 > \lambda_2 > \lambda_3 > \cdots . > \lambda_n$)

4. First Principal Component is PC1=X*V$_1$ , and PC2=X*V$_2$ , ...

5. We may just keep the first few PCs and chop of the rest (shrinking the dimension of the data).

6. Contribution of the $m$ PCs we used to describe the data$=\dfrac{\sum_1^m \lambda_i}{\sum_1^n \lambda_i}$

10 variables $\longrightarrow$ $\lambda_1, \lambda_2, \ldots, \lambda_{10}$

only keep

$\lambda_1$ & $\lambda_2$ $\longrightarrow$ $\dfrac{\lambda_1 + \lambda_2}{\lambda_1 + \lambda_2 + \cdots + \lambda_{10}} \times 100\%$    information preserve.
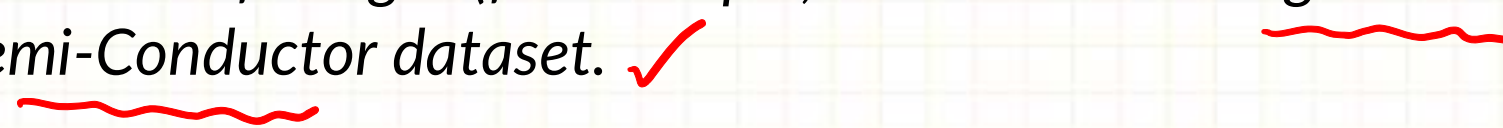
**In Python, SKLearn, there is a built-in Function:**

```python
from sklearn.decomposition import PCA
PrinComp=PCA(n_components=2)
PrinComp.fit(X)
Z=PrinComp.transform(X)
```

Please check the code, Wood_PCA.ipynb on Canvas.

## Practice:

1. Assignment 5. ✓
2. Apply this method to a set of images (for example, the Fashion or Digit MNIST).
3. Apply this to the Semi-Conductor dataset. ✓
4. Attend Tutorial 5.

# Summary

- PCA is a linear transformation which find the direction of maximum variance in the data.

- It reduces the dimension of the data without loosing much information.

- It is a good tool for data visualization.

- In practice, we can simply, use the PCA built-in function in Python to apply PCA to a dataset.