

基礎機器學習與深度學習- 2

黃志勝

義隆電子人工智慧研發部

國立陽明交通大學AI學院合聘助理教授



基礎機器學習

針對前述的介紹，每個**topic**都介紹一個演算算法

1. Regression: Linear regression & Regularization
2. Classification: Linear and Quadratic Discriminant Analysis
3. Clustering: K-means (Unsupervised learning)
4. Dimension Reduction: PCA (Unsupervised learning)
5. Ensemble learning: 不介紹。



k -means Clustering(Unsupervised learning)

- k -means Clustering: 物以類聚(類似歸納法)

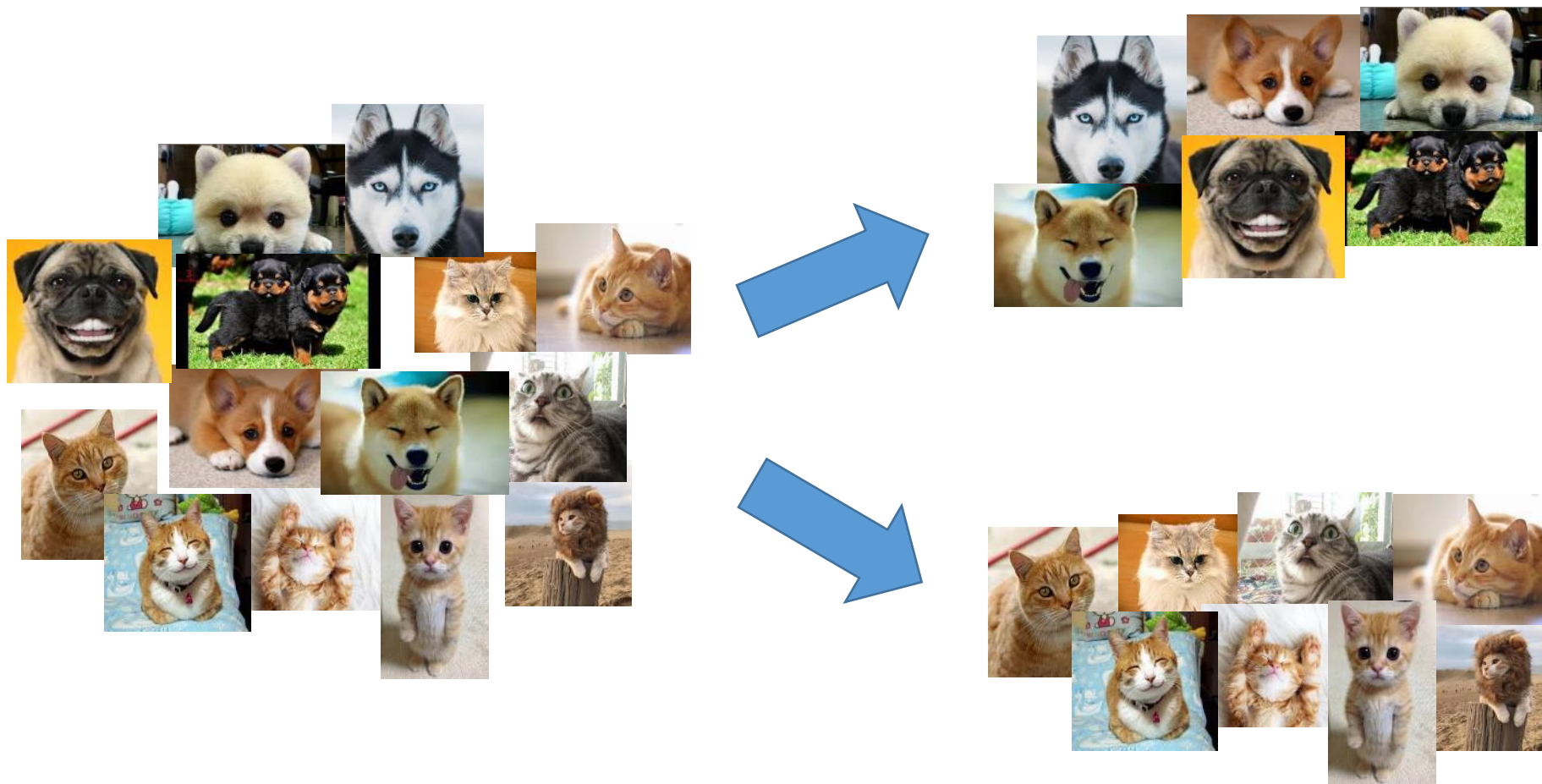
不斷學習(iteration)，直到收斂為止。

為什麼叫 k -means顧名思義就是有 k 個群心，我們將資料學習後判斷這些資料屬於哪個群心。

EX: 給你一組身高和體重資料，但我沒有跟你說這組資料哪些是男生哪些是女生。我希望你用這組資料分出兩群，這種時候就是用非監督式學習。→ 2-means clustering



k-means Clustering(Unsupervised learning)

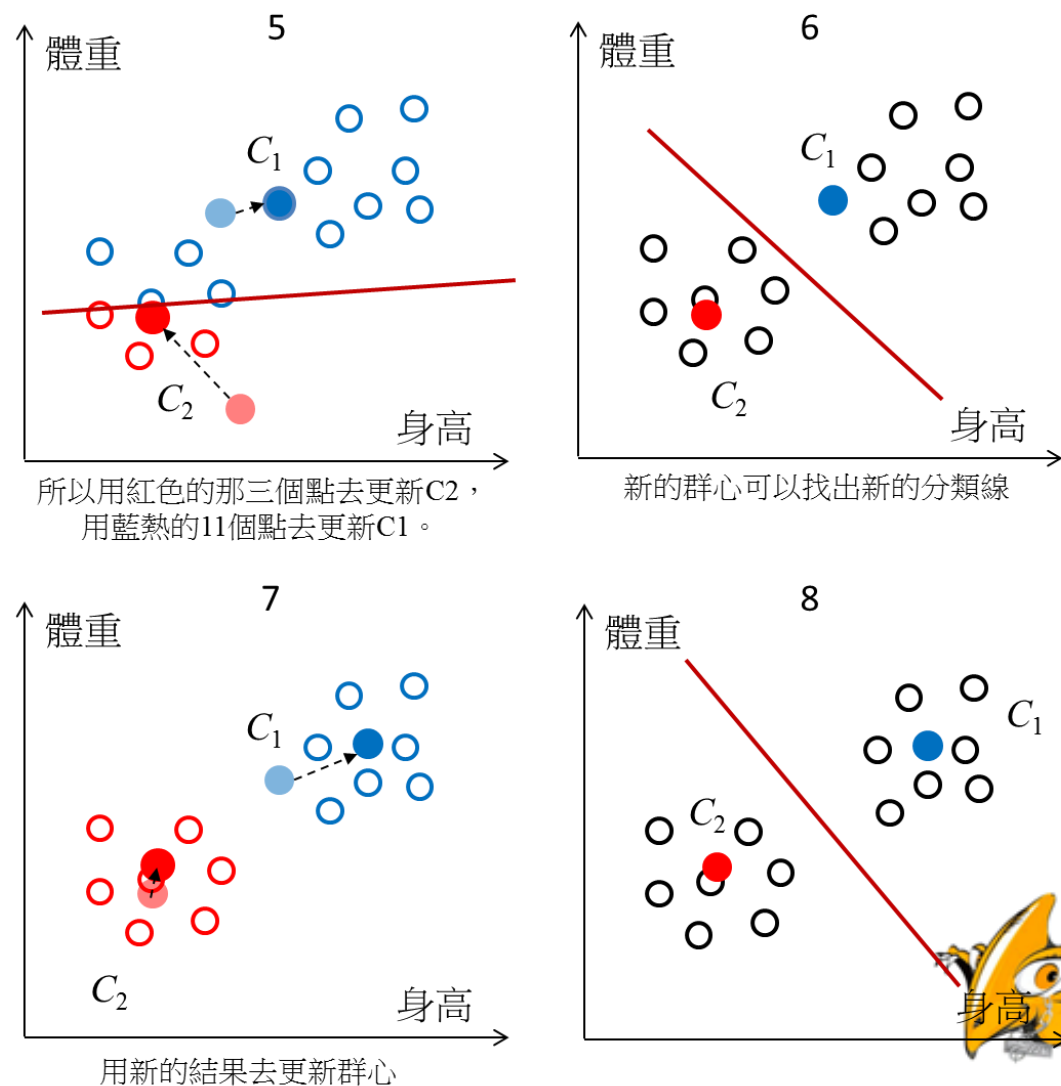
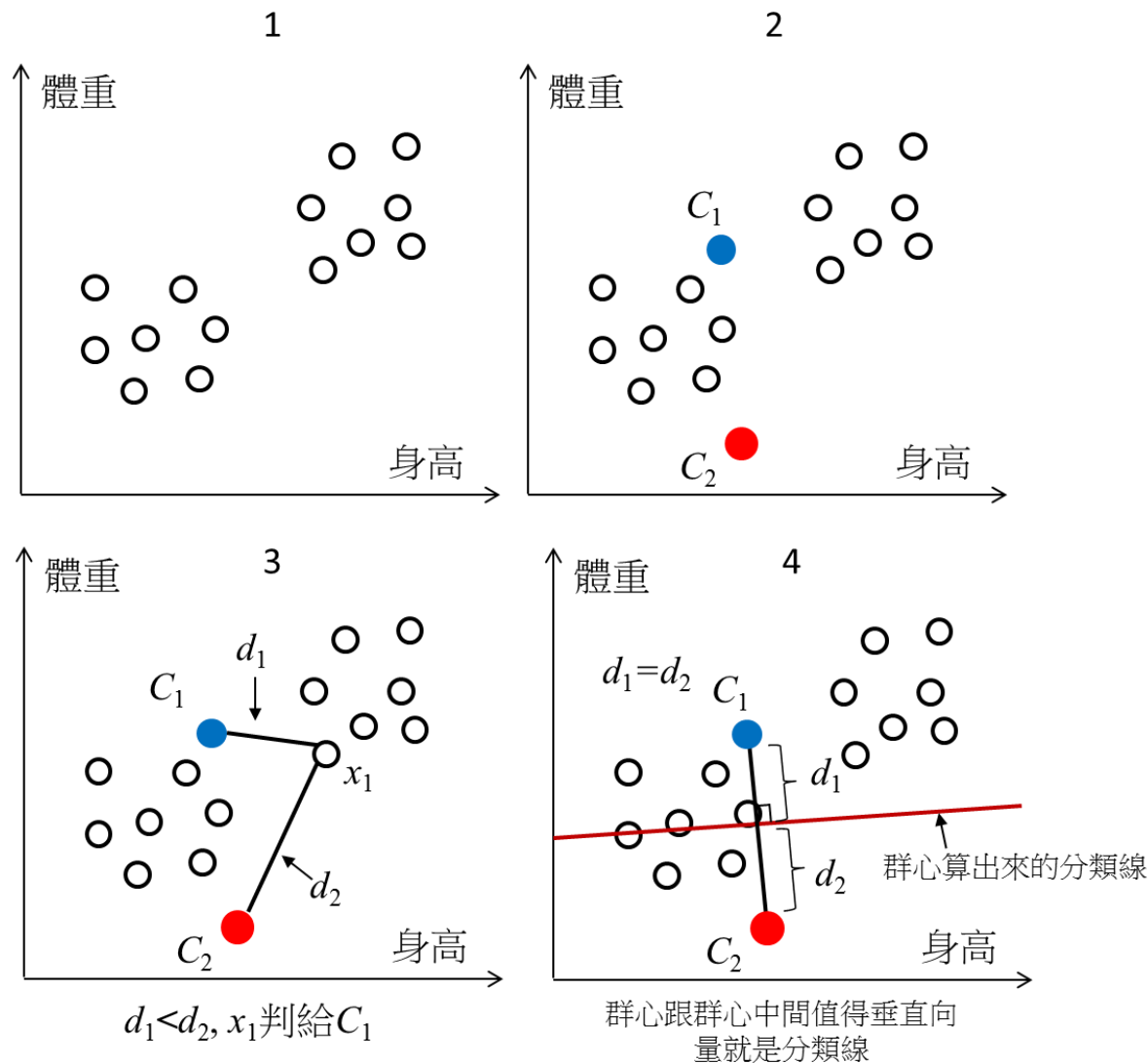


k-means Clustering(Unsupervised learning)

- *k*-means: with a name “means”.
- Most important information is “mean vector”.
- *k*-means is based on learning the mean vector for each cluster.
- Number of cluster is setting by user.



k -means Clustering(Unsupervised learning)



k -means Clustering(Unsupervised learning)

- 設定有 k (必須 $\leq n$)個Clusters $\{S_1, S_2, \dots, S_k\}$ ， k -means clustering就是希望可以最小化群內的資料和群心的誤差平方和越小越好，數學公式如下：

$$\arg \min_{\mu} \sum_{c=1}^K \sum_{i=1}^{n_c} \|x_i - \mu_c\|^2 \Big|_{x_i \in S_c}$$



k -means Clustering

1. 初始隨機設定 k 個群心.

$$\mu_c^{(0)} \in R^d, c = 1, 2, \dots, K.$$

2. 計算分類到每一群體的樣本， (t) 為第 t 次運算

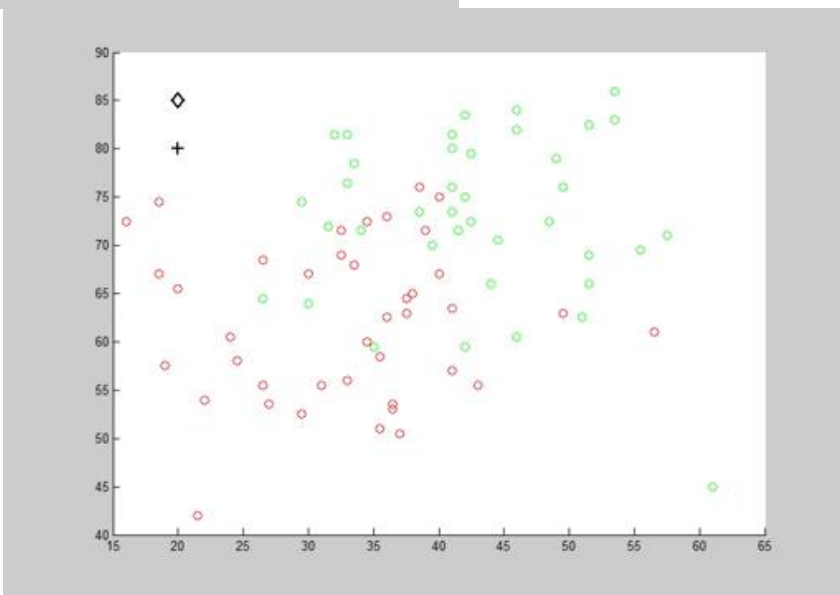
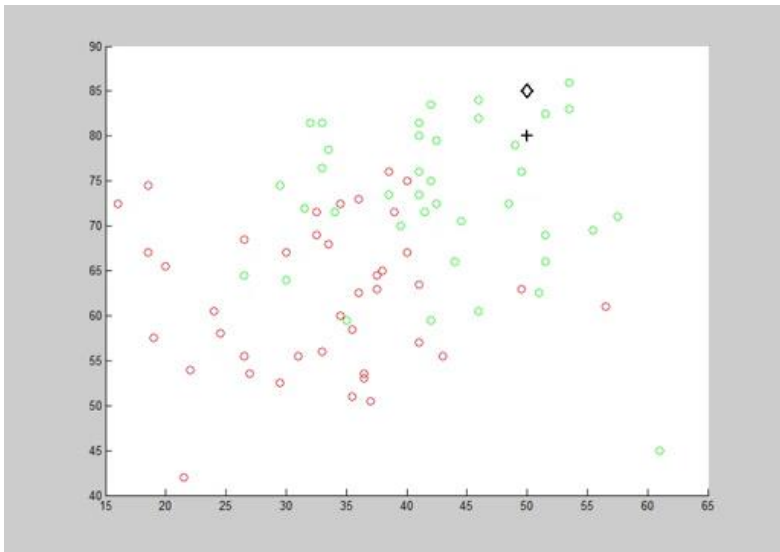
$$S_c^{(t)} = \{x_i: \|x_i - \mu_c^{(t)}\| \leq \|x_i - \mu_{c^*}^{(t)}\|, \forall i = 1, \dots, n\}.$$

3. 更新群心(n_c 個資料在第 c 群內。)

$$\mu_c^{(t+1)} = \frac{\text{sum}(S_c^{(t)})}{n_c} = \sum_{i=1}^{n_c} x_i \Big|_{x_i \in S_c^{(t)}}$$

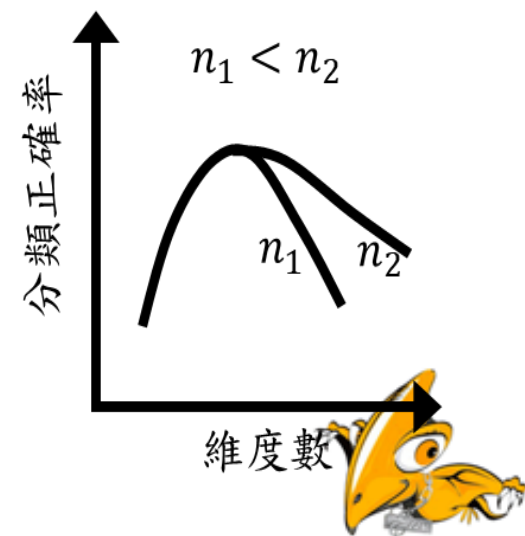
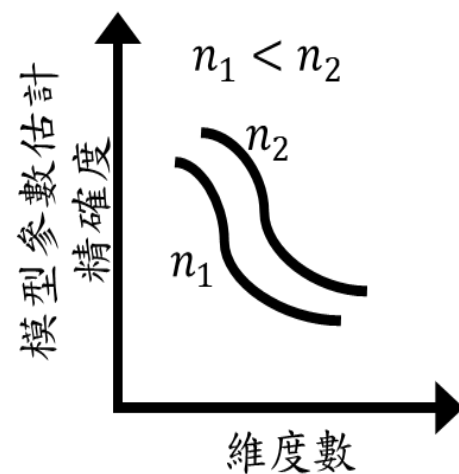
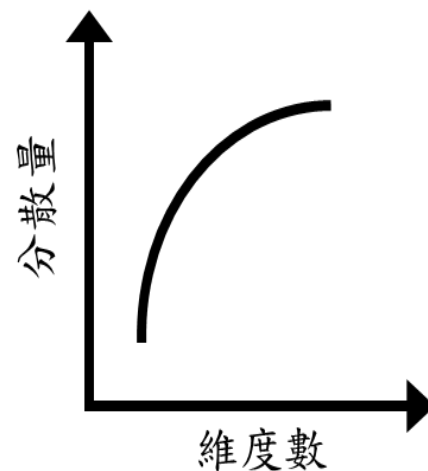
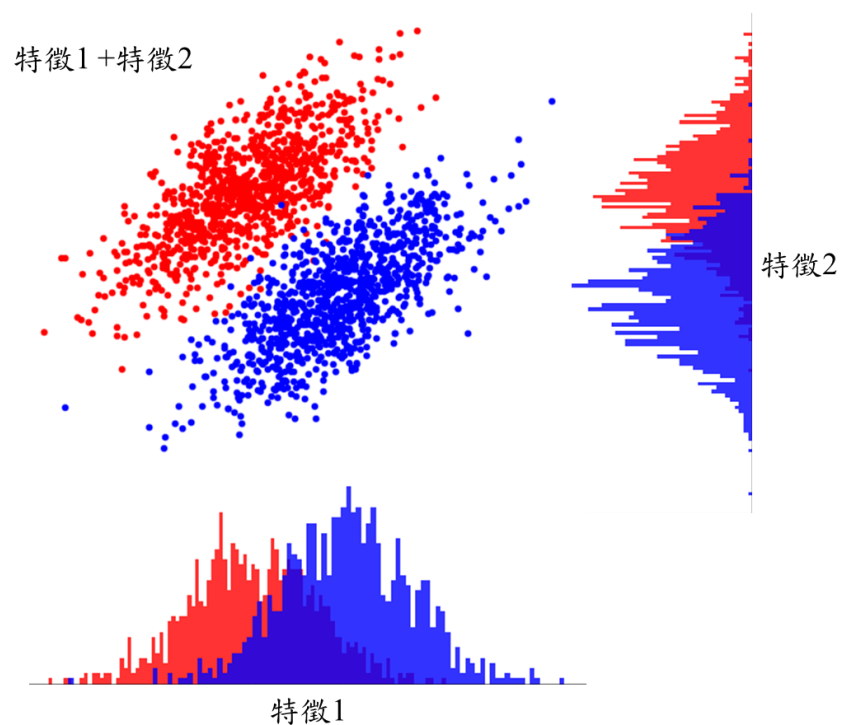
4. 重複2-3，直到群心不變動，也就是

$$S_c^{(t+1)} = S_c^{(t)}, \forall c = 1, \dots, K$$



Dimension Reduction

在建立預測模型時，容易因為特徵數大於資料樣本數造成模型參數估計錯誤，導致模型無法有效進行任務預測，在機器學習稱此現象為「休斯現象(Hughes phenomenon)」，也稱為「維度詛咒(Curse of dimensionality)」，



Dimension Reduction

Example:

Model 1: “body fat (bf)”

Model 2: “body fat (bf)”, “weight (w)”, “hair length (hl)”

$$\text{cov}(\text{model1}) = [\text{cov}(bf, bf)]$$

$$\text{cov}(\text{model2}) = \begin{bmatrix} \text{cov}(bf, bf) & \text{cov}(w, bf) & \text{cov}(hl, bf) \\ \text{cov}(bf, w) & \text{cov}(w, w) & \text{cov}(hl, w) \\ \text{cov}(bf, hl) & \text{cov}(w, hl) & \text{cov}(hl, hl) \end{bmatrix}$$



Example for single variable

If we only get one sample, and try to calculate covariance.

$$\mu = x_i$$

$$\text{cov}(\text{model 1}) = \sigma = \frac{1}{1} \sum_{i=1}^1 (x_i - \mu_x)^2 = 0$$

$$f(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



Example for multi-variables

If we only get two sample, and try to calculate covariance matrix.

$$\Sigma = \begin{bmatrix} \text{cov}(bf, bf) & \text{cov}(w, bf) & \text{cov}(hl, bf) \\ \text{cov}(bf, w) & \text{cov}(w, w) & \text{cov}(hl, w) \\ \text{cov}(bf, hl) & \text{cov}(w, hl) & \text{cov}(hl, hl) \end{bmatrix}$$

The elements in covariance matrix are larger than 0, but the covariance matrix would be singular.

$$f(\mathbf{x}|\boldsymbol{\mu}, \Sigma) = (2\pi)^{-d/2} |\Sigma|^{-0.5} \exp\{-0.5(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu})\}$$



Dimension Reduction

- Dimension Reduction is proposed to overcome this issue.
 1. Feature selection
Using only “import” features.
 2. Feature extraction
Feature Fusion.



Feature Selection

In statistics,

- Forward sequential feature selection
- Backward sequential feature selection
- Stepwise feature selection
- LASSO

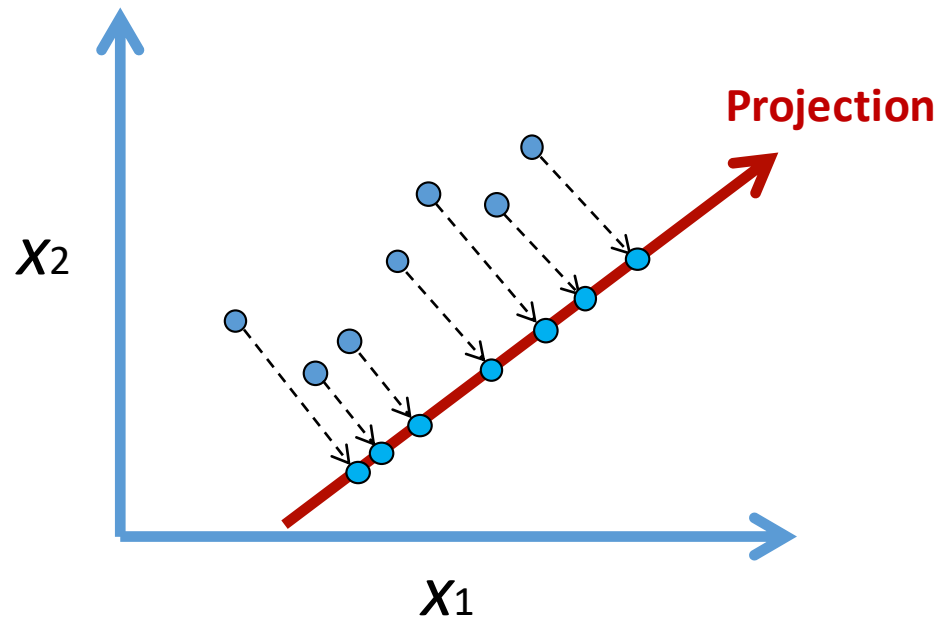
In machine learning,

- Random subspace.



Feature Extraction

- Feature Fusion (Projection)



Feature extraction:
Just finding the projection vectors for input features.



Feature Extraction

- **Principal component analysis (PCA)**
- Independent component analysis (ICA)
- Canonical component analysis (CCA)
- Non-negative matrix factorization
- Discriminant Analysis Feature Extraction (DAFE)
- Neural Network



Principal component analysis (PCA)

Why do I introduce PCA?

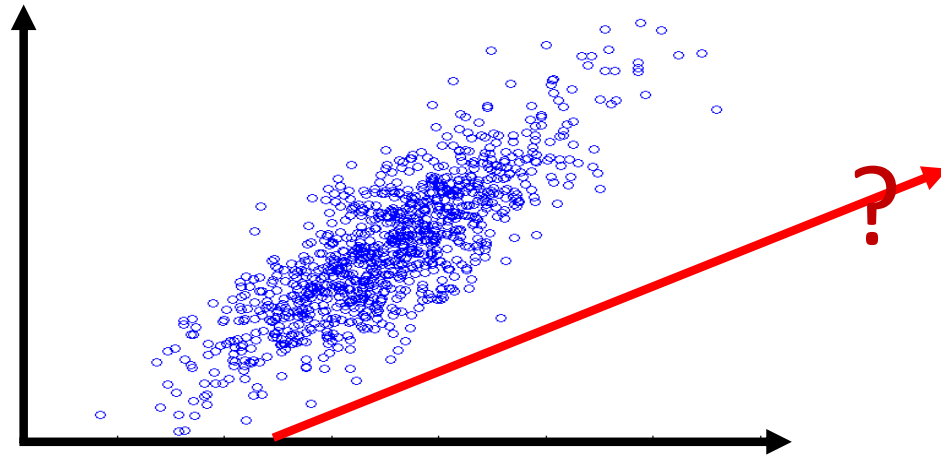
1. Stronger knowledge.
2. Unsupervised.
3. Most popular
4. Basic



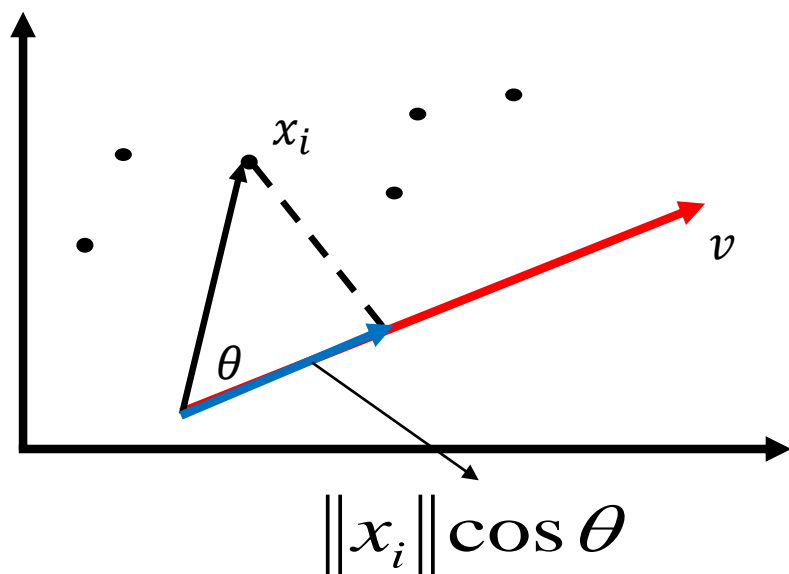
Principle Component Analysis

DL/ML/Statistics are developed by a given goal.

- PCA aims to find a set of vector containing the maximum amount of variance in the data.



Principle Component Analysis



$$\cos \theta = \frac{\langle x_i, v \rangle}{\|x_i\| \|v\|}$$

$$\begin{aligned} & \|x_i\| \cos \theta \frac{v}{\|v\|} \\ &= \|x_i\| \frac{\langle x_i, v \rangle}{\|x_i\| \|v\|} \frac{v}{\|v\|} = \frac{\langle x_i, v \rangle}{\|v\|^2} v \end{aligned}$$

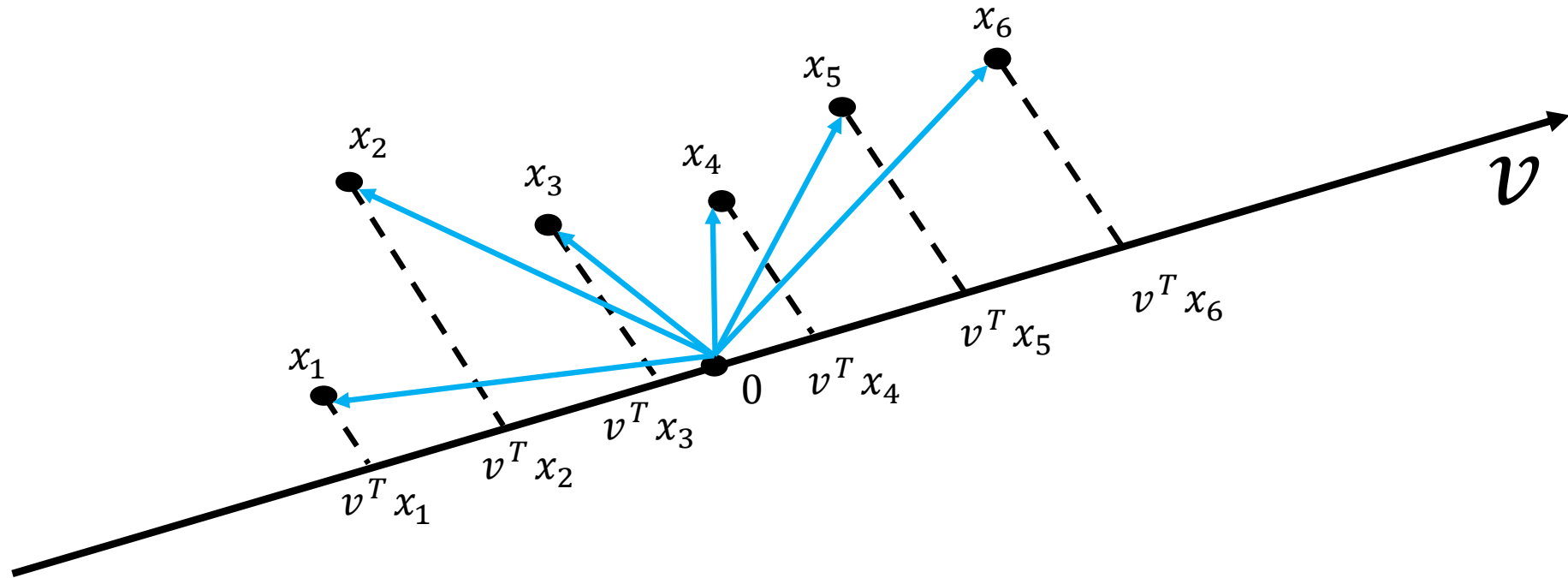
- If $\|v\|$ is unit, then $\langle x_i, v \rangle v$

$$\langle x_i, v \rangle = x_i^T v = v^T x_i$$

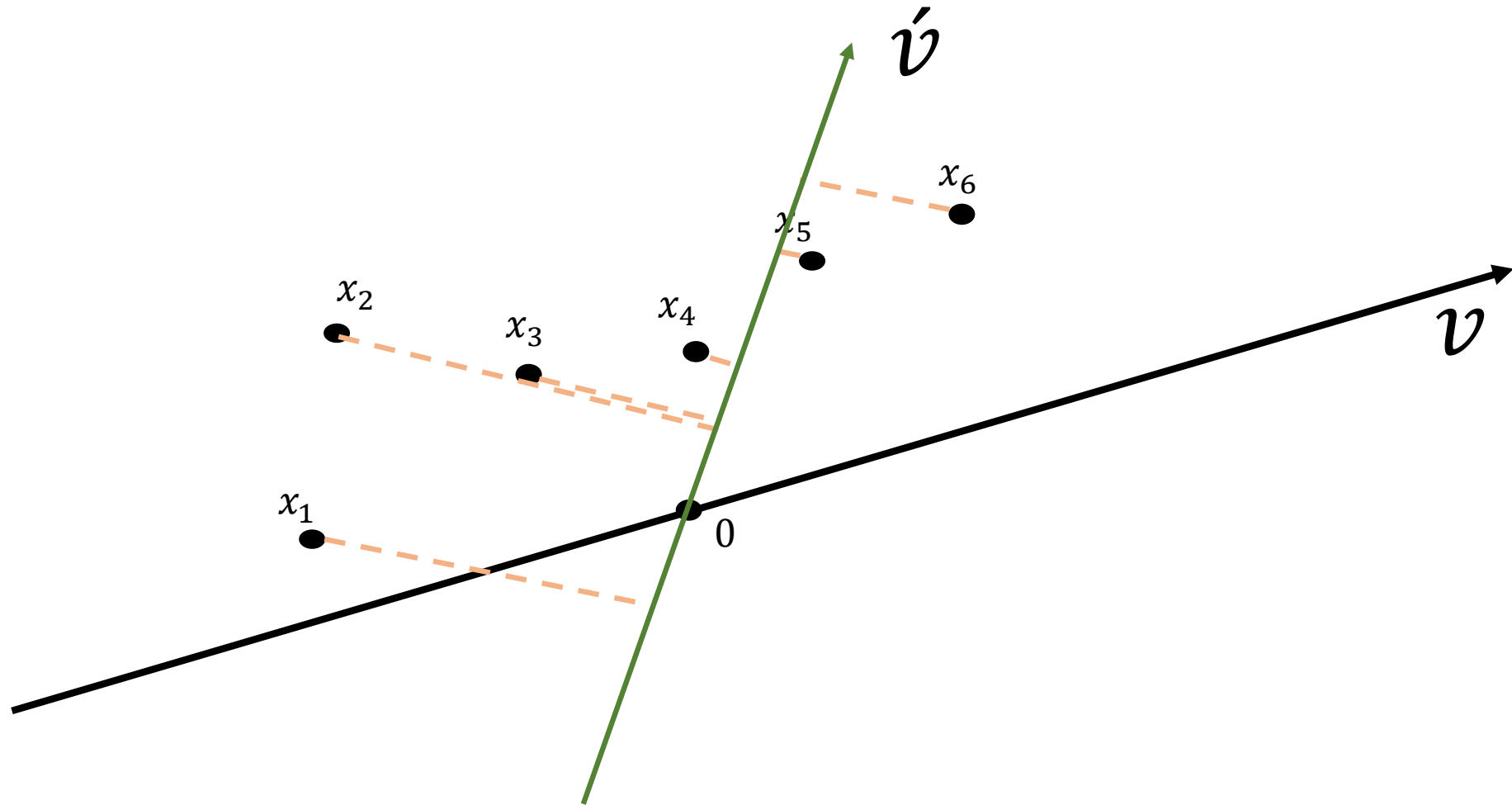
$$y_i = v^T x_i$$



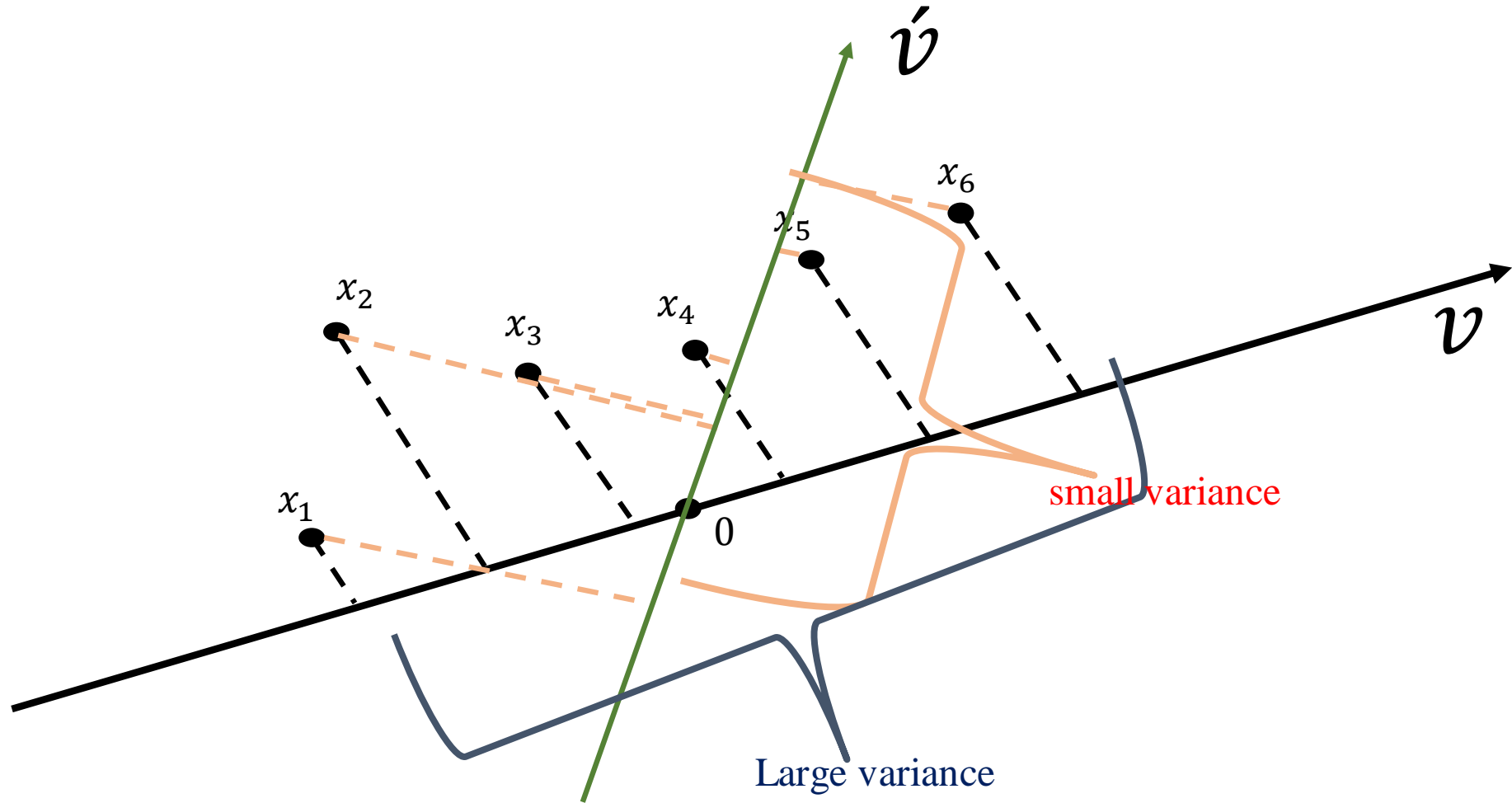
Principle Component Analysis



Principle Component Analysis



Principle Component Analysis



Principle Component Analysis

- The projections of the all points x_i into the direction v are
$$v^T x_1, v^T x_2, \dots, v^T x_N$$

The variance of the projections is

$$\sigma^2 = \frac{1}{N} \sum_{i=1}^N (v^T x_i - \mu)^2 = \frac{1}{N} \sum_{i=1}^N (v^T x_i - 0)^2 = \frac{1}{N} \sum_{i=1}^N (v^T x_i)^2$$

$$\begin{aligned} \Sigma &= \frac{1}{N} \sum_{i=1}^N (v^T x_i)(v^T x_i)^T = \frac{1}{N} \sum_{i=1}^N (v^T x_i x_i^T v) = v^T \left(\frac{1}{N} \sum_{i=1}^N x_i x_i^T \right) v \\ &= v^T C v \end{aligned}$$

C covariance matrix



Principle Component Analysis

- The first principal vector can be found by the following equation:

$$\mathbf{v} = \arg \max_{\mathbf{v} \in \mathbb{R}^d, \|\mathbf{v}\|=1} \mathbf{v}^T \mathbf{C} \mathbf{v}$$



Principle Component Analysis

$$\mathbf{v} = \arg \max_{\mathbf{v} \in \mathbb{R}^d, \|\mathbf{v}\|=1} \mathbf{v}^T \mathbf{C} \mathbf{v}$$

Lagrange function:

$$f(\mathbf{v}, \lambda) = \mathbf{v}^T \mathbf{C} \mathbf{v} - \lambda(\mathbf{v}^T \mathbf{v} - 1)$$

$$\frac{\partial f(\mathbf{v}, \lambda)}{\partial \mathbf{v}} = 0 \Rightarrow 2\mathbf{C} \mathbf{v} - \lambda \mathbf{v} = 0 \Rightarrow \mathbf{C} \mathbf{v} = \lambda \mathbf{v}$$

$$\frac{\partial f(\mathbf{v}, \lambda)}{\partial \lambda} = 0 \Rightarrow \mathbf{v}^T \mathbf{v} - 1 = 0 \Rightarrow \mathbf{v}^T \mathbf{v} = 1$$



Principle Component Analysis

- The first principal vector can be found by the following equation:

$$v = \underset{v \in R^d, \|v\|=1}{\operatorname{argmax}}(v^T C v)$$

- This is equivalent to find the largest eigenvalue of the following eigenvalue problem:

$$Cv = \lambda v$$
$$\|v\| = 1$$

$$C = \frac{1}{N} \sum_{i=1}^N x_i x_i^T = \frac{1}{N} [x_1 \dots x_N] \begin{bmatrix} x_1 \\ \vdots \\ x_N \end{bmatrix} = \frac{1}{N} X^T X$$



Principle Component Analysis

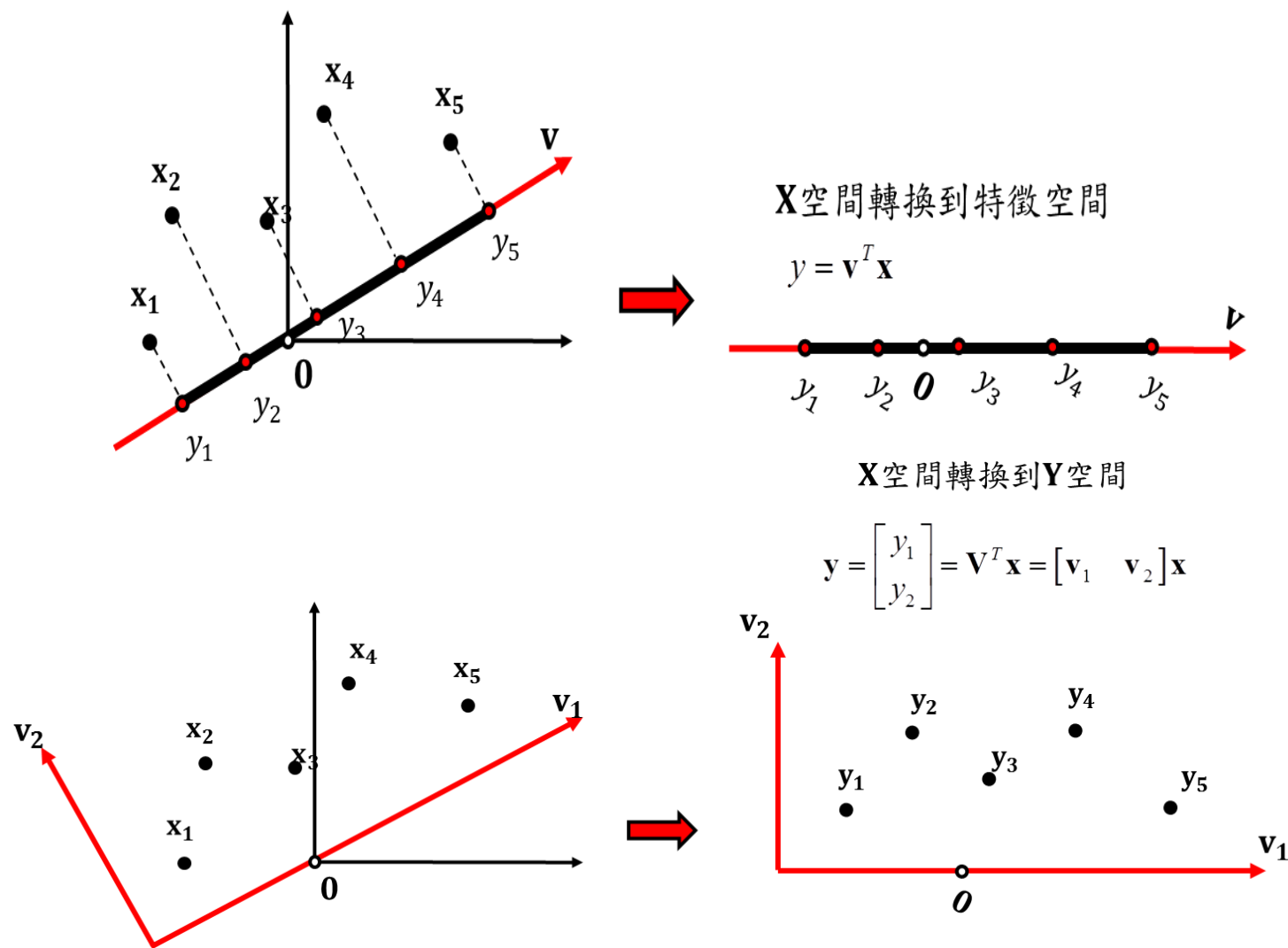
Eigenvalue vector is the corresponding variance vector.

$$\mathbf{v}^T \mathbf{C} \mathbf{v} = \sigma^2$$

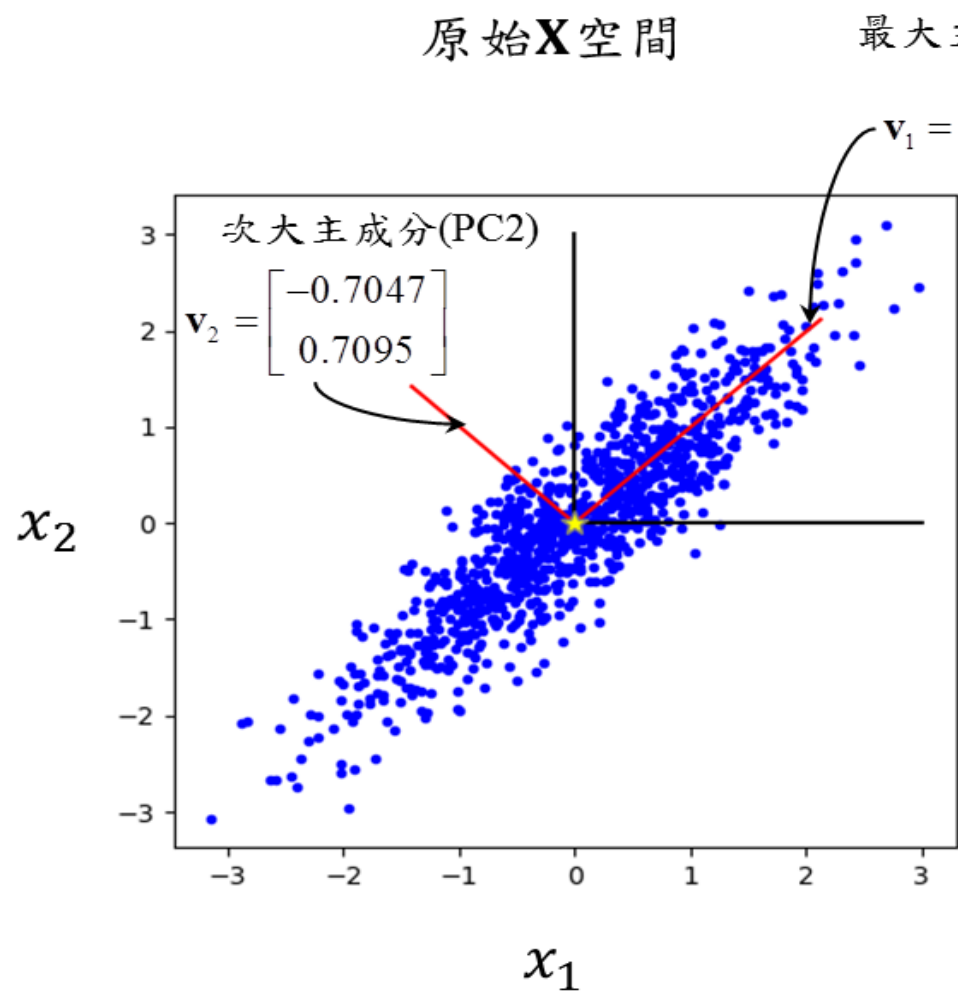
$$\Rightarrow \mathbf{v}^T \lambda \mathbf{v} = \lambda \mathbf{v}^T \mathbf{v} = \lambda = \sigma^2$$



Projection



Exercise



PCA

