

# Data Mining and Machine Learning

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Introducti

ing.

Add WeChat edu\_assist\_pro

Vector Data Analy

Principal

Components Analysis (PCA)



# Objectives

- To introduce Data Mining
- To outline the techniques that we will study in this part of the course – a Data Mining ‘Toolkit’
- To review <https://eduassistpro.github.io/> the notions of mean, variance
- To explain Principal Component Analysis (PCA)
- To present an example of PCA



# What is Data Mining?

- Mining

- *Digging deep into the earth, to find hidden, valuable materials*

- Data Mining <https://eduassistpro.github.io/>

- Analysis of large data c acoustic, video, text,... medical, structure, patterns and relationships
  - Corpora which are too large for human inspection
  - Patterns and structure may be hidden



# Data Mining

- Structure and patterns in large, abstract data sets:
  - Is the data homogeneous or does it consist of several separately identifiable subsets?
  - Are there patterns?
  - If so, do they have a meaningful interpretation?
  - Are there correlations in the data?
  - Is there redundancy in the data?



# Data Mining

- In this part of the course we will develop a basic ‘data mining toolkit’
  - Subspace projection methods (PCA)
  - Clustering
  - Statistical
  - Sequence analysis
    - Dynamic Programming (DP)

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro



# Some example data

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

Fig 1: Single,  
spherical  
cluster centred  
at origin

Fig 2: Sing  
arbitrary elliptical  
cluster

Fig 3: Multiple,  
arbitrary elliptical  
clusters

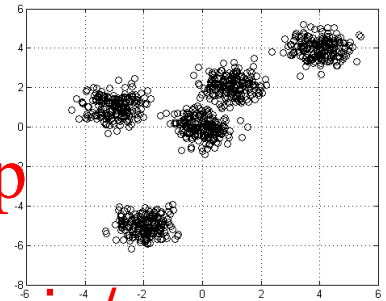


# Objectives

- Fig 3 shows “multiple source” data.

The data is arranged in a set of

“clusters”.



- How do we find the locations of

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

- Remember, in real applications will be many points in a high-dimensional vector space which is difficult to visualise



# Objectives

- Fig 1 shows simplest type of data – single source data centred at origin. Equal variance in both dimensions and no covariance.

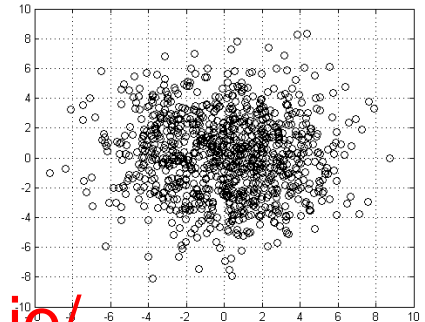


Fig 1

- Fig 2 is again single source data centred at origin. It is correlated and has unequal variance in the two dimensions.

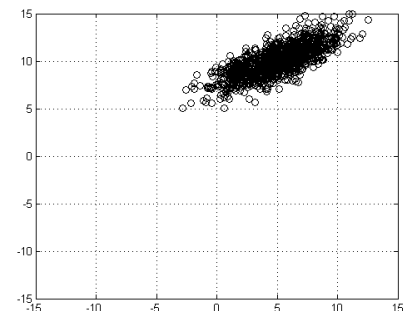


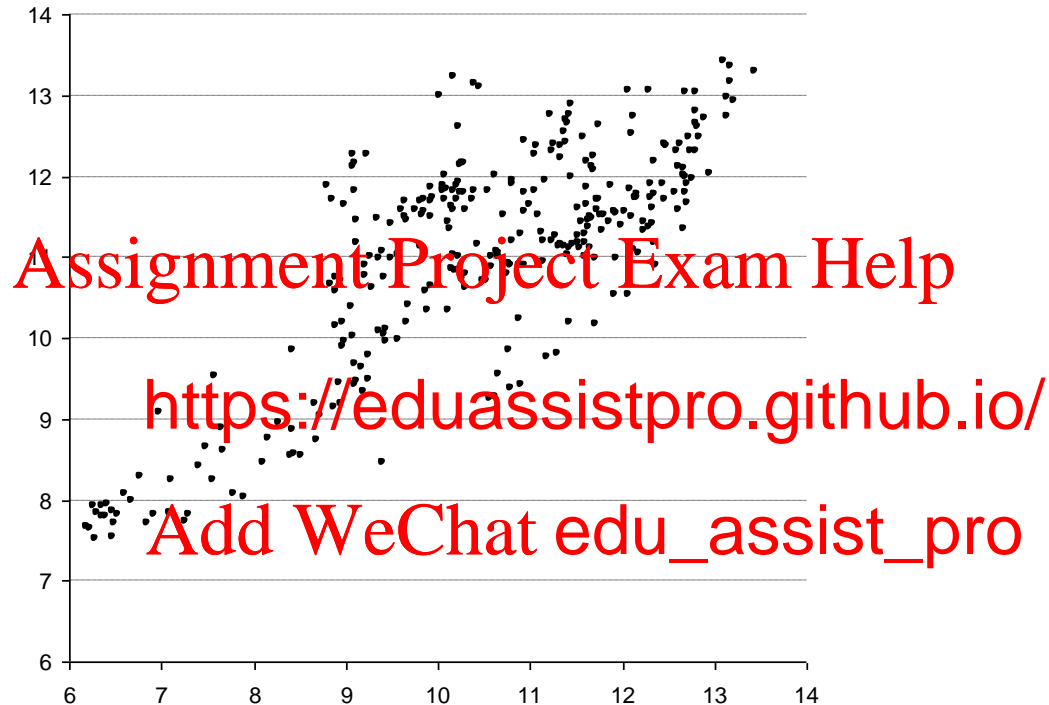
Fig 2

- How do we convert Fig 2 into Fig 1?
- We will start with this problem
- Solution is a technique called Principal Components Analysis (PCA)





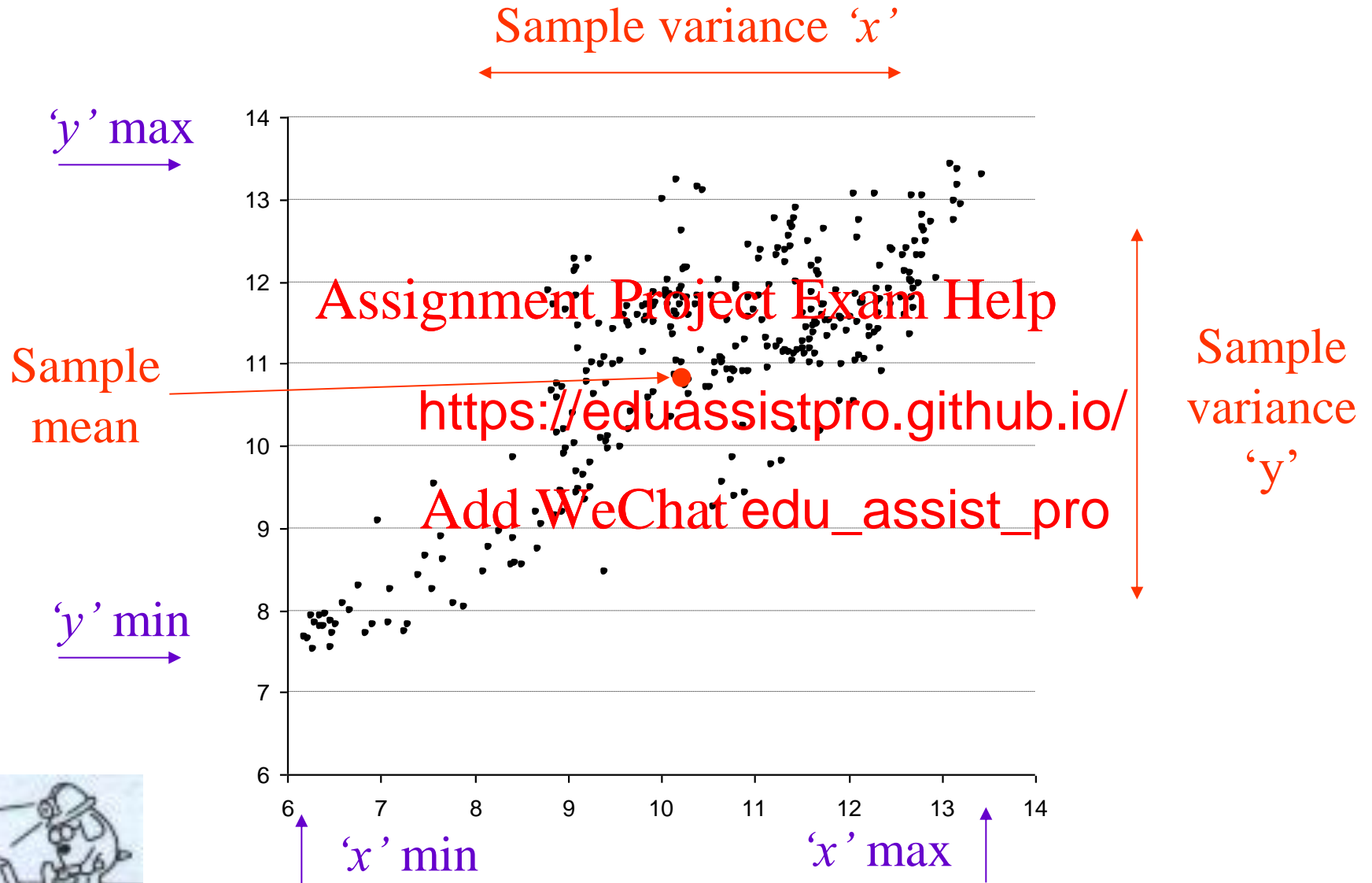
# Example from speech processing



*Plot of high-frequency energy vs low-frequency energy, for 25 ms speech segments, sampled every 10ms*



# Basic statistics



# Basic statistics

- Denote samples by

$$X = x_1, x_2, \dots, x_T$$

where  $x_t = (x_t^1, x_t^2, \dots, x_t^N)$

- The sample mean  $\mu(X)$  vector is given by:

$$\mu^n = \frac{1}{T} \sum_{t=1}^T x_t^n$$

$$\mu = (\mu^1, \mu^2, \dots, \mu^n, \dots, \mu^N)$$



# More basic statistics

- The sample variance  $\sigma$  (more correctly  $\sigma(X)$ ) vector is given by:

Assignment Project Exam Help

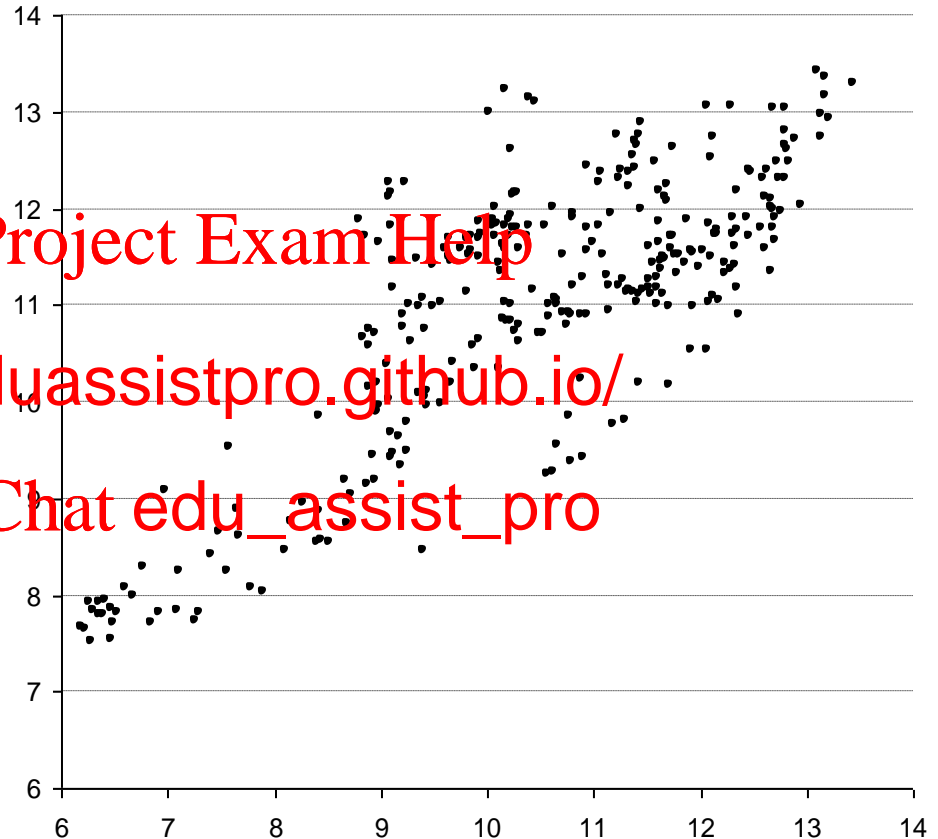
<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro



# Covariance

- In this data, as the  $x$  value increases, the  $y$  value also increases
- This is (positive) co-variance
- If  $y$  decreases as  $x$  increases, the result is negative covariance



# Definition of covariance

- The covariance between the  $m^{th}$  and  $n^{th}$  components of the sample data is defined by:

Assignment Project Exam Help

<https://eduassistpro.github.io/>

- In practice it is useful to subtract the mean  $\mu$  from each of the data points  $x_t$ . The mean is then 0 and

$$\sigma^{m,n} = \frac{1}{T-1} \sum_{t=1}^T x_t^m x_t^n,$$



# The covariance matrix

$$\sigma = \begin{bmatrix} \sigma^{1,1} & \sigma^{1,2} & \dots & \sigma^{1,n} & \dots & \sigma^{1,N} \\ \sigma^{2,1} & \sigma^{2,2} & \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots & \dots & \dots \\ \sigma^{m,1} & \dots & \dots & \sigma^{m,n} & \dots & \dots \\ \dots & \dots & \dots & \dots & \dots & \dots \\ \sigma^{N,1} & \dots & \dots & \dots & \dots & \sigma^{N,N} \end{bmatrix}$$

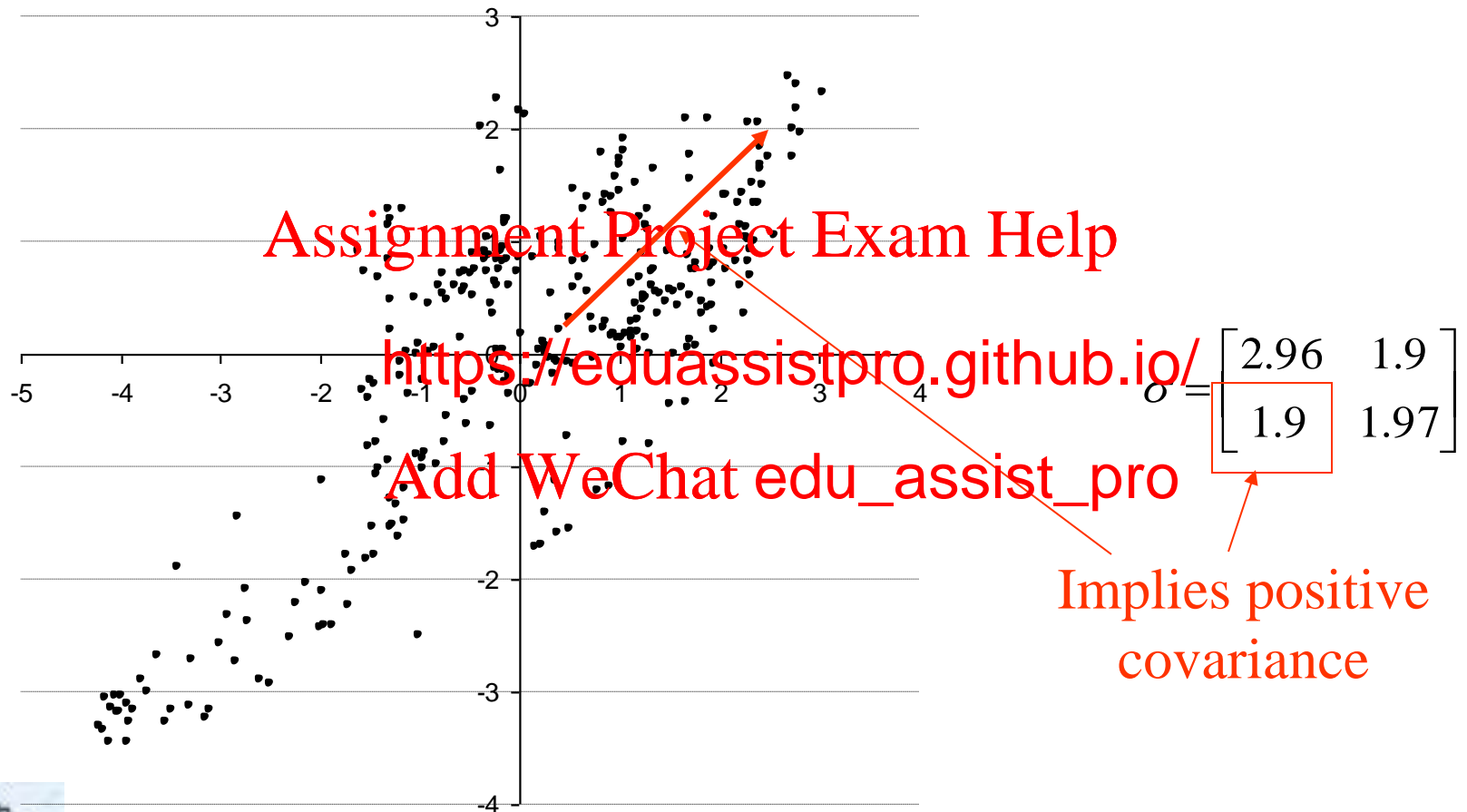
Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

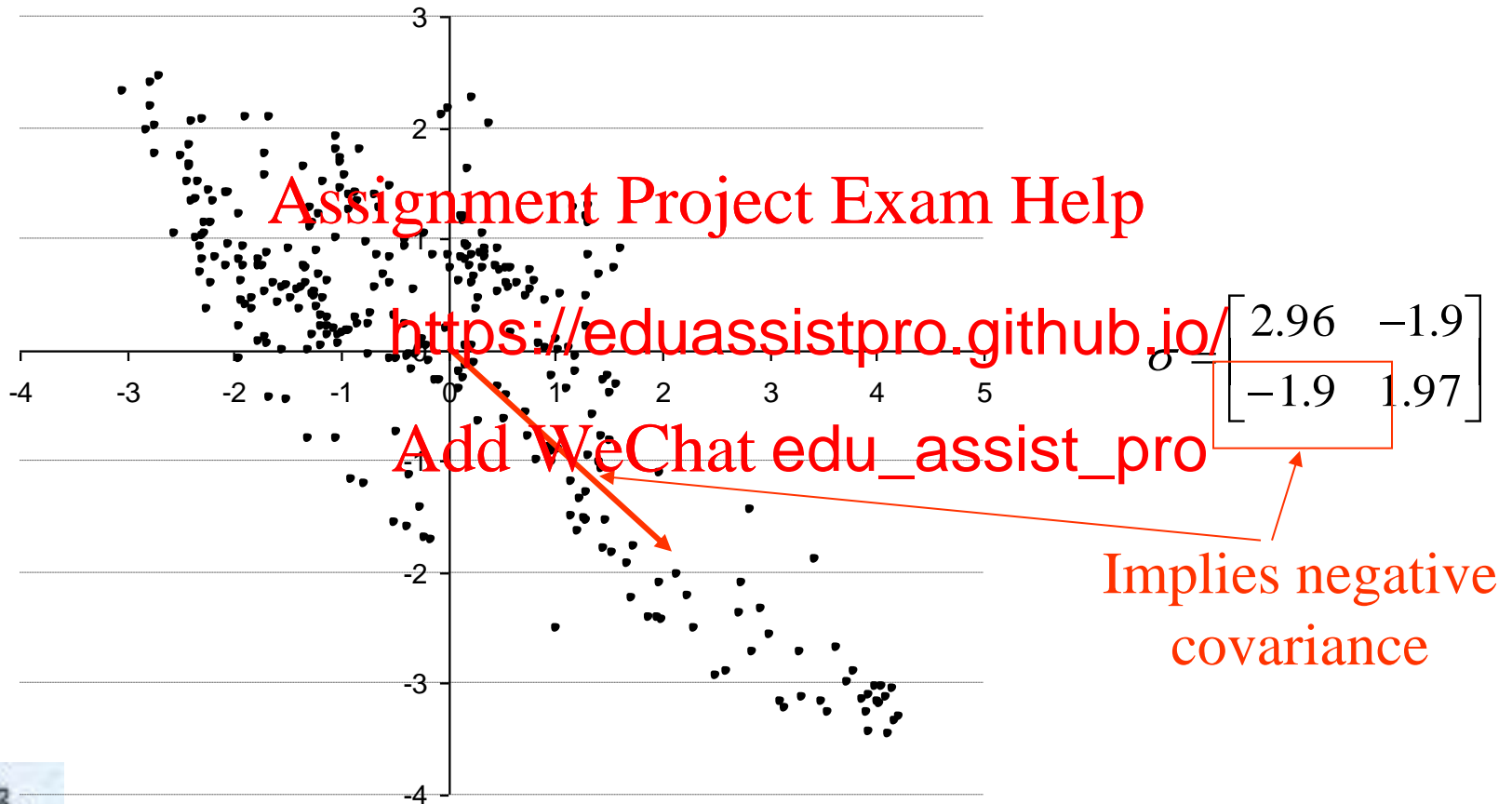


# Data with mean subtracted

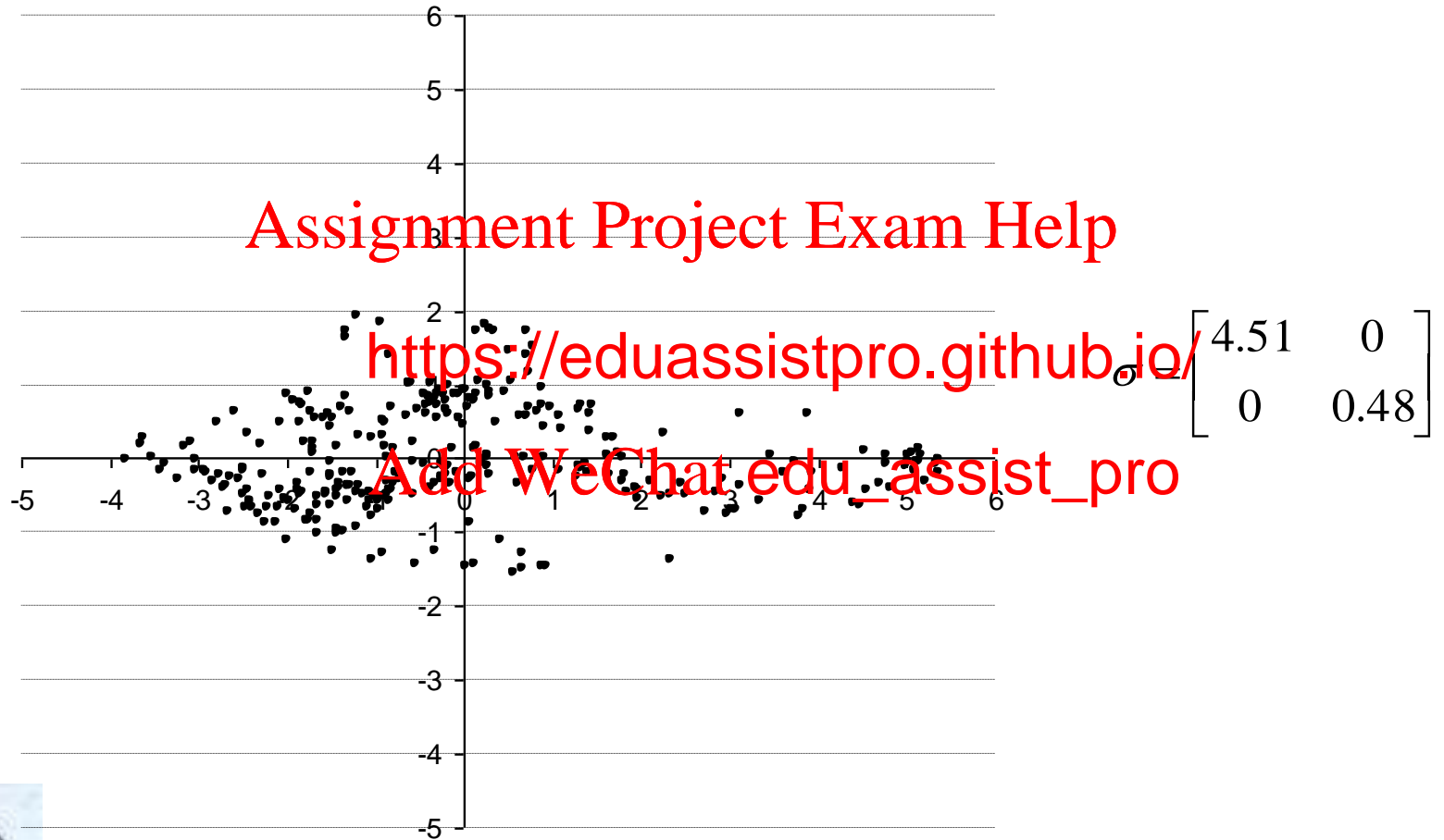




# Sample data rotated



# Data with covariance removed



# Principal Components Analysis

- PCA is the technique which I used to diagonalise the sample covariance matrix
- The first step is to write the covariance matrix in the form:

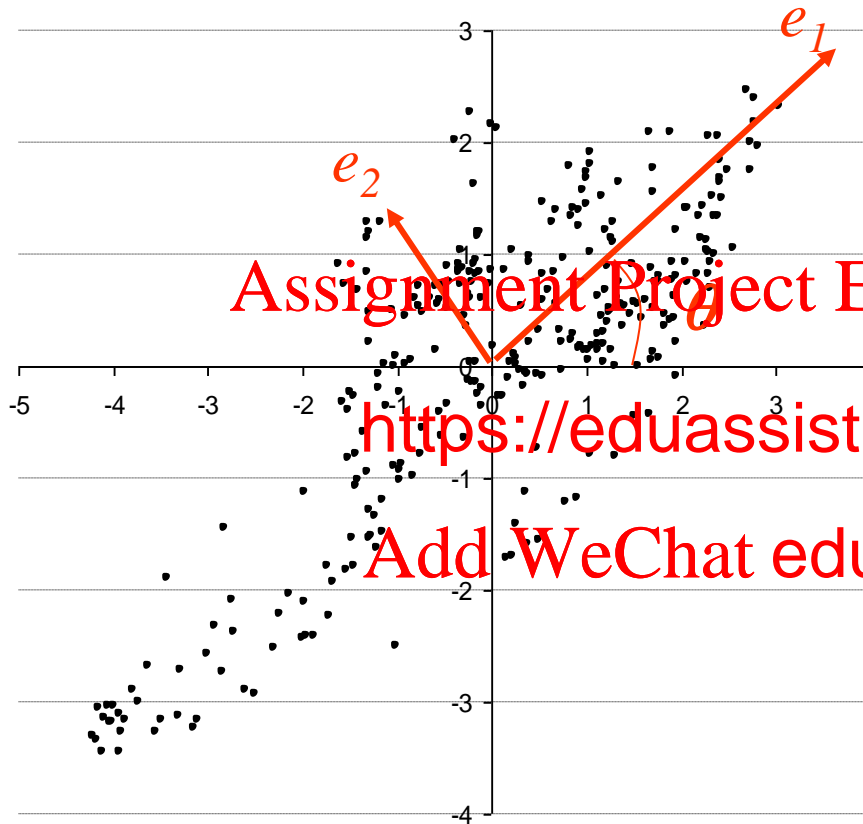
$$\sigma = UDU^T$$

where  $D$  is diagonal and  $U$  is a matrix corresponding to a rotation

- You can do this using SVD (see lecture on LSI) or Eigenvalue Decomposition



# PCA continued



$U$  implements rotation through angle  $\theta$

$e_1$  is the first column of  $U$

$$e_1 = \begin{bmatrix} u_{11} \\ u_{21} \end{bmatrix}$$

is the variance in the

direction  $e_1$

the 2<sup>nd</sup> column of  $U$

$d_{22}$  is the variance in the direction  $e_2$

$$\sigma = UDU^T = \begin{bmatrix} u_{11} & u_{12} \\ u_{21} & u_{22} \end{bmatrix} \begin{bmatrix} d_{11} & 0 \\ 0 & d_{22} \end{bmatrix} \begin{bmatrix} u_{11} & u_{21} \\ u_{12} & u_{22} \end{bmatrix}$$



# PCA Example

- Abstract data set

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro



# PCA Example (continued)

- Step 1: load the data into MATLAB:
  - `A=load('data4');`
- Step 2: Calculate the mean and subtract this from each sample
  - `M=ones(` <https://eduassistpro.github.io/>
  - `N=mean(A);` [Add WeChat edu\\_assist\\_pro](#)
  - `M(:,1)=M(:,1)*N(1)`
  - `M(:,2)=M(:,2)*N(2);`
  - `B=A-M;`

■ Plot B



# PCA Example (continued)

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro



# PCA Example (continued)

- Calculate the covariance matrix of B (or A)

- $S = (B' * B) / \text{size}(B, 1) ;$

- or **Assignment Project Exam Help**

- $S = \text{cov}(B)$  **<https://eduassistpro.github.io/>**

$$S = \begin{bmatrix} 6.78 & 3.27 \\ 3.27 & 2.76 \end{bmatrix}$$

**Add WeChat edu\_assist\_pro**

- Difficult to deduce much about the data from this covariance matrix





# PCA Example (continued)

- Calculate the eigenvalue decomposition of  $S$

- $[U, E] = \text{eig}(S)$  ;

Assignment Project Exam Help

$$U = \begin{bmatrix} 0.4884 & -0.8726 \\ -0.8726 & -0.4884 \end{bmatrix}, E = \begin{bmatrix} 0.9307 & 0 \\ 0 & 8.6079 \end{bmatrix}$$

<https://eduassistpro.github.io/>  
Add WeChat edu\_assist\_pro

- After transforming the data using  $U$  its covariance matrix becomes  $E$ . You can confirm this by plotting the transformed data:



# PCA Example (continued)

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro



# PCA Example (continued)

- After transformation by the matrix  $U$ , the covariance matrix has been diagonalized and is now equal to  $E$ 
  - variance in the  $x$  direction is 0.93
  - variance in
- This tells us that the data is contained in the (new)  $y$  direction
- There is much less variation in the  $x$  direction, and we could get a 1 dimensional approximation to the data by discarding this dimension
- None of this is obvious from the original covariance matrix

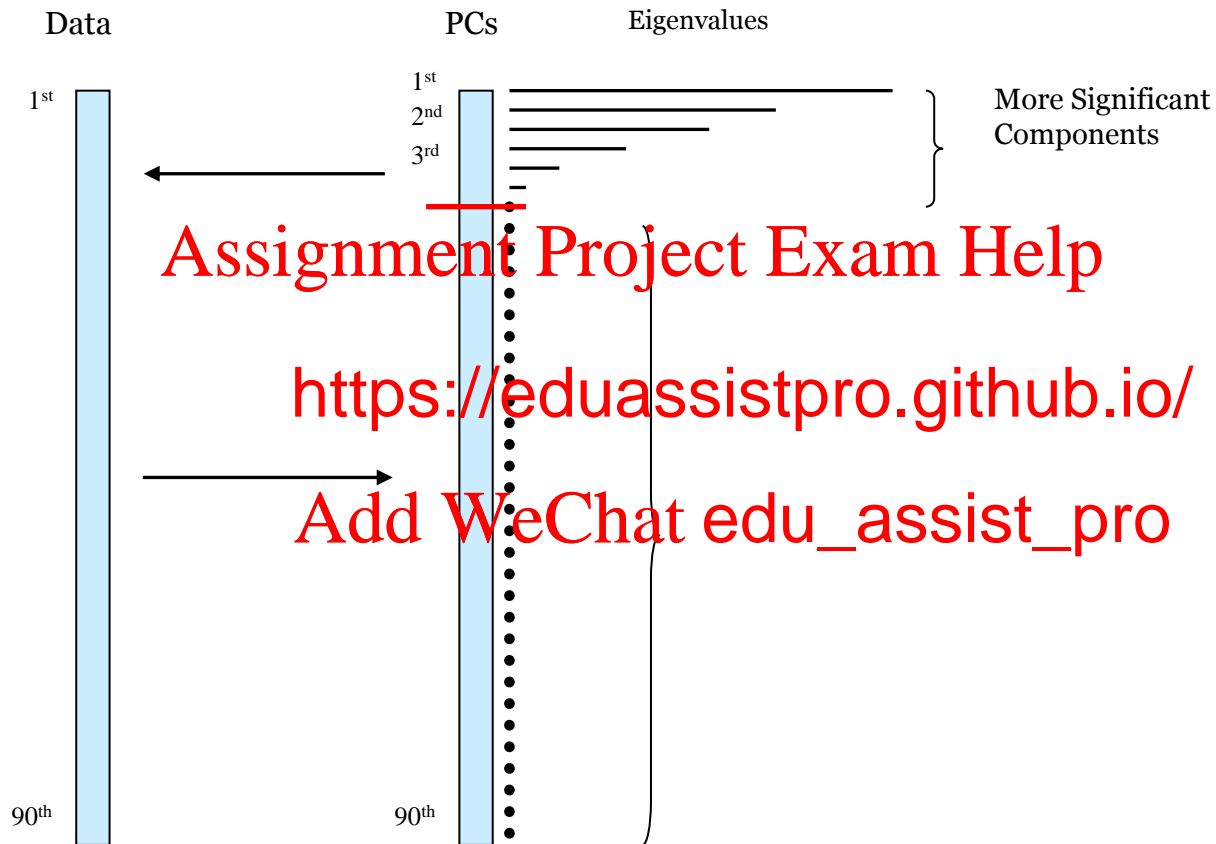


# Final notes

- Each column of  $U$  is a principal vector
- The corresponding eigenvalue indicates the variance of the data along that dimension
  - Large eigenvalues indicate the most components of the data
  - Small eigenvalues indicate the least variation along the corresponding eigenvectors
- It may be advantageous to ignore dimensions which correspond to small eigenvalues and only consider the projection of the data onto the most significant eigenvectors – this way the dimension of the data can be reduced

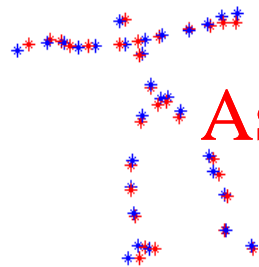


# Eigenvalues



# Visualising PCA

Original pattern (blue)



$U$

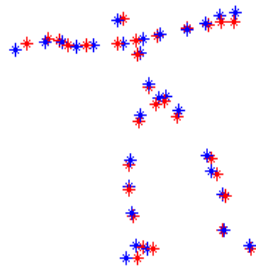
Eigenspace

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Set coordinates  
 $n \rightarrow 90$  to zero

Reduced pattern (red)



$U^{-1}$

Eigenspace

Add WeChat edu\_assist\_pro



# Summary

- Review of basic data analysis (mean, variance and covariance)

Assignment Project Exam Help

- Introduction <https://eduassistpro.github.io/> Analysis (PCA)

Add WeChat edu\_assist\_pro

- Example of PCA

