

Big Data Fundamentals and Applications

# Principal Component Analysis – PCA

**Asst. Prof. Chan, Chun-Hsiang**

*Master program in Intelligent Computing and Big Data, Chung Yuan Christian University, Taoyuan, Taiwan*

*Undergraduate program in Intelligent Computing and Big Data, Chung Yuan Christian University, Taoyuan, Taiwan*

*Undergraduate program in Applied Artificial Intelligence, Chung Yuan Christian University, Taoyuan, Taiwan*

# Outlines

1. Review
2. Why Do We Need Dimension Reduction?
3. PCA – Mathematical Explanation
4. PCA – Programming in Python
5. Part VIII Basic numerical methods (Scipy)  
Part IX Advanced numerical methods (Scipy)
6. Notice

# Review

- Before we explain PCA, we need to review the mathematical meaning of three basic descriptive statistics, including expectation, variance, and covariance.
- In previous courses or your understanding, these three parameters usually perform as above equations.

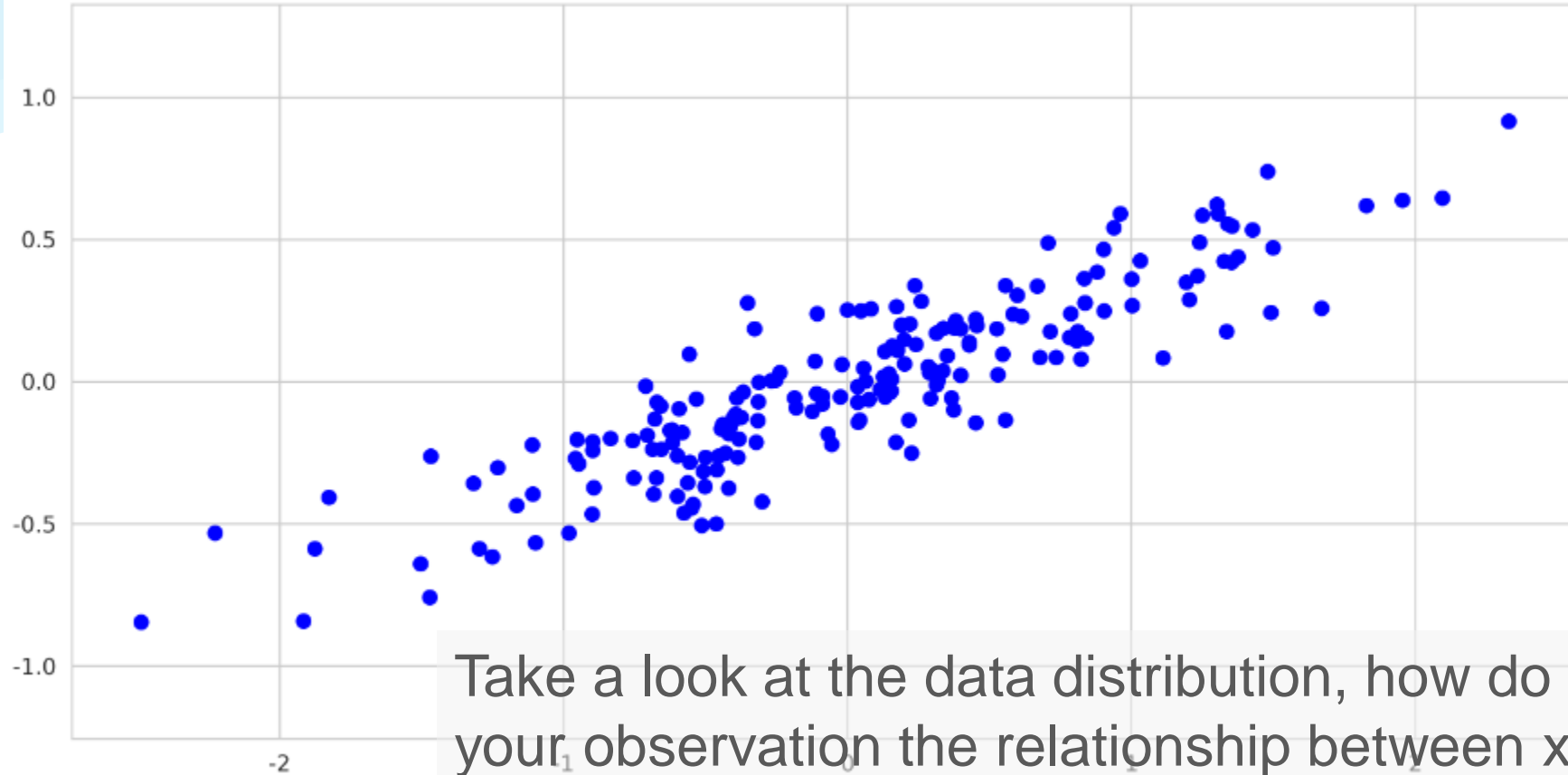
$$\mu = \frac{1}{n} \sum_{i=1}^n x_i$$

$$\text{var}(x) = \sigma^2 = \left( \sqrt{\frac{\sum_{i=1}^n (x_i - \mu)^2}{n}} \right)^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2$$

$$\text{Cov}(x, y) = \frac{1}{n} \sum_{i=1}^n (x_i - \mu_x)(y_i - \mu_y)$$

# Review

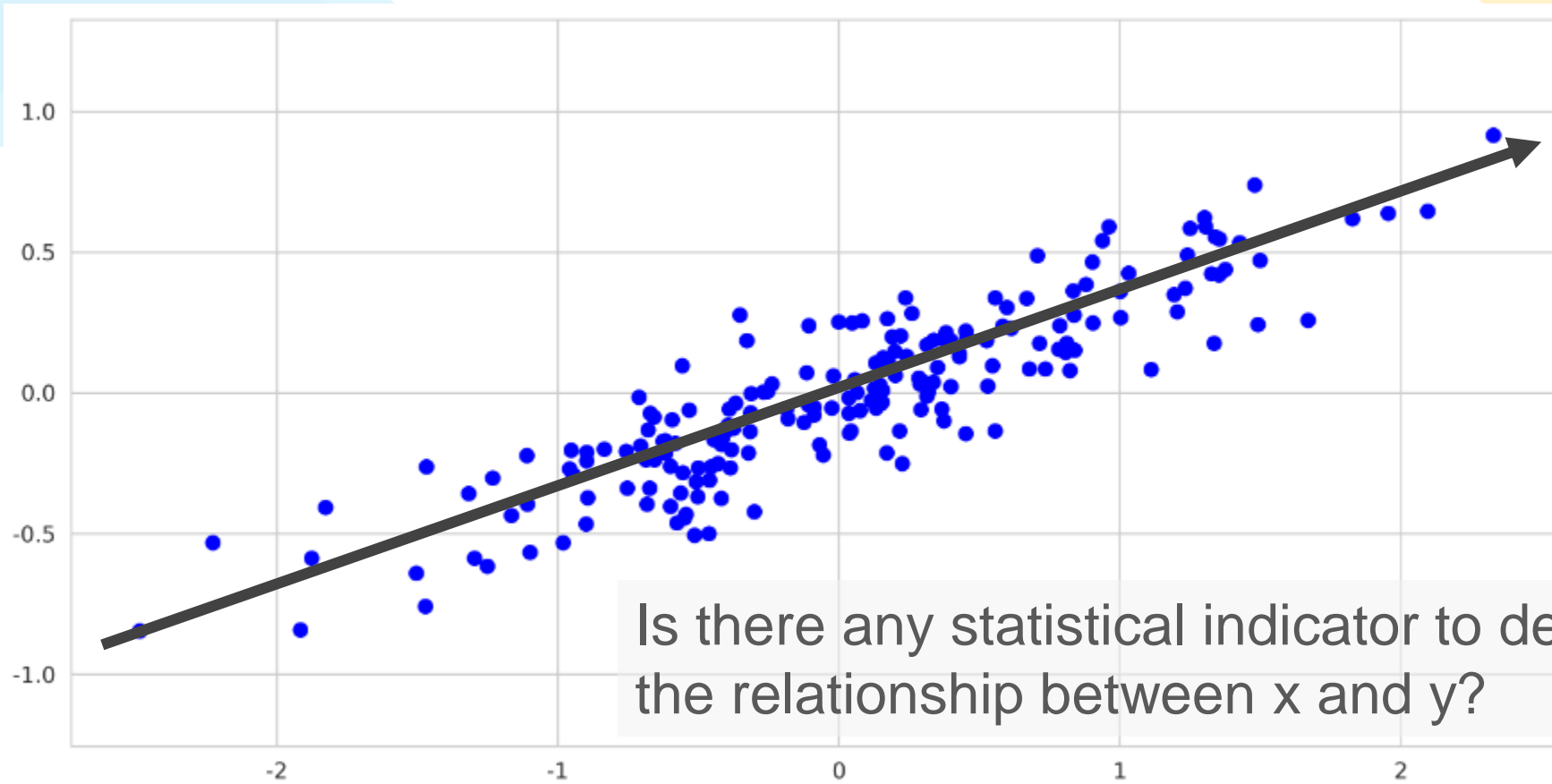
$$\mu = \frac{1}{n} \sum_{i=1}^n x_i \quad | \quad \text{var}(x) = \sigma^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2 \quad | \quad \text{Cov}(x, y) = \frac{1}{n} \sum_{i=1}^n (x_i - \mu_x)(y_i - \mu_y)$$



Take a look at the data distribution, how do you explain your observation the relationship between x and y?

# Review

$$\mu = \frac{1}{n} \sum_{i=1}^n x_i \quad | \quad \text{var}(x) = \sigma^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2 \quad | \quad \text{Cov}(x, y) = \frac{1}{n} \sum_{i=1}^n (x_i - \mu_x)(y_i - \mu_y)$$



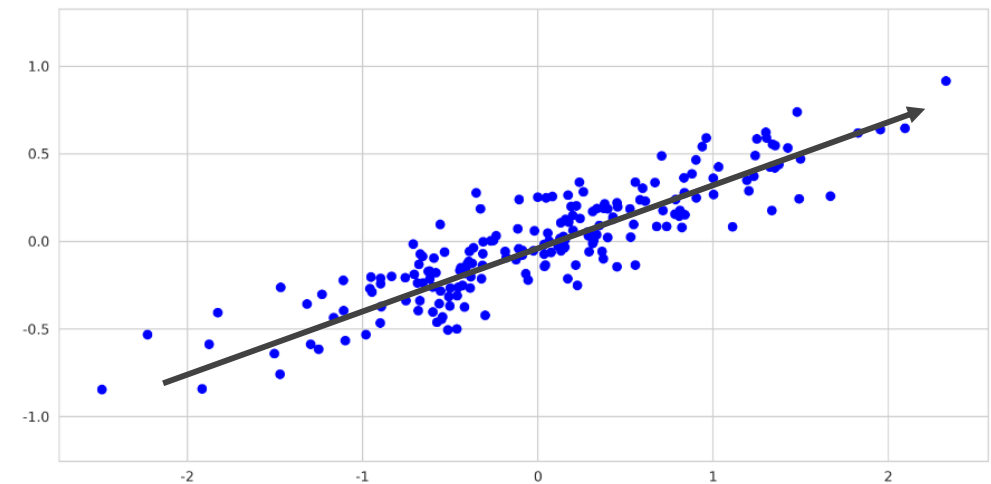
Is there any statistical indicator to describe the relationship between x and y?

# Review

$$\mu = \frac{1}{n} \sum_{i=1}^n x_i \quad | \quad \text{var}(x) = \sigma^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2 \quad | \quad \text{Cov}(x, y) = \frac{1}{n} \sum_{i=1}^n (x_i - \mu_x)(y_i - \mu_y)$$

- Pearson's correlation coefficient
- Given two parameters  $x_i$  and  $y_i$ , where  $i$  ranges from 1 to  $n$ . Then, Pearson's correlation coefficient could be defined as follows.

$$\rho = \frac{\sum_{i=1}^n (x_i - \mu_x)(y_i - \mu_y)}{\sqrt{\sum_{i=1}^n (x_i - \mu_x)^2 \sum_{i=1}^n (y_i - \mu_y)^2}}$$



## Question 1

If  $x$  is highly correlated with  $y$ , and then what do you expect from their covariance and standard deviations?

# Review

- **Expectation**

$$E(X) = \sum_x xP(X = x) = \mu$$

- **Variance**

$$\begin{aligned} \text{var}(X) &= E([X - \mu]^2) \\ &= E(X^2 - 2\mu X + \mu^2) \\ &= E(X^2) - 2\mu E(X) + \mu^2 \\ &= E(X^2) - 2\mu^2 + \mu^2 \\ &= E(X^2) - \mu^2 \\ &= E(X^2) - E(X)^2 \end{aligned}$$

## Characteristics of Expectation

$E(aX + bY) = aE(X) + bE(Y)$ ,  $a, b \in \mathbb{R}$   
X and Y can be independent or dependent.

$$E(XY) = E(X)E(Y)$$

Where  $\text{cov}(X, Y) = 0$

# Review

- Covariance
- If  $x$  and  $y$  are independent...

$$\text{var}(X + Y) = \text{var}(X) + \text{var}(Y)$$

- If  $x$  and  $y$  are dependent...

$$\text{var}(X + Y) = E([(X + Y) - E(X + Y)]^2)$$

$$= E\left(\left[(X + Y) - (E(X) + E(Y))\right]^2\right)$$

$$= E\left(\left[(X - E(X)) + (Y - E(Y))\right]^2\right)$$

$$= E\left(\left(X - E(X)\right)^2 + 2(X - E(X))(Y - E(Y)) + \left(Y - E(Y)\right)^2\right)$$

$$= E\left[\left(X - E(X)\right)^2\right] + E\left[\left(Y - E(Y)\right)^2\right] + 2E\left[\left(X - E(X)\right)\left(Y - E(Y)\right)\right]$$

$$= \text{var}(X) + \text{var}(Y) + 2\text{cov}(X, Y)$$

$E(aX + bY) = aE(X) + bE(Y)$ ,  $a, b \in \mathbb{R}$   
 $X$  and  $Y$  can be independent or dependent.

$$E(XY) = E(X)E(Y)$$

Where  $\text{cov}(X, Y) = 0$

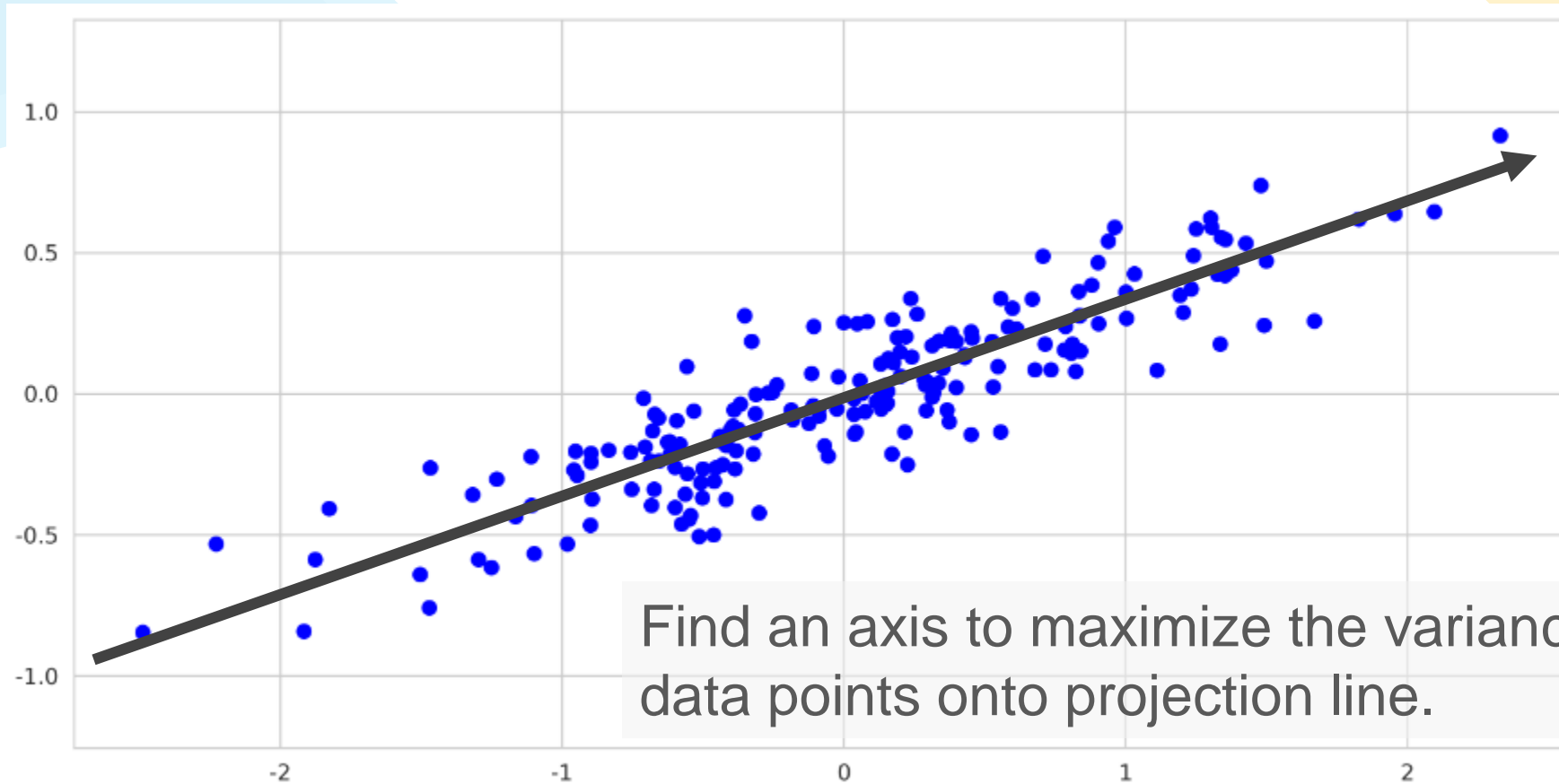


# Why Do We Need Dimension Reduction?

- Here comes the first question into your mind.
  - Why do we need dimension reduction?
  - What's the importance of dimension reduction?
  - Can we directly import all datasets into your model without dimension reduction?
- Statistical models (e.g., linear regression) have several assumptions when you adopt them. One of them is “all variables have to be linearly independent,” indicating no collinearity.
- To achieve this goal, various methods were developed for orthogonalizing parameters and reducing the dimension of the dataset, such as PCA, LDA, LLE, and Laplacian Eigenmaps.

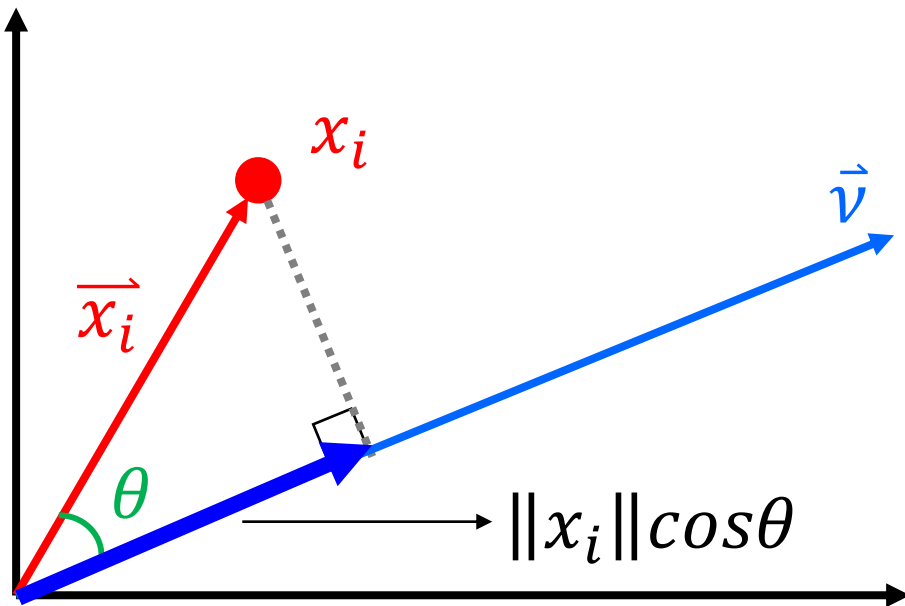
# PCA – Math

$$\mu = \frac{1}{n} \sum_{i=1}^n x_i \quad | \quad \text{var}(x) = \sigma^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2 \quad | \quad \text{Cov}(x, y) = \frac{1}{n} \sum_{i=1}^n (x_i - \mu_x)(y_i - \mu_y)$$



# PCA – Math

- Given a point “ $x$ ” and project onto a vector “ $v$ ”.



$$\cos \theta = \frac{x_i^T \cdot v}{\|x_i\| \|v\|}$$

$$\|x_i\| \cos \theta = \|x_i\| \frac{x_i^T \cdot v}{\|x_i\| \|v\|} = \frac{x_i^T \cdot v}{\|v\|}$$

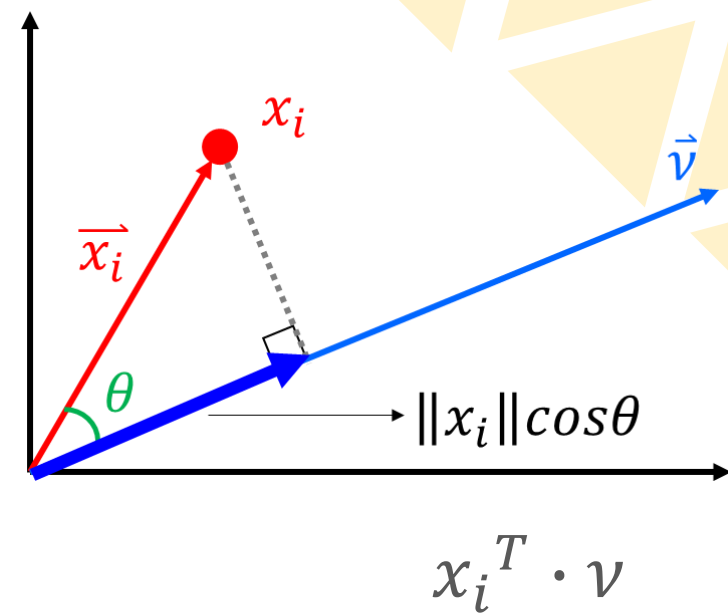
if  $v$  is unit vector ...  $\|v\| = 1$

$$= \frac{x_i^T \cdot v}{\|v\|} = x_i^T \cdot v$$

# PCA – Math

$$X = \begin{bmatrix} | & | & \dots & | \\ x_1 & x_2 & \dots & x_n \\ | & | & \dots & | \end{bmatrix} \rightarrow X^T = \begin{bmatrix} - & x_1 & - \\ - & x_2 & - \\ - & \vdots & - \\ - & x_n & - \end{bmatrix}$$

$$P = \begin{bmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_n^T \end{bmatrix} = X^T u \Rightarrow \text{solve } P$$



# PCA – Math

$$J(u) = \|P^2\| = P^T P = (X^T u)^T (X^T u) = u^T X X^T u$$

$$\operatorname{argmax}_u J(u) = u^T X X^T u, \text{ where subject to } u^T u = 1$$

Add Lagrange multiplier

$$\operatorname{argmax}_{u, \lambda} J(u, \lambda) = u^T X X^T u + \lambda(1 - u^T u)$$

$$\nabla_u J(u, \lambda) = \nabla_u (u^T X X^T u + \lambda(1 - u^T u)) = 0$$

$$\Rightarrow 2X X^T u - 2\lambda u = 0$$

$$\Rightarrow \boxed{X X^T} u = \boxed{\lambda u} \rightarrow \text{eigenvector} \longrightarrow Au = \lambda u$$

$\downarrow$   $\text{cov}(X)$   $\downarrow$  eigenvalue

# PCA – Math

*when  $u$  is the eigen vector*

$$J(u) = \|P^2\| = u^T X X^T u = u^T \lambda u = \lambda u^T u = \lambda$$

***Given an eigenvector, the total square of projected values is the eigenvalue =  $\lambda$***

$$X X^T u = \lambda u$$

Eigenvector is a symmetry matrix  
 $u$  is an unit vector

$$u u^T = u^T u = 1$$

# PCA – Math

A is a square symmetric matrix has orthogonal eigenvectors with different eigenvalues.

$$[x_1, \lambda_1], [x_2, \lambda_2]$$

$$\begin{cases} Ax_1 = \lambda_1 x_1 \\ Ax_2 = \lambda_2 x_2 \end{cases}$$

$$x_1^T Ax_2 = x_1^T \lambda_2 x_2 = \lambda_2 x_1^T x_2 \quad \text{Equal}$$

$$x_1^T A^T x_2 = (Ax_1)^T x_2 = (\lambda_1 x_1)^T x_2 = \lambda_1 x_1^T x_2$$

$$(\because A \in \text{symmetric matrix}, \therefore A = A^T)$$

$$XX^T u = \lambda u$$

Eigenvector is a symmetry matrix  
 $u$  is an unit vector

$$uu^T = u^T u = 1$$

$$\begin{aligned} \lambda_2 x_1^T x_2 &= \lambda_1 x_1^T x_2 \\ \underline{x_1^T x_2} (\lambda_2 - \lambda_1) &= 0 \end{aligned}$$

Orthogonal  $x_1^T x_2 = 0$  All eigenvalues are different

# PCA – Math

- Conversion between orthogonal bases

$$u_i \cdot u_j = u_i^T \cdot u_j = \begin{cases} 1, & \text{if } i = j \\ 0, & \text{otherwise} \end{cases}$$

$$U = \begin{bmatrix} | & | & & | \\ u_1 & u_2 & \dots & u_d \\ | & | & & | \end{bmatrix} \Rightarrow U^T U = I = U^{-1} = U^T$$

$$x = y_1 u_1 + y_2 u_2 + \dots + y_d u_d = \begin{bmatrix} | & | & & | \\ u_1 & u_2 & \dots & u_d \\ | & | & & | \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_d \end{bmatrix} = Uy$$

$$\Rightarrow y = U^{-1}x = U^T x$$



# PCA – Math

$$XX^T v = \lambda v \Rightarrow AV = \lambda V$$

where  $V$  is eigenvector and  $\lambda$  is eigenvalue.

## Principal component $PC$

$$AV = \lambda V$$

## From Singular Vector Decomposition (SVD)

$$A = U\Sigma V^T \Rightarrow AV = U\Sigma \Rightarrow \lambda = \frac{\Sigma^2}{N}, \therefore AV = U\Sigma$$

## Variable loading $L$ (A onto PC)

$$\text{Standardize PC} = \sqrt{N}U$$

$$\text{cov}(A, PC) = \frac{A^T \sqrt{N}U}{N} = \frac{V\Sigma U^T U}{\sqrt{N}} = V \frac{\Sigma}{\sqrt{N}} = V\sqrt{E} = L$$

## PCA – Math

$$J(u) = \|P^2\| = P^T P = (X^T u)^T (X^T u) = u^T X X^T u$$

$$\operatorname{argmax}_u J(u) = u^T X X^T u, \text{ where subject to } u^T u = 1$$

Add Lagrange multiplier

$$\operatorname{argmax}_{u, \lambda} J(u, \lambda) = u^T X X^T u - \lambda(1 - u^T u)$$

$$\nabla_u J(u, \lambda) = \nabla_u (u^T X X^T u - \lambda(1 - u^T u)) = 0$$

$$\Rightarrow 2X X^T u - 2\lambda u = 0$$

$$\Rightarrow \boxed{X X^T u} = \boxed{\lambda u} \rightarrow \text{eigenvector} \longrightarrow Au = \lambda u$$

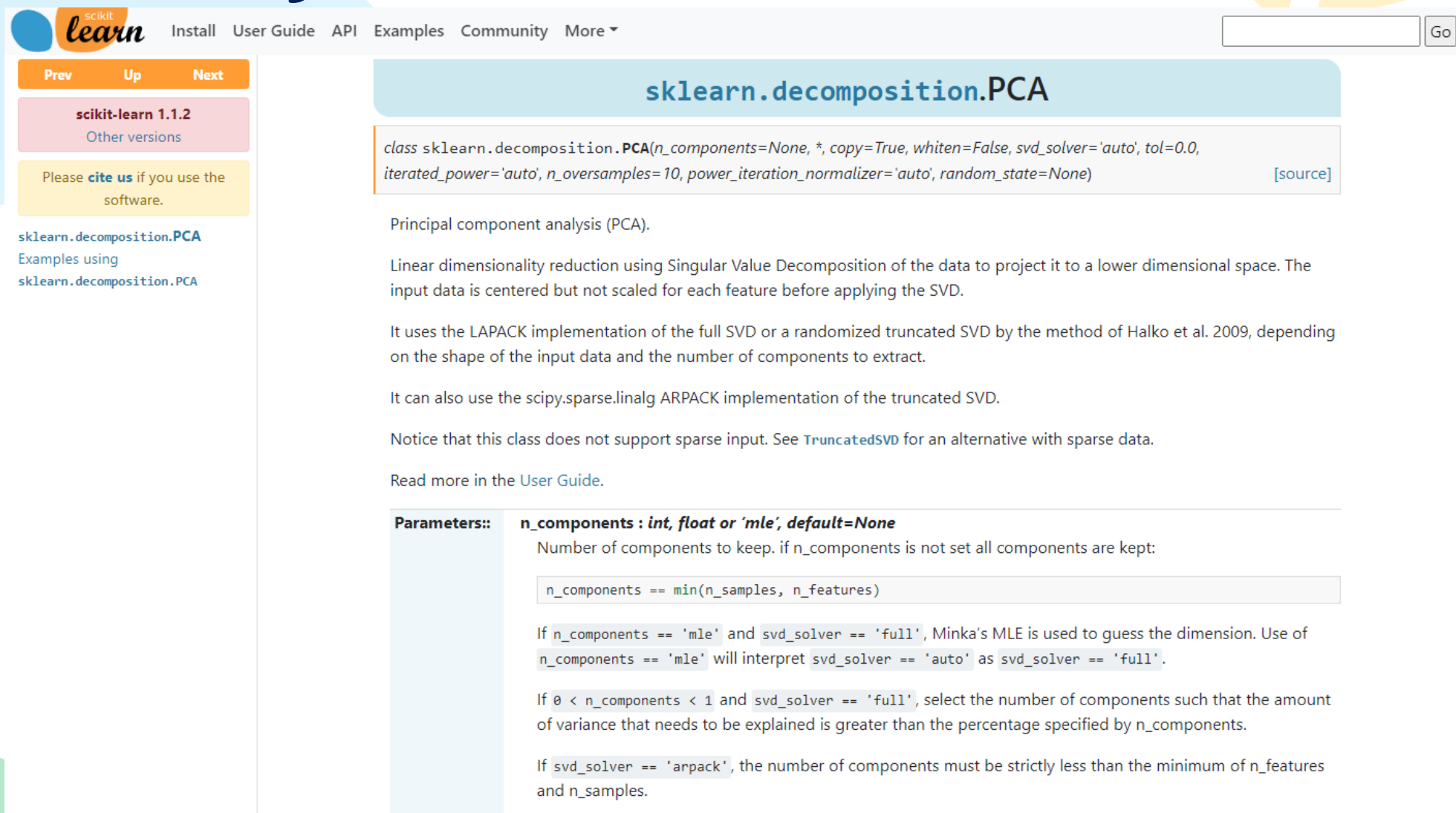
$\downarrow$  cov(X)  $\downarrow$  eigenvalue

# PCA – Math

## The steps of PCA

1. Find the sample mean  $\mu = \frac{1}{n} \sum_{i=1}^n x_i$
2. Subtract mean
3. Compute covariance matrix  $C = \frac{1}{n} X X^T = \frac{1}{n} \sum_{i=1}^n (x_i - \mu) (x_i - \mu)$
4. Find the eigenvalues of C and arrange them into descending order  
 $\lambda_1 > \lambda_2 > \dots > \lambda_d, \{u_1, u_2, \dots, u_d\}$
5. The transformation is  $y = U^T X$ .

# PCA – Python



The image shows a screenshot of the scikit-learn documentation page for the `sklearn.decomposition.PCA` class. The page layout includes a top navigation bar with links for 'Install', 'User Guide', 'API', 'Examples', 'Community', and 'More'. A search bar is located in the top right corner. On the left side, there are navigation buttons for 'Prev', 'Up', and 'Next', along with a section for 'scikit-learn 1.1.2' and a 'cite us' notice. The main content area features a blue header for the class name, followed by the class signature in a code block. Below this, there are several paragraphs of text explaining the PCA method, its implementation, and usage. A 'Parameters:' section is also present, detailing the `n_components` parameter and its various settings.

scikit-learn Install User Guide API Examples Community More  Go

Prev Up Next

scikit-learn 1.1.2  
Other versions

Please [cite us](#) if you use the software.

[sklearn.decomposition.PCA](#)  
Examples using  
[sklearn.decomposition.PCA](#)

## sklearn.decomposition.PCA

```
class sklearn.decomposition.PCA(n_components=None, *, copy=True, whiten=False, svd_solver='auto', tol=0.0, iterated_power='auto', n_oversamples=10, power_iteration_normalizer='auto', random_state=None) [source]
```

Principal component analysis (PCA).

Linear dimensionality reduction using Singular Value Decomposition of the data to project it to a lower dimensional space. The input data is centered but not scaled for each feature before applying the SVD.

It uses the LAPACK implementation of the full SVD or a randomized truncated SVD by the method of Halko et al. 2009, depending on the shape of the input data and the number of components to extract.

It can also use the `scipy.sparse.linalg` ARPACK implementation of the truncated SVD.

Notice that this class does not support sparse input. See [TruncatedSVD](#) for an alternative with sparse data.

Read more in the [User Guide](#).

**Parameters:**

**n\_components** : *int, float or 'mle', default=None*  
Number of components to keep. if `n_components` is not set all components are kept:

```
n_components == min(n_samples, n_features)
```

If `n_components == 'mle'` and `svd_solver == 'full'`, Minka's MLE is used to guess the dimension. Use of `n_components == 'mle'` will interpret `svd_solver == 'auto'` as `svd_solver == 'full'`.

If  $0 < n\_components < 1$  and `svd_solver == 'full'`, select the number of components such that the amount of variance that needs to be explained is greater than the percentage specified by `n_components`.

If `svd_solver == 'arpark'`, the number of components must be strictly less than the minimum of `n_features` and `n_samples`.

# PCA – Python

```
# linear algebra
```

```
from sklearn.decomposition import PCA
```

```
# call PCA
```

```
pca = PCA(n_components=2) # number of preserved components
```

```
pca.fit(X)
```

```
# show results
```

```
print(pca.mean_) # mean
```

```
print(pca.explained_variance_) # eigenvalues
```

```
print(pca.components_) # eigenvectors
```

# Notice

- Advantages of PCA
  - Easy to use
  - Accelerate model fitting
  - Avoid overfitting
- Disadvantage of PCA
  - Low interpretability
  - The trade-off between information loss and features (dimensions)

- Limitations of PCA
  - Linear relationship between features
  - Correlation between features
  - Sensitive to the scale of data
  - Not robust to outliers
  - Not accept missing values

# Question Time

If you have any questions, please do not hesitate to ask me.

## References:

張智星 科學計算講義

Sklearn

# The End

*Thank you for your attention ))*