# 国際融合科学論/先端融合科学論

## LECTURE 02
## Machine Learning I: conventional machine learning

Dr. Suyong Eum

OSAKA UNIVERSITY

1) Principal Component Analysis (PCA)

- Feature selections

- Dimension reduction

2) Support Vector Machine (SVM)
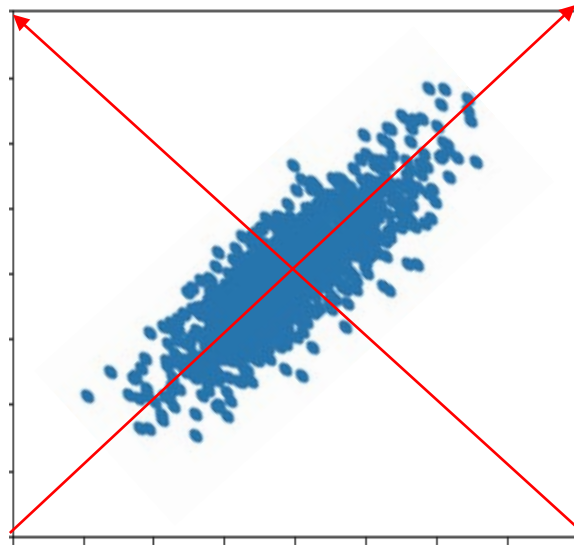
- Hard margin SVM: linear classification

- Kernel trick: nonlinear classification
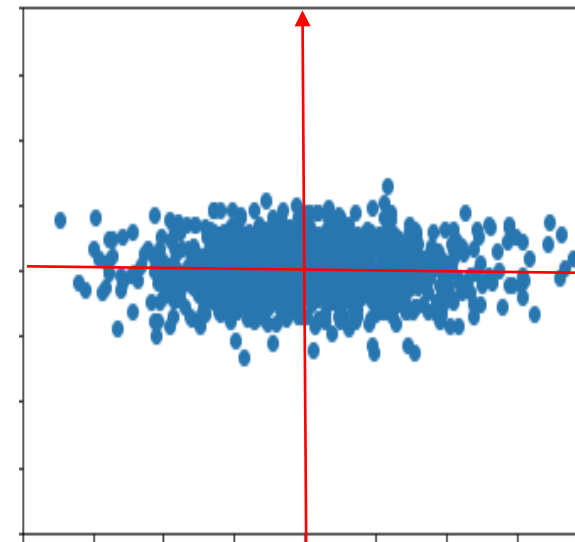
# Principal Component Analysis (PCA)

A statistical procedure that uses an orthogonal transformation to convert a set of observations of possibly correlated variables into a set of values of linearly uncorrelated variables called principal components.
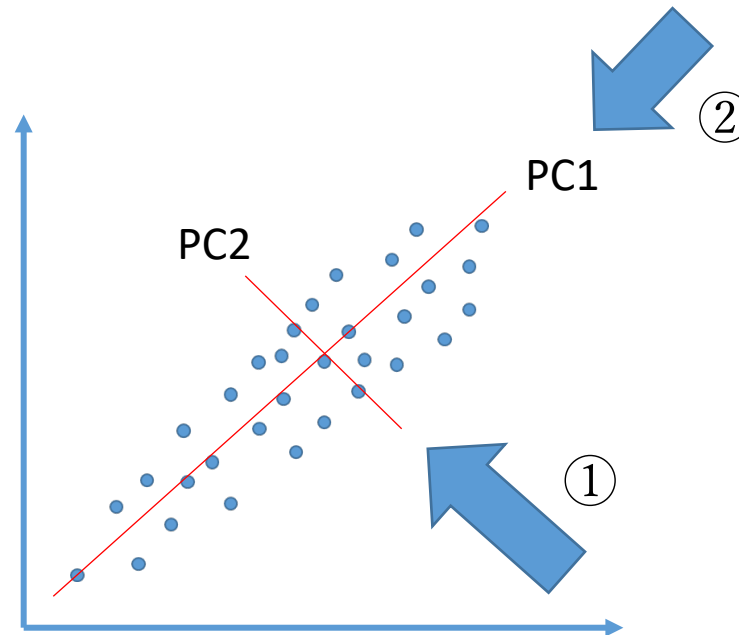
*In Wikipedia:*

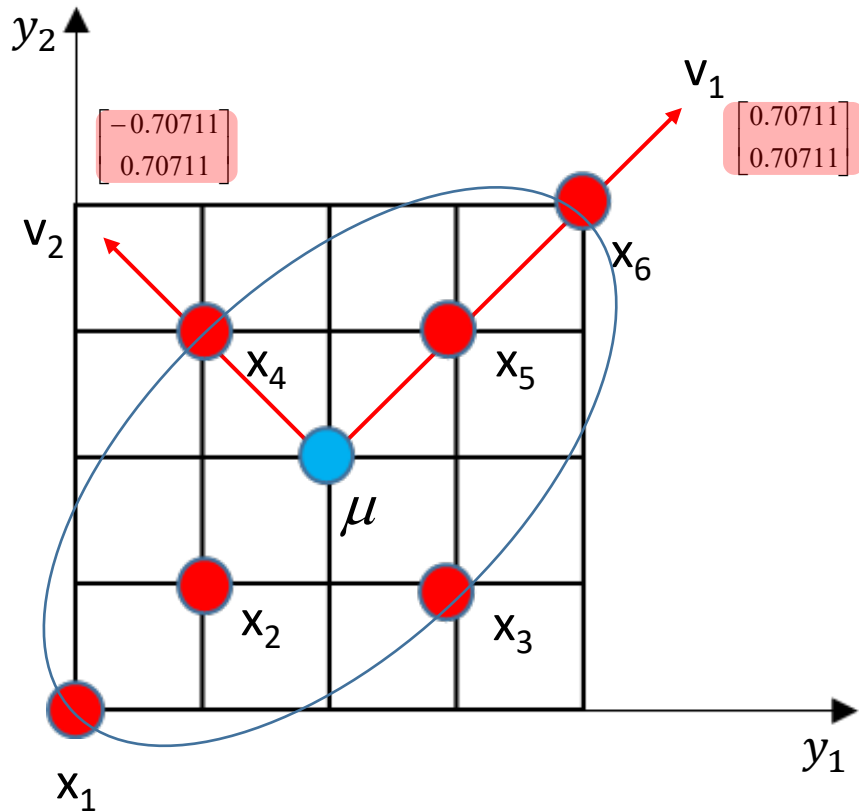$$\begin{bmatrix} 8 & 1 \\ 1 & 8 \end{bmatrix} \qquad \begin{bmatrix} 10 & 0 \\ 0 & 1 \end{bmatrix}$$

❑ How to select principal components?

- One that captures the largest variance of the data points

- Intuitively speaking, you can observe more data from the direction ① than any other direction, and then from the direction ②, you can observe the data with the least redundancy compared to the direction ① .

PC1

PC2

②

①

☐ Distance to data points from the mean along the axis of **"$v_1$"**

$$= [2\sqrt{2}, \ \sqrt{2}, \ 0, \ 0, \ \sqrt{2}, \ 2\sqrt{2}]$$

**Variance = 4**

☐ Distance to data points from the mean along the axis of **"$v_2$"**

$$= [0, \ 0, \ \sqrt{2}, \ \sqrt{2}, 0,0]$$

**Variance = 0.8**

☐ $\text{cov}(X) = \begin{bmatrix} 2.4 & 1.6 \\ 1.6 & 2.4 \end{bmatrix}$

$$X = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \\ x_6 \end{bmatrix} = \begin{bmatrix} -2 & -2 \\ -1 & -1 \\ 1 & -1 \\ -1 & 1 \\ 1 & 1 \\ 2 & 2 \end{bmatrix}$$
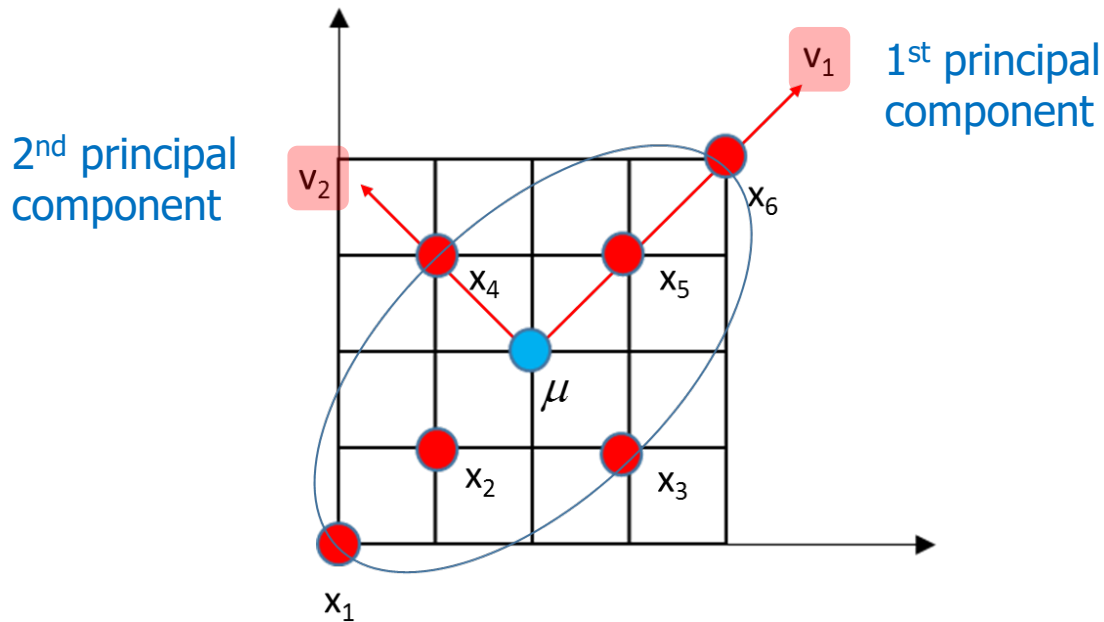
☐ $\text{cov}(x) = V\Lambda V^T$

```
>> [vec, val] = eig(cov(x))
vec =

   -0.70711    0.70711
    0.70711    0.70711
```
$V$

```
val =

Diagonal Matrix

    0.80000          0
          0    4.00000
```
$\Lambda$

1) Find the covariance matrix of data points.

2) Obtain the eigen values and vectors of the covariance matrix: eigen value decomposition.

3) Sort the eigen vectors in descending order in terms of their corresponding eigen values.

   - an eigen vector with the largest eigen value becomes the first principal component.

2nd principal component

1st principal component

$V_1$

$V_2$

$x_6$

$x_4$

$x_5$

$\mu$

$x_2$

$x_3$

$x_1$

```
>> x
x =

    -2  -2
    -1  -1
     1  -1
    -1   1
     1   1
     2   2

>> cov(x)
ans =

    2.4000   1.6000
    1.6000   2.4000
```

```
>> [vec, val] = eig(cov(x))
vec =
```

2nd principal component

1st principal component

```
    -0.70711   0.70711
     0.70711   0.70711

val =

Diagonal Matrix

     0.80000        0
          0   4.00000
```

❑  Actually, there is a more convenient way of doing it, which is called "Singular Value Decomposition" or SVD.

Eigen decomposition

$$X^T X = \boxed{V \Lambda V^T}$$

Singular Value Decomposition (SVD)

$$\boxed{X = U\Sigma V^T}$$

```
>> x
x =

  -2  -2
  -1  -1
   1  -1
  -1   1
   1   1
   2   2

>> cov(x)
ans =

  2.4000  1.6000
  1.6000  2.4000
```

```
>> [vec, val] = eig(cov(x))
vec =

  -0.70711   0.70711
   0.70711   0.70711

val =

Diagonal Matrix

   0.80000        0
        0   4.00000
```

```
>> [vec, val]=eig(transpose(x)*x)
vec =

  -0.70711   0.70711
   0.70711   0.70711

val =

Diagonal Matrix

   4.0000         0
        0   20.0000
```

$$X^T X = (U\Sigma V^T)^T (U\Sigma V^T)$$

$$= V\Sigma^T U^T U\Sigma V^T$$

$$= \boxed{V\Sigma^2 V^T} \qquad \Lambda = \Sigma^2$$

$$\Lambda = \begin{bmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_n \end{bmatrix} \qquad \Sigma = \begin{bmatrix} \sqrt{\lambda_1} & & 0 \\ & \ddots & \\ 0 & & \sqrt{\lambda_n} \end{bmatrix}$$

Eigen value                    Singular value

❑ Actually, there is a more convenient way of doing it, which is called "Singular Value Decomposition" or SVD.

Eigen decomposition

$$X^T X = V \Lambda V^T$$

Singular Value Decomposition (SVD)

$$X = U \Sigma V^T$$

```
>> x
X =

   -2   -2
   -1   -1
    1   -1
   -1    1
    1    1
    2    2

>> cov(x)
ans =

   2.4000   1.6000
   1.6000   2.4000
```

```
>> [vec, val] = eig(cov(x))
vec =

   -0.70711   0.70711
    0.70711   0.70711

val =

Diagonal Matrix

   0.80000        0
        0   4.00000
```

```
>> [vec, val]=eig(transpose(x)*x)
vec =

   -0.70711   0.70711
    0.70711   0.70711

val =

Diagonal Matrix

   4.0000        0
        0   20.0000
```

$$\Lambda = \Sigma^2$$

$$4.4721^2 = 20$$
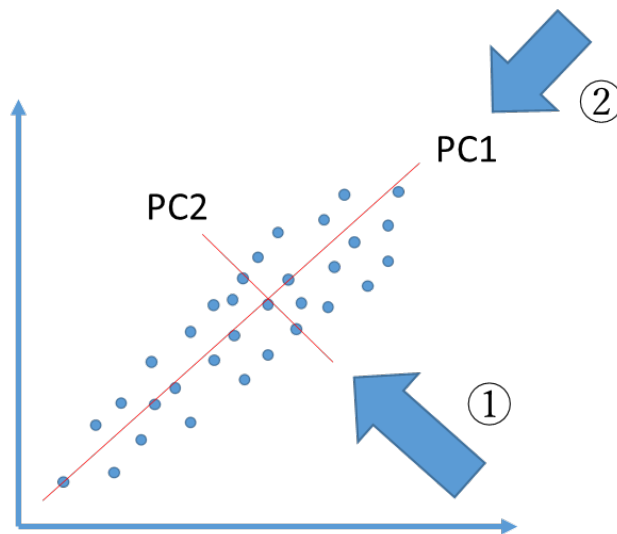
```
>> [U,S,V]=svd(x)
U =

   -0.63246   0.00000   0.30819  -0.30819   0.28637   0.57274
   -0.31623  -0.00000  -0.63635   0.63635   0.13426   0.26851
    0.00000  -0.70711   0.50000   0.50000   0.00000   0.00000
   -0.00000   0.70711   0.50000   0.50000  -0.00000  -0.00000
    0.31623   0.00000  -0.00399   0.00399   0.94140  -0.11720
    0.63246   0.00000  -0.00799   0.00799  -0.11720   0.76560

S =

Diagonal Matrix

   4.4721        0
        0   2.0000
        0        0
        0        0
        0        0
        0        0

V =

   0.70711  -0.70711
   0.70711   0.70711
```
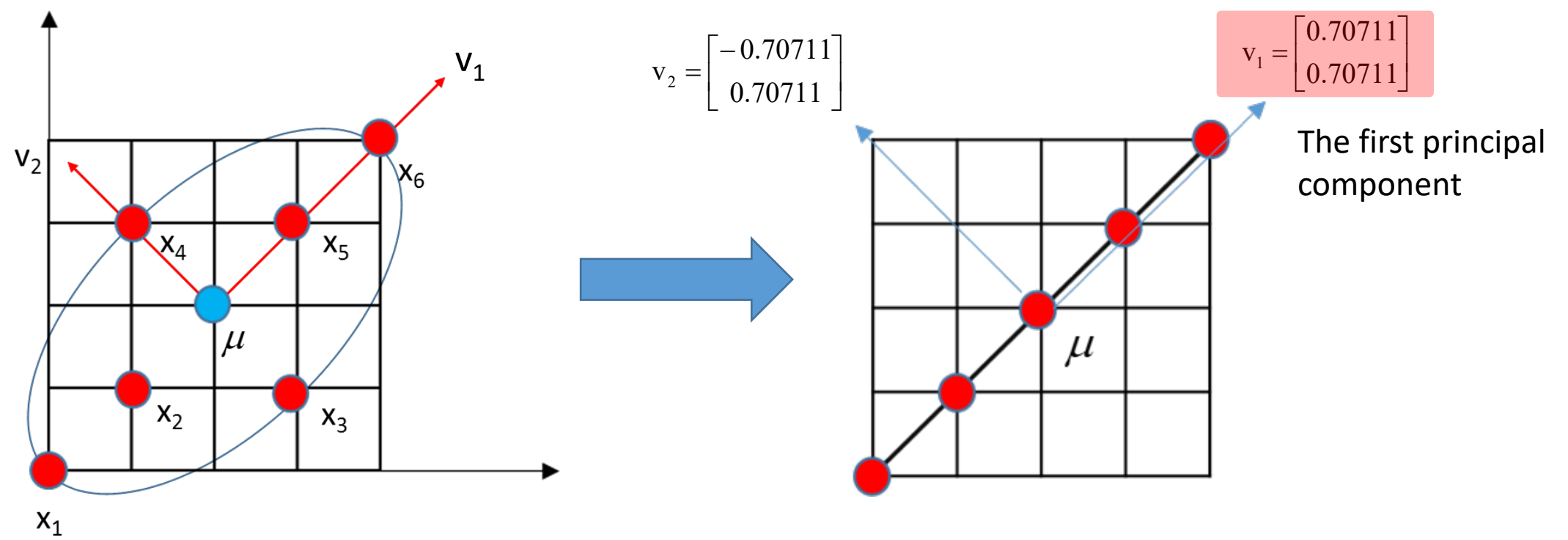
❑ Principal Component Analysis (PCA) is nothing but finding principal components of a given data set,

- Principal components are the directions where you look at the data set, which provides the most information of the data set.

- They're equivalent to eigen vectors which can be found by SVD or EVD.

- The eigen value corresponding to each eigen vector represents how widely the data set is spread along the direction which is perpendicular to the eigen vector.
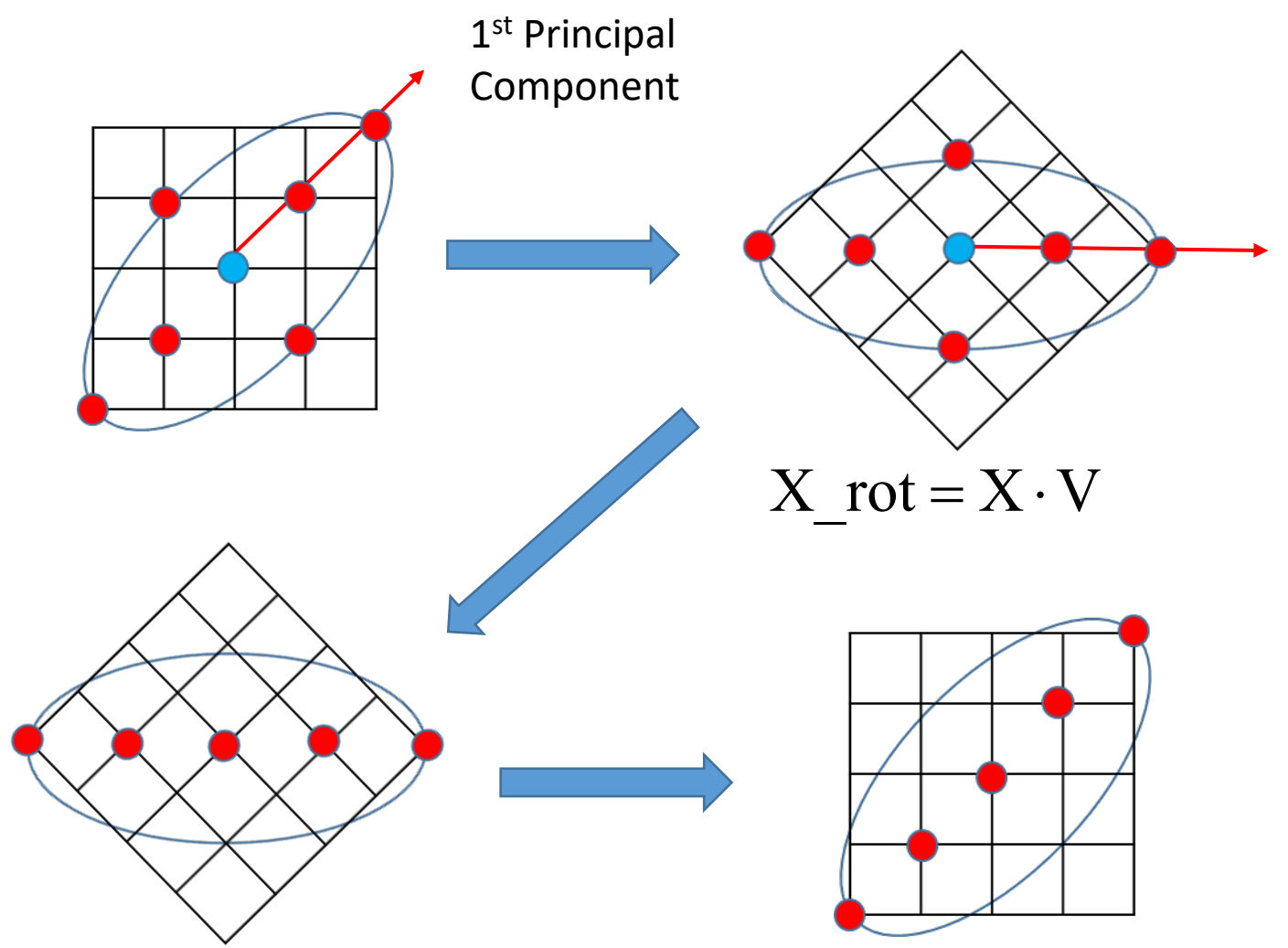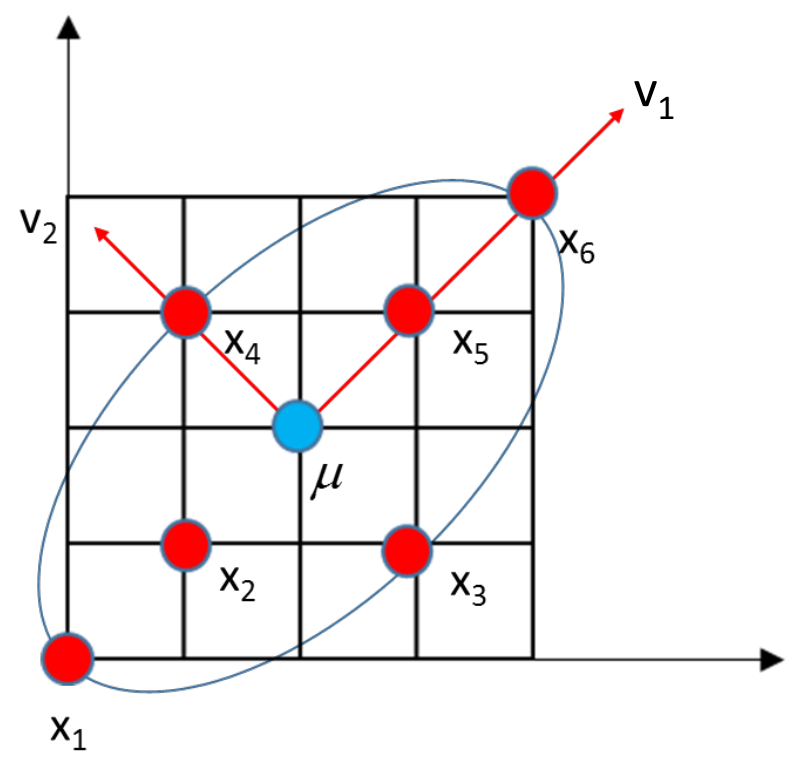
❑ A data point is defined by several, let's say, features,

❑ The number of features to define a data point is called the dimension of the data,

❑ High dimension data implies that it contains much information,

❑ Sometimes, we reduce its dimension, e.g., to visualize the data or to efficiently analyze them,

❑ PCA can reduce the dimension without losing relatively less information of the data.

- The previous example shows the case of two-dimensional data

- How can we reduce the two-dimensional data to one dimension?

- Yes, just project the data points onto the eigenvector space!



$$v_2 = \begin{bmatrix} -0.70711 \\ 0.70711 \end{bmatrix}$$

$$v_1 = \begin{bmatrix} 0.70711 \\ 0.70711 \end{bmatrix}$$

The first principal component

1st Principal Component

$$X\_rot = X \cdot V$$

Set the "$v_2$" into zero

$$X\_rot\_zero = X\_rot \cdot V^{-1}$$

$$X' = X\_rot\_zero \cdot V^{-1}$$

$$= X\_rot\_zero \cdot V^{T}$$

$$X = \begin{bmatrix} -2 & -2 \\ -1 & -1 \\ 1 & -1 \\ -1 & 1 \\ 1 & 1 \\ 2 & 2 \end{bmatrix}$$
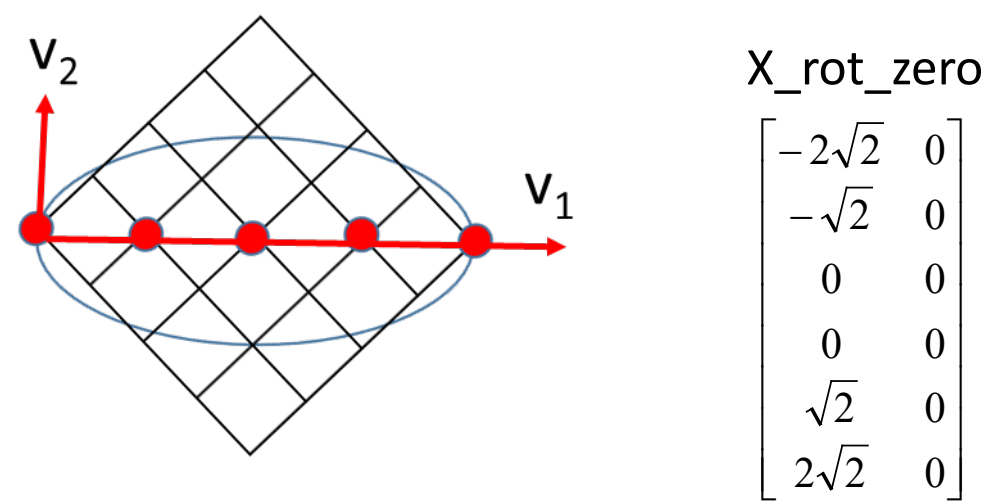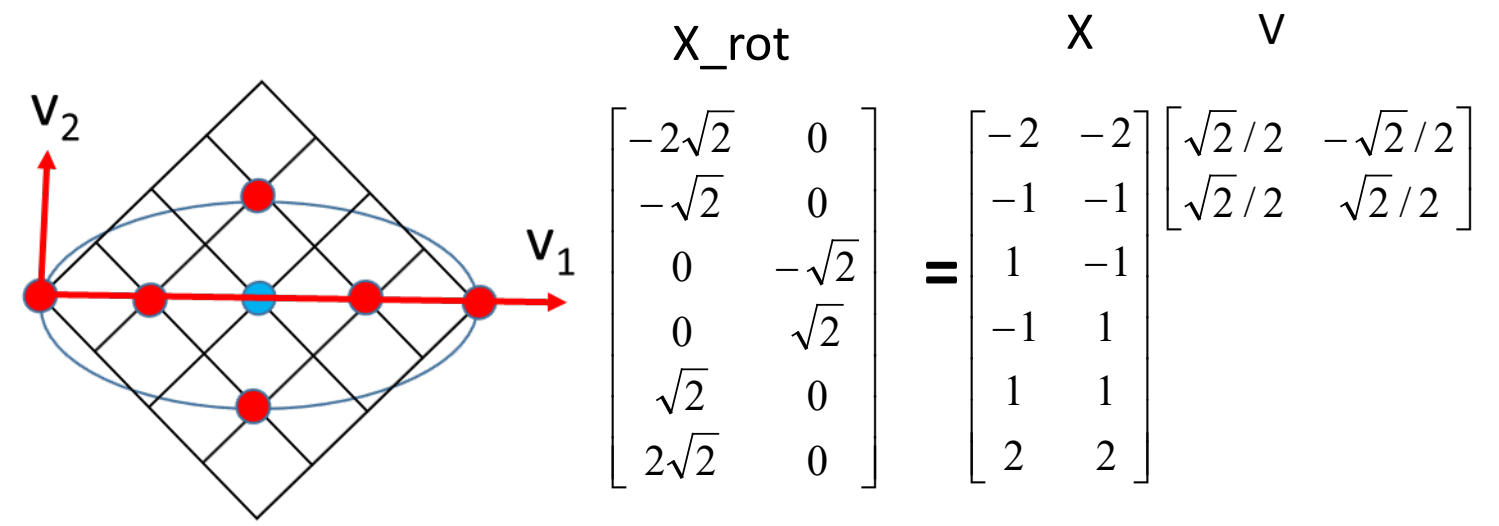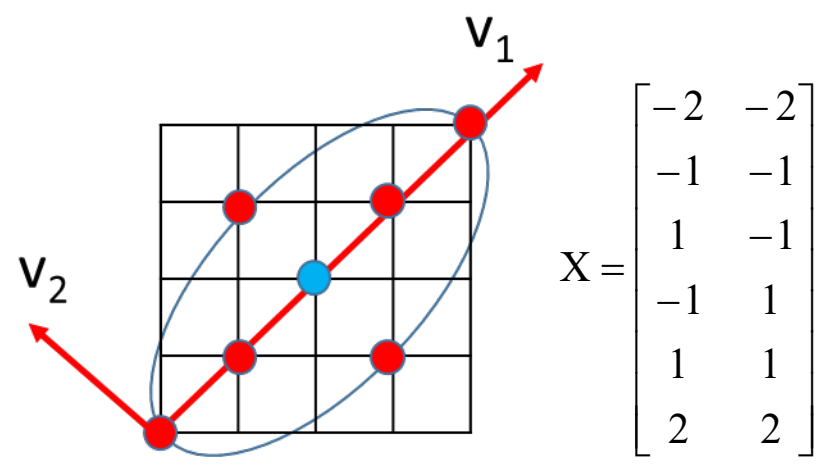
$$X = U\Sigma V^T$$

```
>> [U,S,V]=svd(x)
U =

  -0.63246    0.00000    0.30819   -0.30819    0.28637    0.57274
  -0.31623   -0.00000   -0.63635    0.63635    0.13426    0.26851
   0.00000   -0.70711    0.50000    0.50000    0.00000    0.00000
  -0.00000    0.70711    0.50000    0.50000   -0.00000   -0.00000
   0.31623    0.00000   -0.00399    0.00399    0.94140   -0.11720
   0.63246    0.00000   -0.00799    0.00799   -0.11720    0.76560

S =

Diagonal Matrix

   4.4721         0
        0    2.0000
        0         0
        0         0
        0         0
        0         0

V =

   0.70711   -0.70711
   0.70711    0.70711
```
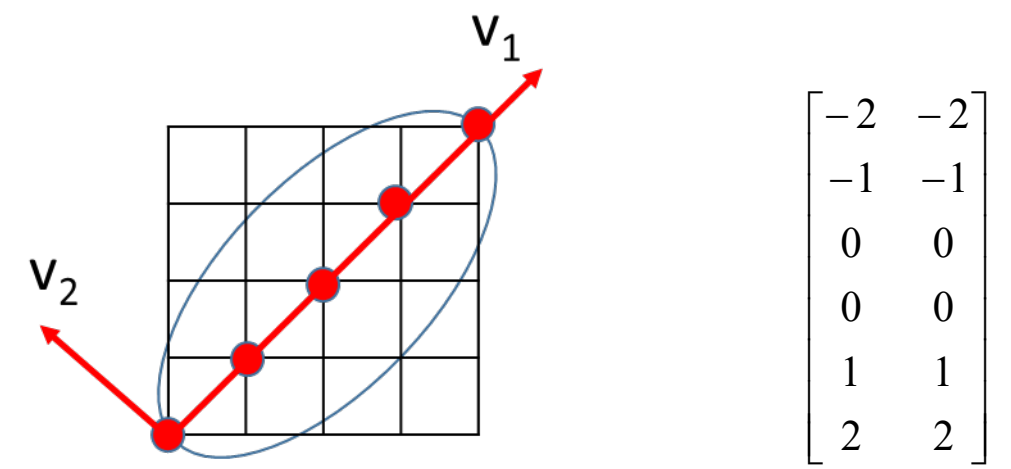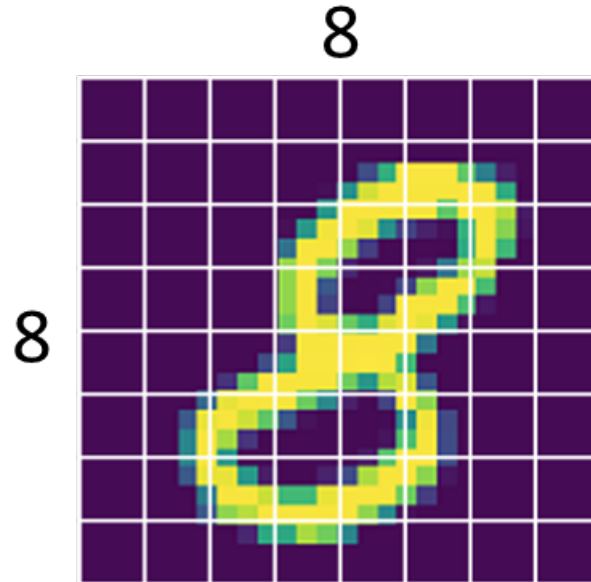
$$X = \begin{bmatrix} -2 & -2 \\ -1 & -1 \\ 1 & -1 \\ -1 & 1 \\ 1 & 1 \\ 2 & 2 \end{bmatrix}$$

X_rot       X        V

$$\begin{bmatrix} -2\sqrt{2} & 0 \\ -\sqrt{2} & 0 \\ 0 & -\sqrt{2} \\ 0 & \sqrt{2} \\ \sqrt{2} & 0 \\ 2\sqrt{2} & 0 \end{bmatrix} = \begin{bmatrix} -2 & -2 \\ -1 & -1 \\ 1 & -1 \\ -1 & 1 \\ 1 & 1 \\ 2 & 2 \end{bmatrix} \begin{bmatrix} \sqrt{2}/2 & -\sqrt{2}/2 \\ \sqrt{2}/2 & \sqrt{2}/2 \end{bmatrix}$$

X_rot_zero

$$\begin{bmatrix} -2\sqrt{2} & 0 \\ -\sqrt{2} & 0 \\ 0 & 0 \\ 0 & 0 \\ \sqrt{2} & 0 \\ 2\sqrt{2} & 0 \end{bmatrix}$$

Set the "$v_2$" elements into zero

$$\begin{bmatrix} -2 & -2 \\ -1 & -1 \\ 0 & 0 \\ 0 & 0 \\ 1 & 1 \\ 2 & 2 \end{bmatrix}$$

$$X^{'} = X\_rot\_zero \cdot V^{-1} = X\_rot\_zero \cdot V^{T}$$

❑ Let's say, we have one image representing one data point as shown below,

❑ Then, we decide to present the data by all pixels which are 64 in this case, in other words, it is 64-dimensional data,

❑ What happens if we reduce its dimension to 2 dimension?
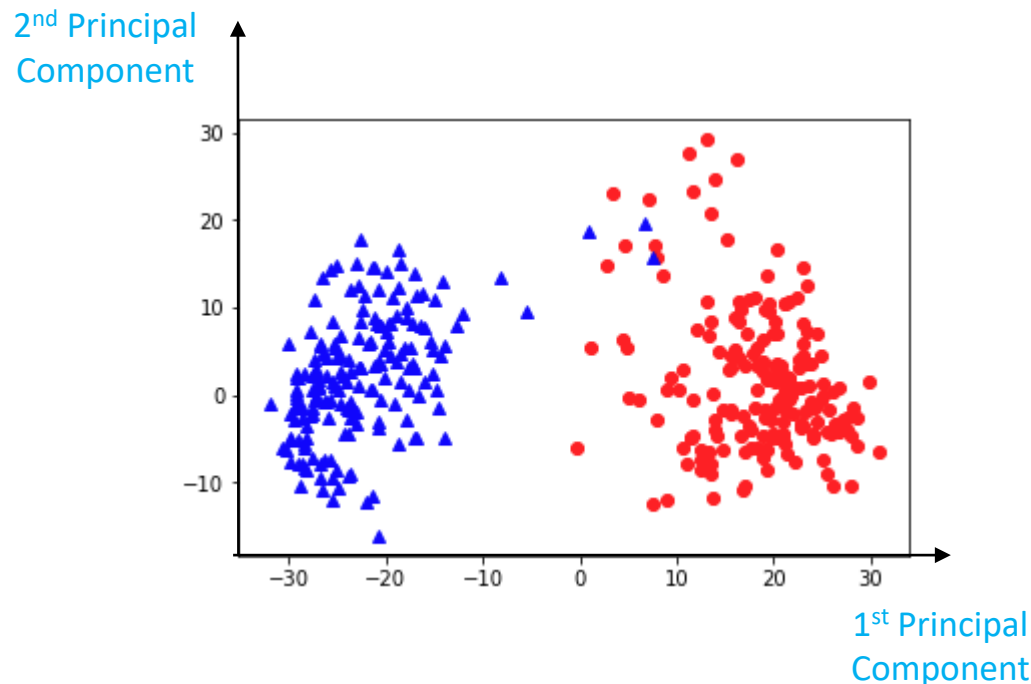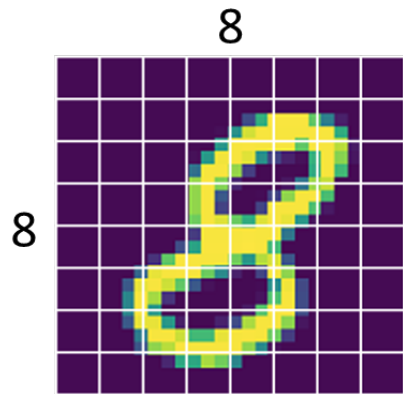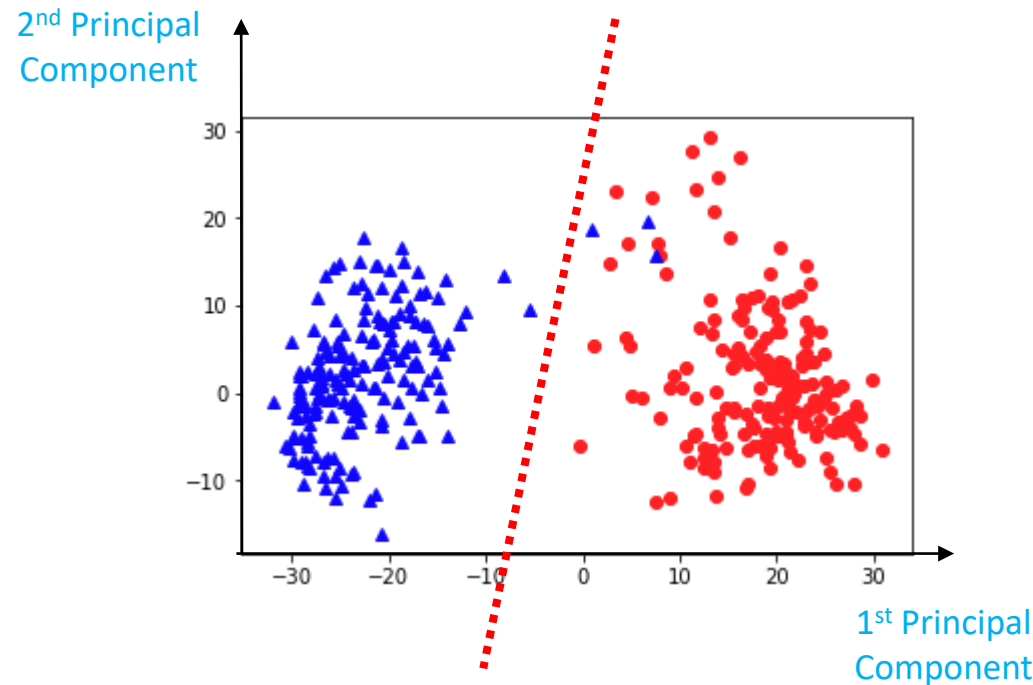
- ❑ Let's say, we have one image representing one data point as shown below,

- ❑ Then, we decide to present the data by all pixels which are 64 in this case, in other words, it is 64-dimensional data,

- ❑ What happens if we reduce its dimension to 2 dimension?

❑ Well, now we have a new set of data which have two dimension, so they can be presented in the two-dimensional space. Data visualization!

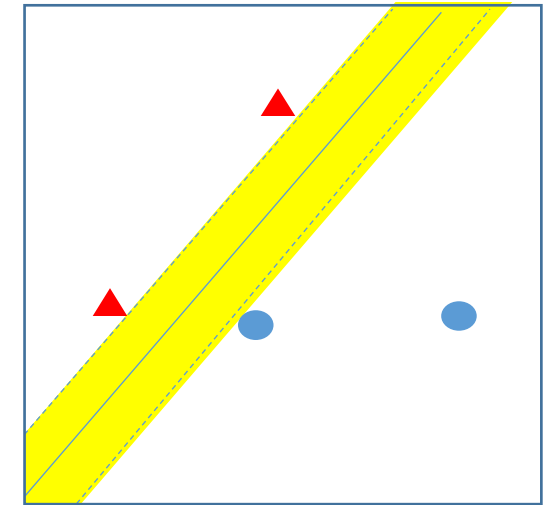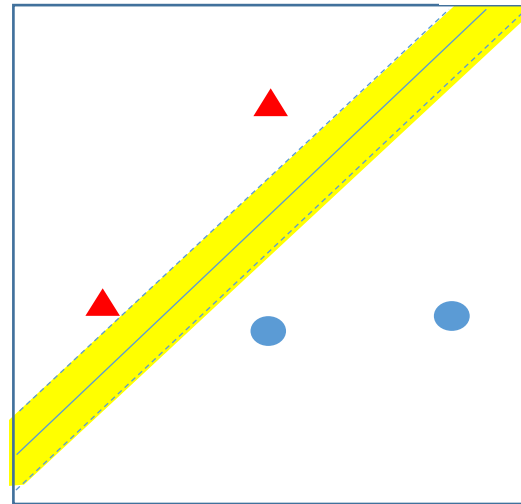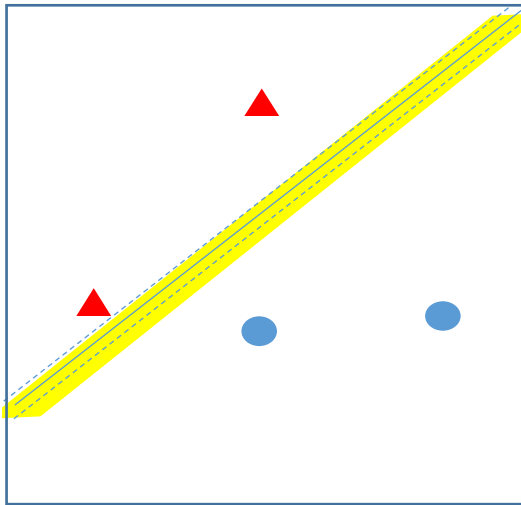❑ Also, we may be able to classify those data by drawing a line???
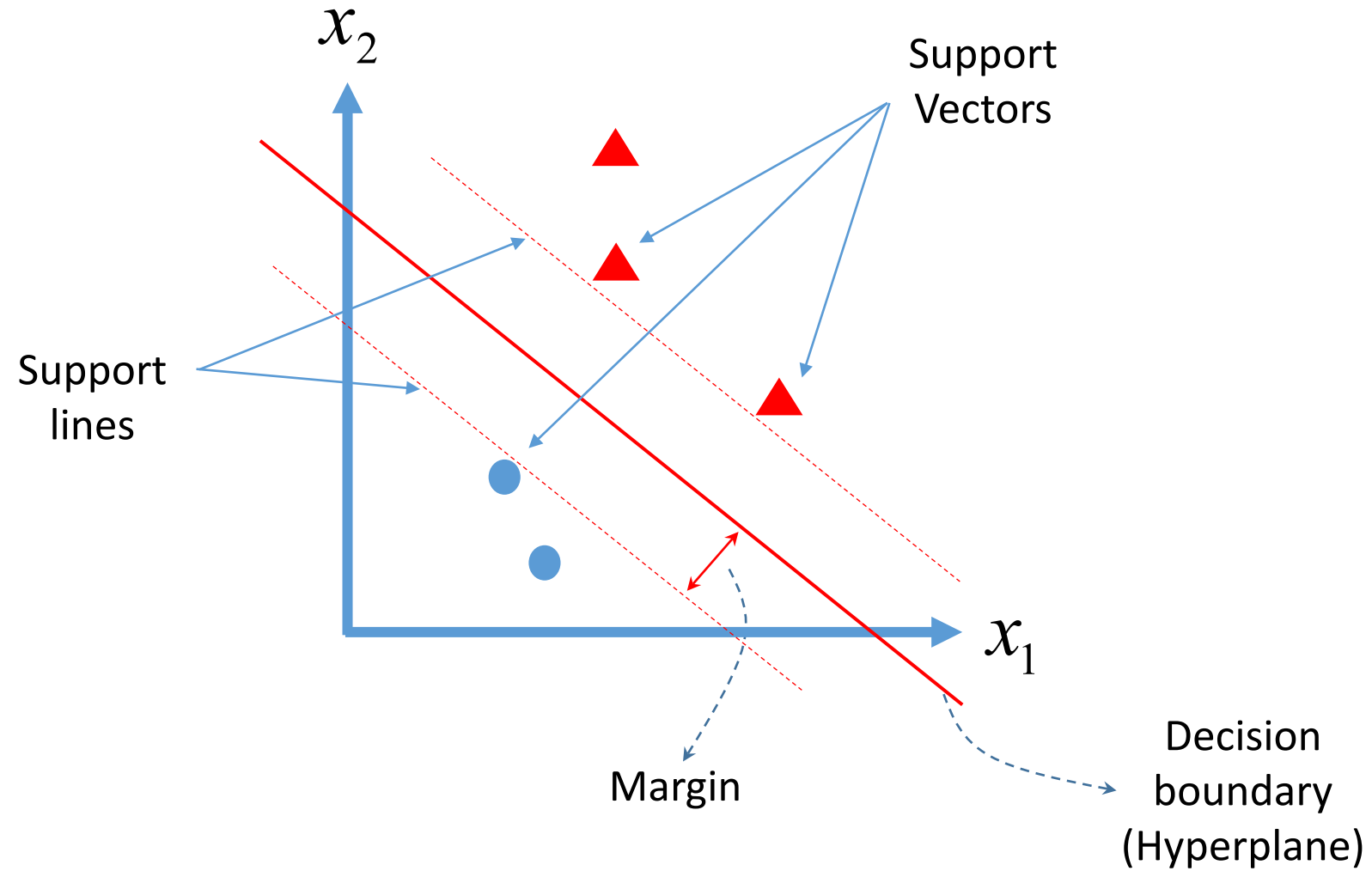
# Support Vector Machine (SVM)

❑    Most widely used classification approach (practical)

   -  Linearly separable data set

   -  Non-linearly separable data set

❑    Supported by well defined mathematical theories

   -  Geometry

   -  Optimization

Support Vectors

$x_2$

Support lines

Margin

Decision boundary (Hyperplane)

$x_1$

$$w_2 x_2 + w_1 x_1 + w_0 = y(\mathrm{x})$$

$$\mathrm{w}^\mathrm{T} \mathrm{x}^c + w_0 = y(\mathrm{x}^c)$$

$$\mathrm{w}^\mathrm{T} \mathrm{x}^b + w_0 = 0$$

Size of the vector

$$\mathrm{x}^c - \mathrm{x}^b = \|r\| \frac{\mathrm{w}}{\|\mathrm{w}\|}$$

**Margin**

Unit vector showing
the direction only

$$x^c = x^b + \|r\| \frac{w}{\|w\|}$$

$$w_2 x_2 + w_1 x_1 + w_0 = y(x)$$

$$w^T x^c + w_0 = y(x^c)$$

$$w^T x^b + w_0 = 0$$

❑ Let's multiply $w^T$ and add $w_0$ in both sides.

$$w^T x^c + w_0 = w^T x^b + w_0 + w^T \|r\| \frac{w}{\|w\|}$$

$$y(x^c) = w^T \|r\| \frac{w}{\|w\|}$$

$$\|r\| = \frac{y(x^c)}{\|w\|}$$

$$\|r\| = \frac{1}{\|w\|}$$

Let's say

$$|y(x^c)| = 1$$

We use it later …

❑ Finding a decision boundary which maximizes the margin.

$$\max \| r \| = \frac{1}{\| \mathbf{w} \|}$$

s.t.

$$t_n y(\mathbf{x}_n) > 0$$

➡ Every data points are classified correctly.

$$\begin{cases} t_n = +1, & y(\mathbf{x}_n) > 0 \\ t_n = -1, & y(\mathbf{x}_n) < 0 \end{cases}$$

$x_2$

$\mathbf{x}^c$

$\mathbf{w}$

$\mathbf{x}^b$ $\| r \|$

$x_1$

$$w_2 x_2 + w_1 x_1 + w_0 = 0$$

**25**

❑ Let's modify the optimization problem a bit.

$$\max \frac{1}{\parallel w \parallel}$$

$$s.t. \quad t_n y(x_n) > 0, \quad \forall n$$

❑ Do you remember?

$$\max \frac{1}{\parallel w \parallel}$$

Let's say

$$| y(x^c) | = 1$$

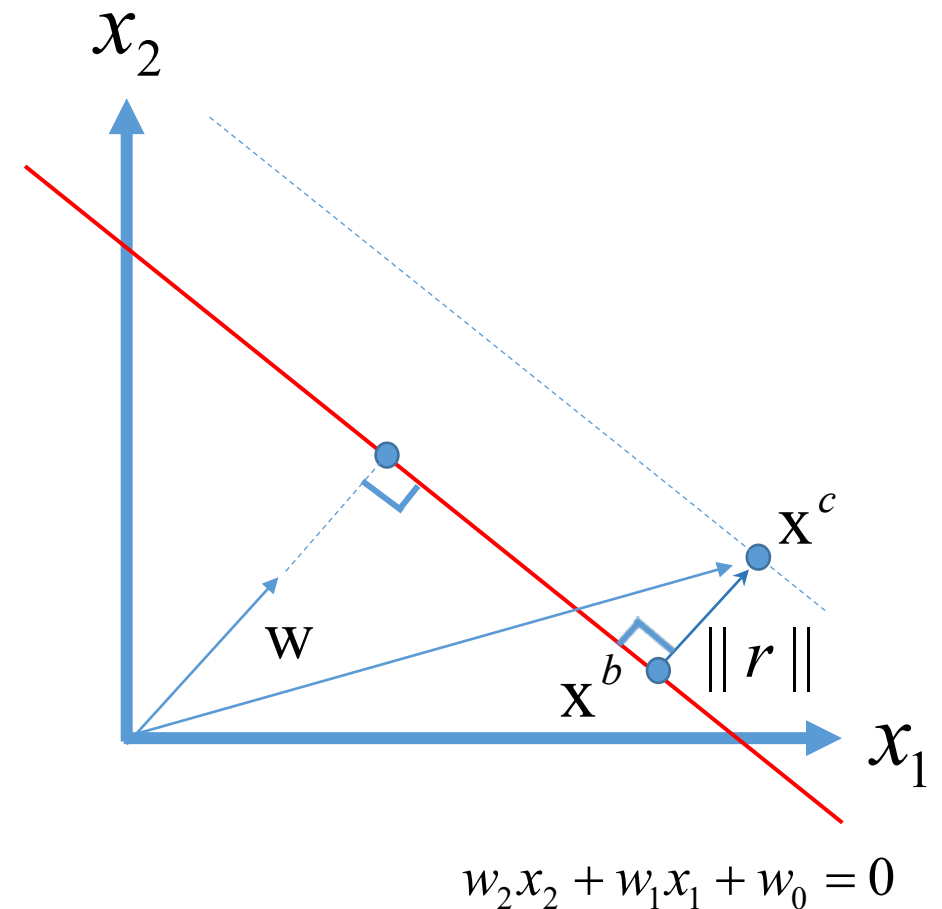$$s.t. \quad t_n y(x_n) \geq 1, \quad \forall n$$

meaning that any data point is away
from the decision boundary at least 1

❑ Finally

$$\min \frac{1}{2} \parallel w \parallel^2$$

$$s.t. \quad t_n (w^T x_n + w_0) \geq 1, \quad \forall n$$

Quadratic programming

$x_2$

$x^c$

$w$

$x^b$ $\parallel r \parallel$

$x_1$

$$w_2 x_2 + w_1 x_1 + w_0 = 0$$

$$\min \quad \frac{1}{2}||w||^2$$

$$s.t. \quad t(x_n)(W \cdot x_n + w_0) \geq 1, \forall n$$

# Kernel Trick

$$\min x^2$$
$$s.t. \quad x \geq b$$

**?**
**=**

$$\min_{x} \max_{\lambda} x^2 - \lambda(x - b)$$
$$s.t. \quad \lambda \geq 0$$

$x^2$



$x^2$



- ❑ If b ≤ 0, the minima is 0 … so λ=0

- ❑ If b > 0, the minima is $b^2$ … so x=b

- ❑ So, either λ or (x − b) becomes zero, in other words,
  - λ(x-b) = 0 (complementary slackness)

- ❑ Since x ≥ b,  maximizing λ minimizes the objective value
  - λ ≥ 0

**29**

$$\min \frac{1}{2} \mathbf{w}^T \mathbf{w}$$

$$s.t. \quad t_n(\mathbf{w}^{\mathrm{T}} x_n + w_0) \geq 1$$

→

$$\min_{\mathbf{w}} \max_{\lambda} \frac{1}{2} \mathbf{w}^T \mathbf{w} - \sum_{n=1}^{n} \lambda_n (t_n(\mathbf{w}^{\mathrm{T}} x_n + w_0) - 1)$$

$$s.t. \quad \lambda_n \geq 0$$

**30**

# Proof begins

$$\min \frac{1}{2} \mathbf{w}^T \mathbf{w}$$

$$s.t. \quad t_n(\mathbf{w}^\mathrm{T} x_n + w_0) \geq 1$$

$$\min_{\mathbf{w}} \max_{\lambda} \frac{1}{2} \mathbf{w}^T \mathbf{w} - \sum_{n=1}^{n} \lambda_n (t_n(\mathbf{w}^\mathrm{T} x_n + w_0) - 1)$$

$$s.t. \quad \lambda_n \geq 0$$

- ❑ We would like to convert again the optimization problem above into another form, which provides same results.
  - Because we want to solve the optimization problem in term of "lagrange multiplier ($\lambda_n$)".

$$\max_{\lambda} \min_{\mathbf{w}} \frac{1}{2} \mathbf{w}^T \mathbf{w} - \sum_{n=1}^{n} \lambda_n (t_n(\mathbf{w}^\mathrm{T} x_n + w_0) - 1)$$

$$s.t. \quad \lambda_n \geq 0$$

Primal problem

$$\min_{\mathbf{w}} \max_{\lambda} \frac{1}{2}\mathbf{w}^T\mathbf{w} - \sum_{n=1}^{n} \lambda_n(t_n(\mathbf{w}^T x_n + w_0) - 1)$$

$$s.t. \quad \lambda_n \geq 0$$

❑ We would like to convert again the optimization problem above into another form, which provides same results.

- Because we want to solve the optimization problem in term of "lagrange multiplier ($\lambda_n$)".

$$\max_{\lambda} \min_{\mathbf{w}} \frac{1}{2}\mathbf{w}^T\mathbf{w} - \sum_{n=1}^{n} \lambda_n(t_n(\mathbf{w}^T x_n + w_0) - 1)$$

$$s.t. \quad \lambda_n \geq 0$$

Dual problem

HALF FULL

HALF EMPTY

<u>Karush–Kuhn–Tucker conditions</u>

**KKT conditions**

1) Stationarity condition

$$\frac{\partial}{\partial \mathbf{w}} \frac{1}{2} \mathbf{w}^T \mathbf{w} - \frac{\partial}{\partial \mathbf{w}} \sum_{n=1}^{n} \lambda_n (t_n (\mathbf{w}^T x_n + w_0) - 1) = 0$$

2) Complementary slackness condition

$$\lambda_n (t_n (\mathbf{w}^T x_n + w_0) - 1) = 0$$

3) Duality feasibility condition

$$\lambda_n \geq 0$$

**Primal problem**

$$\min_{\mathbf{w}} \max_{\lambda} \frac{1}{2} \mathbf{w}^T \mathbf{w} - \sum_{n=1}^{n} \lambda_n (t_n (\mathbf{w}^T x_n + w_0) - 1)$$

$$s.t. \quad \lambda_n \geq 0$$

❑ We would like to convert again the optimization problem above into another form, which provides same results.
- Because we want to solve the optimization problem in term of "lagrange multiplier ($\lambda_n$)".

$$\max_{\lambda} \min_{\mathbf{w}} \frac{1}{2} \mathbf{w}^T \mathbf{w} - \sum_{n=1}^{n} \lambda_n (t_n (\mathbf{w}^T x_n + w_0) - 1)$$

$$s.t. \quad \lambda_n \geq 0$$

**Dual problem**

$$\max_{\lambda} \min_{w,w_0} \quad L(\mathrm{w}, \mathrm{w}_0, \lambda) = \frac{1}{2}\mathrm{w}^{\mathrm{T}}\mathrm{w} - \sum_{n=1}^{N} \lambda_n \left( t_n (\mathrm{w}^{\mathrm{T}} x_n + w_0) - 1 \right)$$

$$\mathrm{w} = \sum_{n=1}^{N} \lambda_n t_n x_n$$

$$\frac{\partial L}{\partial w} = \mathrm{w} - \sum_{n=1}^{N} \lambda_n t_n x_n = 0$$

❖ The first one is called stationarity condition.

➢ when we partial differentiate the problem with respect to its parameter "w", each of them should be zero.

$$\max_{\lambda} \min_{w,w_0} \quad L(\mathrm{w}, \mathrm{w}_0, \lambda) = \frac{1}{2}\mathrm{w}^{\mathrm{T}}\mathrm{w} - \sum_{n=1}^{N} \lambda_n \left( t_n (\mathrm{w}^{\mathrm{T}} x_n + w_0) - 1 \right)$$

$$\mathrm{w} = \sum_{n=1}^{N} \lambda_n t_n x_n \qquad\qquad \sum_{n=1}^{N} \lambda_n t_n = 0$$

$$\frac{\partial L}{\partial w} = \mathrm{w} - \sum_{n=1}^{N} \lambda_n t_n x_n = 0 \qquad\qquad \frac{\partial L}{\partial w_0} = -\sum_{n=1}^{N} \lambda_n t_n = 0$$

❖ The first one is called stationarity condition.

➤ Again, this time in terms of "$w_0$"

$$\max_{\lambda} \min_{w,w_0} \quad L(\mathrm{w}, \mathrm{w}_0, \lambda) = \frac{1}{2}\mathrm{w}^{\mathrm{T}}\mathrm{w} \; - \sum_{n=1}^{N} \lambda_n \left( t_n(\mathrm{w}^{\mathrm{T}}x_n + w_0) - 1 \right)$$

$$\mathrm{w} = \sum_{n=1}^{N} \lambda_n t_n x_n$$

$$\sum_{n=1}^{N} \lambda_n t_n = 0$$

$$L(\lambda) = \frac{1}{2}\sum_{n=1}^{N}\sum_{m=1}^{N} t_n t_m \lambda_n \lambda_m \mathrm{x}_n^T \mathrm{x}_m \; - \sum_{n=1}^{N}\sum_{m=1}^{N} t_n t_m \lambda_n \lambda_m \mathrm{x}_n^T \mathrm{x}_m \; - \sum_{n=1}^{N} \lambda_n t_n w_0 + \sum_{n=1}^{N} \lambda_n$$

$$\max_{\lambda} \min_{w,w_0} \quad L(\mathrm{w}, \mathrm{w}_0, \lambda) = \frac{1}{2}\mathrm{w}^{\mathrm{T}}\mathrm{w} \quad - \sum_{n=1}^{N} \lambda_n\left(t_n(\mathrm{w}^{\mathrm{T}}x_n + w_0) - 1\right)$$

$$\mathrm{w} = \sum_{n=1}^{N} \lambda_n t_n x_n$$

$$\sum_{n=1}^{N} \lambda_n t_n = 0$$

$$L(\lambda) = \frac{1}{2}\sum_{n=1}^{N}\sum_{m=1}^{N} t_n t_m \lambda_n \lambda_m \mathrm{x}_n^{T}\mathrm{x}_m - \sum_{n=1}^{N}\sum_{m=1}^{N} t_n t_m \lambda_n \lambda_m \mathrm{x}_n^{T}\mathrm{x}_m - \sum_{n=1}^{N} \lambda_n t_n w_0 + \sum_{n=1}^{N} \lambda_n$$

$$\max_{\lambda} \quad L(\lambda) = \sum_{n=1}^{N} \lambda_n - \frac{1}{2}\sum_{n=1}^{N}\sum_{m=1}^{N} t_n t_m \lambda_n \lambda_m \mathrm{x}_n^{T}\mathrm{x}_m$$

$$\max_{\lambda} \min_{w,w_0} \quad L(\mathrm{w}, \mathrm{w}_0, \lambda) = \frac{1}{2}\mathrm{w}^{\mathrm{T}}\mathrm{w} - \sum_{n=1}^{N}\lambda_n\left(t_n(\mathrm{w}^{\mathrm{T}}x_n + w_0) - 1\right)$$

$$\mathrm{w} = \sum_{n=1}^{N}\lambda_n t_n x_n \qquad \sum_{n=1}^{N}\lambda_n t_n = 0$$

$$L(\lambda) = \frac{1}{2}\sum_{n=1}^{N}\sum_{m=1}^{N}t_n t_m \lambda_n \lambda_m \mathrm{x}_n^T \mathrm{x}_m - \sum_{n=1}^{N}\sum_{m=1}^{N}t_n t_m \lambda_n \lambda_m \mathrm{x}_n^T \mathrm{x}_m - \sum_{n=1}^{N}\lambda_n t_n w_0 + \sum_{n=1}^{N}\lambda_n$$

$$\max_{\lambda} L(\lambda) = \sum_{n=1}^{N}\lambda_n - \frac{1}{2}\sum_{n=1}^{N}\sum_{m=1}^{N}t_n t_m \lambda_n \lambda_m \mathrm{x}_n^T \mathrm{x}_m$$

$$\lambda_n \geq 0 \qquad \sum_{n=1}^{N}\lambda_n t_n = 0 \qquad \mathrm{w} = \sum_{n=1}^{N}\lambda_n t_n x_n$$

"w" does not appear in the equation, and so we do not use this constraint anymore

$$\max_{\lambda} \quad L(\lambda) = \sum_{n=1}^{N} \lambda_n - \frac{1}{2} \sum_{n=1}^{N} \sum_{m=1}^{N} t_n t_m \lambda_n \lambda_m \mathbf{x}_n^T \mathbf{x}_m$$

$$s.t. \quad \lambda_n \geq 0, \quad \sum_{n=1}^{N} \lambda_n t_n = 0$$

- ❑ Let's change it to a quadratic programming again.
- ❑ As mentioned previously, a quadratic programming problem needs to be minimized

40

$$\max_{\lambda} \quad L(\lambda) = \sum_{n=1}^{N} \lambda_n - \frac{1}{2} \sum_{n=1}^{N} \sum_{m=1}^{N} t_n t_m \lambda_n \lambda_m \mathbf{x}_n^T \mathbf{x}_m$$

$$s.t. \quad \lambda_n \geq 0, \quad \sum_{n=1}^{N} \lambda_n t_n = 0$$

$$\min_{\lambda} \quad L(\lambda) = \frac{1}{2} \sum_{n=1}^{N} \sum_{m=1}^{N} t_n t_m \lambda_n \lambda_m \mathbf{x}_n^T \mathbf{x}_m - \sum_{n=1}^{N} \lambda_n$$

$$s.t. \quad \lambda_n \geq 0, \quad \sum_{n=1}^{N} \lambda_n t_n = 0$$

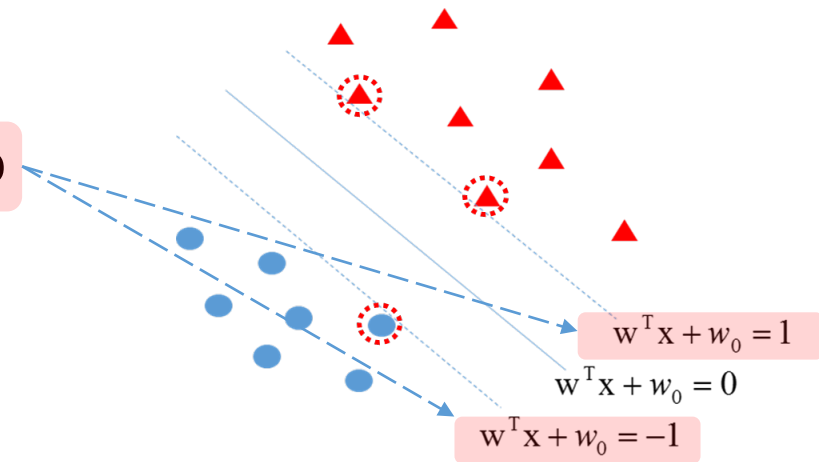❑ Again, the optimization problem becomes a quadratic programming problem.

# Let's summarize

$$\min_{\lambda} \quad L(\lambda) = \frac{1}{2}\sum_{n=1}^{N}\sum_{m=1}^{N} t_n t_m \lambda_n \lambda_m \mathbf{x}_n^T \mathbf{x}_m - \sum_{n=1}^{N}\lambda_n$$

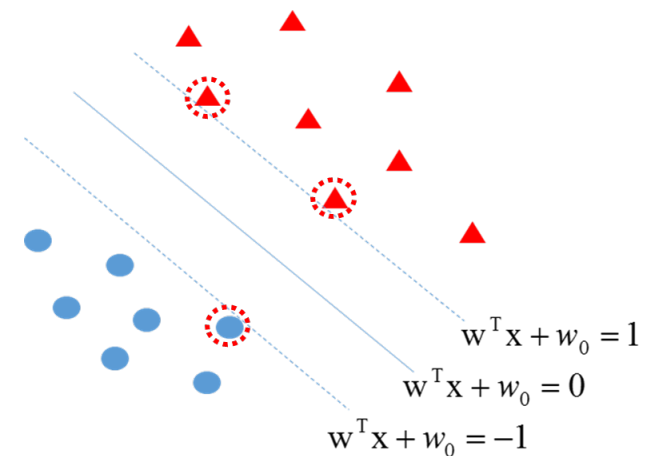$$s.t. \quad \lambda \geq 0, \quad t^T \lambda = 0$$

- ❑ The solution from the quadratic programming is "lagrange multipliers"($\lambda_n$)
- ❑ Many of the solutions (lagrange multipliers) are zero
- ❑ Complementary slackness (one of KKT conditions) should be satisfied.

$$\lambda_n(t_n(\mathbf{w}^T x_n + w_0) - 1) = 0$$

- ❑ In other words, if $\lambda_n$ are not zero, $(t_n(w_t x_n + w_0) - 1)$ should be zero
  where corresponding data points should be support vectors.

$$\mathbf{w}^T \mathbf{x} + w_0 = 1$$
$$\mathbf{w}^T \mathbf{x} + w_0 = 0$$
$$\mathbf{w}^T \mathbf{x} + w_0 = -1$$

# Let's summarize

$$\min_{\lambda} \quad L(\lambda) = \frac{1}{2}\sum_{n=1}^{N}\sum_{m=1}^{N} t_n t_m \lambda_n \lambda_m x_n^T x_m - \sum_{n=1}^{N} \lambda_n$$

$$s.t. \quad \lambda \geq 0, \quad t^T \lambda = 0$$

- ❑ The solution from the quadratic programming is "lagrange multipliers"($\lambda_n$)
- ❑ Many of the solutions (lagrange multipliers) are zero
- ❑ Complementary slackness (one of KKT conditions) should be satisfied.

$$\lambda_n (t_n(w^T x_n + w_0) - 1) = 0$$

- ❑ In other words, if $\lambda_n$ are not zero, $(t_n(w_t x_n + w_0)-1)$ should be zero where corresponding data points should be support vectors.
- ❑ With the non-zero $\lambda_n$, w and $w_0$ can be calculated using $t_n(w_t x_n + w_0)=1$

$$w = \sum_{n=1}^{N} \lambda_n t_n x_n \qquad w_0 = t_n - \sum_{n=1}^{N} \lambda_n t_n x_n x_n$$

$$w^T x + w_0 = 1$$
$$w^T x + w_0 = 0$$
$$w^T x + w_0 = -1$$

$$\min_{\lambda} \quad L(\lambda) = \frac{1}{2}\sum_{n=1}^{N}\sum_{m=1}^{N} t_n t_m \lambda_n \lambda_m \mathbf{x}_n^T \mathbf{x}_m - \sum_{n=1}^{N} \lambda_n$$

$$s.t. \quad \lambda \geq 0, \quad t^T \lambda = 0$$

We obtained previously

❑ $t_n(w_t x_n + w_0) = 1$

$\Rightarrow (w_t x_n + w_0) = t_n$

$\Rightarrow w_0 = t_n - w_t x_n$

$$w_0 = t_n - \sum_{n=1}^{N} \lambda_n t_n x_n x_n$$

❑ The solution from the quadratic programming is "lagrange multipliers"($\lambda_n$)

❑ Many of the solutions (lagrange multipliers) are zero

❑ Complementary slackness (one of KKT conditions) should be satisfied.

$$\lambda_n(t_n(\mathbf{w}^T x_n + w_0) - 1) = 0$$

❑ In other words, if $\lambda_n$ are not zero, $(t_n(w_t x_n + w_0) - 1)$ should be zero
   where corresponding data points should be support vectors.

❑ With the non-zero $\lambda_n$, w and $w_0$ can be calculated using $t_n(w_t x_n + w_0) = 1$
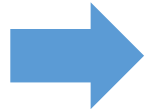
$$\mathbf{w} = \sum_{n=1}^{N} \lambda_n t_n x_n$$

$$w_0 = t_n - \sum_{n=1}^{N} \lambda_n t_n x_n x_n$$

Proof ends

$$\min \frac{1}{2}\mathbf{w}^T\mathbf{w}$$

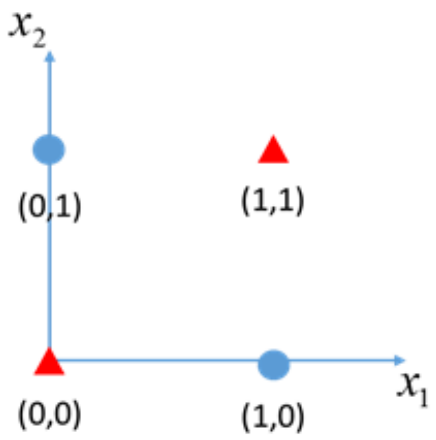$$s.t. \quad t_n(\mathbf{w}^T x_n + w_0) \geq 1$$

$$\min_{\lambda} \quad L(\lambda) = \frac{1}{2}\sum_{n=1}^{N}\sum_{m=1}^{N} t_n t_m \lambda_n \lambda_m \mathbf{x}_n^T \mathbf{x}_m - \sum_{n=1}^{N} \lambda_n$$

$$s.t. \quad \lambda \geq 0, \quad t^T \lambda = 0$$

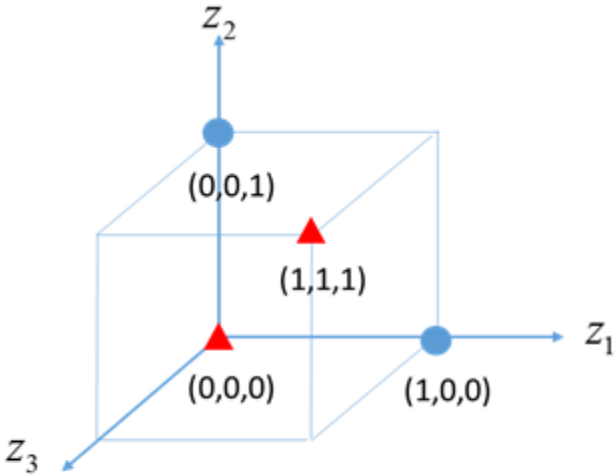❑ If data $x_n$ are not linearly separable, what should we do?



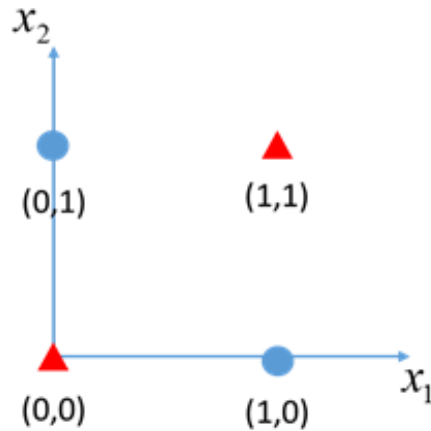$$\phi(\mathbf{x}) = \begin{pmatrix} x_1^2 \\ x_1 x_2 \\ x_2^2 \end{pmatrix}$$
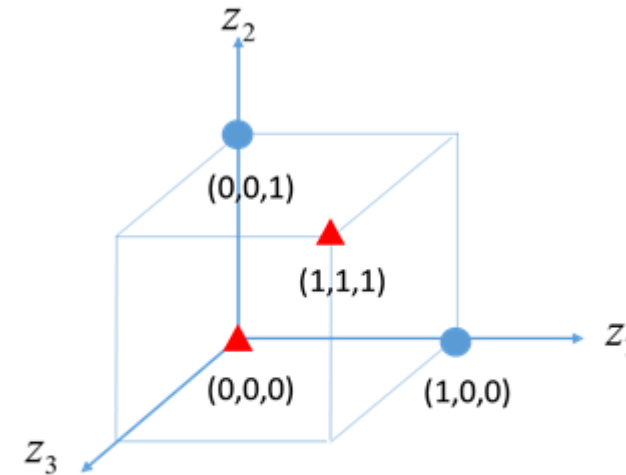
$$\phi(\mathbf{x})$$
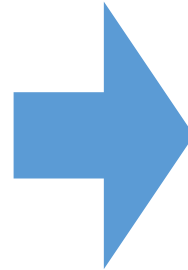
$$X \rightarrow Z$$

Space X

Space Z

❑ The idea of Kernel trick begins from here: to find the scalar values (the inner product of two vectors: $z_n$ and $z_m$ ) and so we can formulate the quadratic problem which can be linearly separable.
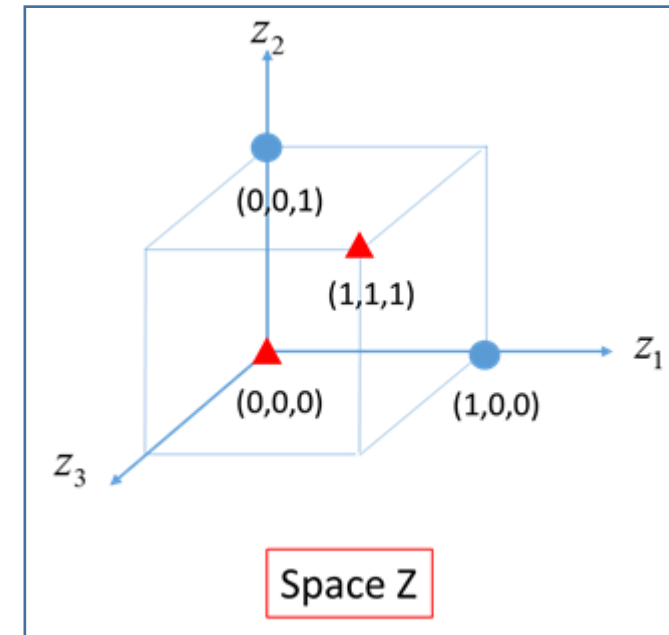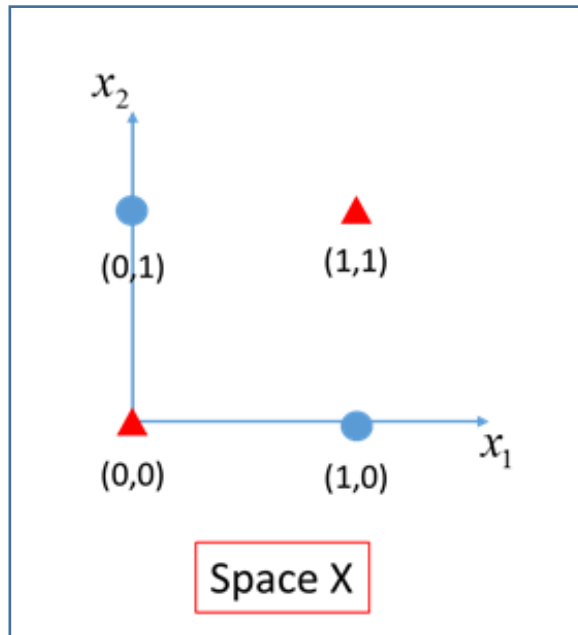


Space X

Space Z

$$\min_{\lambda} \ L(\lambda) = \ \frac{1}{2}\sum_{n=1}^{N}\sum_{m=1}^{N} t_n t_m \lambda_n \lambda_m \mathbf{x}_n^T \mathbf{x}_m - \sum_{n=1}^{N} \lambda_n$$

$$s.t. \quad \lambda \geq 0, \quad t^T \lambda = 0$$

$$\min_{\lambda} \ L(\lambda) = \ \frac{1}{2}\sum_{n=1}^{N}\sum_{m=1}^{N} t_n t_m \lambda_n \lambda_m \mathbf{z}_n^T \mathbf{z}_m - \sum_{n=1}^{N} \lambda_n$$

$$s.t. \quad \lambda \geq 0, \quad t^T \lambda = 0$$

❑ Kernel function K() is a function which returns the scalar values (<span style="color:red">the inner product of two vectors: $z_n$ and $z_m$ in Z space</span>) when the data points (<span style="color:red">$x_n$ and $x_m$ in X space</span>) are given.

$$K(\mathrm{x}_n^T, \mathrm{x}_m) = \mathrm{z}_n^T \mathrm{z}_m$$

❑ With the Kernel function defined previously, we want to change the quadratic problem as follows:

- Because the Kernel function is a function of data points ($x_n$ and $x_m$) which we already have.

$$\min_{\lambda} \quad L(\lambda) = \frac{1}{2}\sum_{n=1}^{N}\sum_{m=1}^{N} t_n t_m \lambda_n \lambda_m \mathbf{z}_n^T \mathbf{z}_m - \sum_{n=1}^{N} \lambda_n$$

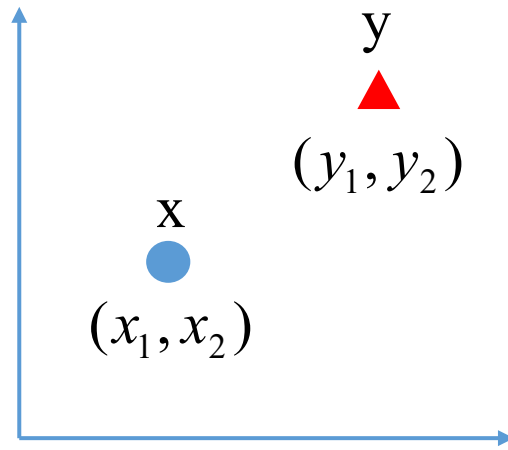$$s.t. \quad \lambda \geq 0, \quad t^T \lambda = 0$$

$$\min_{\lambda} \quad L(\lambda) = \frac{1}{2}\sum_{n=1}^{N}\sum_{m=1}^{N} t_n t_m \lambda_n \lambda_m \mathrm{K}(\mathbf{x}_n^T \mathbf{x}_m) - \sum_{n=1}^{N} \lambda_n$$

$$s.t. \quad \lambda \geq 0, \quad t^T \lambda = 0$$

❑ With the Kernel function defined previously, we want to change the quadratic problem as follows:
- Because the Kernel function is a function of data points ($x_n$ and $x_m$ ) which we already have.

$$\min_{\lambda} \quad L(\lambda) = \frac{1}{2}\sum_{n=1}^{N}\sum_{m=1}^{N}t_n t_m \lambda_n \lambda_m z_n^T z_m - \sum_{n=1}^{N}\lambda_n$$
$$s.t. \quad \lambda \geq 0, \quad t^T\lambda = 0$$

$$\min_{\lambda} \quad L(\lambda) = \frac{1}{2}\sum_{n=1}^{N}\sum_{m=1}^{N}t_n t_m \lambda_n \lambda_m K(x_n^T x_m) - \sum_{n=1}^{N}\lambda_n$$
$$s.t. \quad \lambda \geq 0, \quad t^T\lambda = 0$$

$$\min_{\lambda} \quad L(\lambda) = \frac{1}{2}\lambda^T\begin{bmatrix} t_1 t_1 K(x_1,x_1) & t_1 t_2 K(x_1^T,x_2) & \cdots & t_1 t_N K(x_1^T,x_N) \\ t_2 t_1 K(x_2,x_1) & t_2 t_2 K(x_2^T,x_2) & \cdots & t_2 t_N K(x_2^T,x_N) \\ \ldots & \ldots & \ldots & \ldots \\ t_N t_1 K(x_N x_1) & t_N t_2 K(x_N^T,x_2) & \cdots & t_N t_N K(x_N^T,x_N) \end{bmatrix}\lambda + (-1^T)\lambda$$

$\boxed{\text{Space X}}$

$$K(\mathrm{x}, y) = (\mathrm{xy})^2$$

$$= \left( (x_1, x_2) \cdot (y_1, y_2) \right)^2$$

$$= \left( x_1 y_1 + x_2 y_2 \right)^2$$

$$= x_1^2 y_1^2 + 2 x_1 x_2 y_1 y_2 + x_2^2 y_2^2$$

y

▲

$(y_1, y_2)$

x

●

$(x_1, x_2)$

$\boxed{\text{Space X}}$

$\phi(\text{x})$

●

$(x_1^2, \sqrt{2}x_1x_2, x_2^2)$

$\phi(y)$

▲

$(y_1^2, \sqrt{2}y_1y_2, y_2^2)$

$\boxed{\text{Space Z}}$

$$K(\text{x}, y) = (\text{xy})^2$$
$$= ((x_1, x_2) \cdot (y_1, y_2))^2$$
$$= (x_1y_1 + x_2y_2)^2$$
$$= x_1^2 y_1^2 + 2x_1x_2y_1y_2 + x_2^2 y_2^2$$

$$\phi(\text{x})\phi(y) = (x_1^2, \sqrt{2}x_1x_2, x_2^2) \cdot (y_1^2, \sqrt{2}y_1y_2, y_2^2)$$
$$= x_1^2 y_1^2 + 2x_1x_2y_1y_2 + x_2^2 y_2^2$$

$$K(\mathrm{x}_n, \mathrm{x}_m) = \exp\left(-\alpha \parallel \mathrm{x_n} - \mathrm{x_m} \parallel^2\right)$$

$$= \exp\left(-\alpha \mathrm{x}_n^2\right)\exp\left(-\alpha \mathrm{x}_m^2\right)\exp\left(2\alpha \mathrm{x}_n \mathrm{x}_m\right)$$

Taylor series expansion of an exponential function

$$\exp(x) = \frac{x^0}{0!} + \frac{x}{1!} + \frac{x^2}{2!} + \frac{x^3}{3!} + \cdots$$

$$= \exp\left(-\alpha \mathrm{x}_n^2\right)\exp\left(-\alpha \mathrm{x}_m^2\right)\sum_{k=0}^{\infty} \frac{(2\alpha)^k (\mathrm{x_n})^k (\mathrm{x_m})^k}{k!}$$
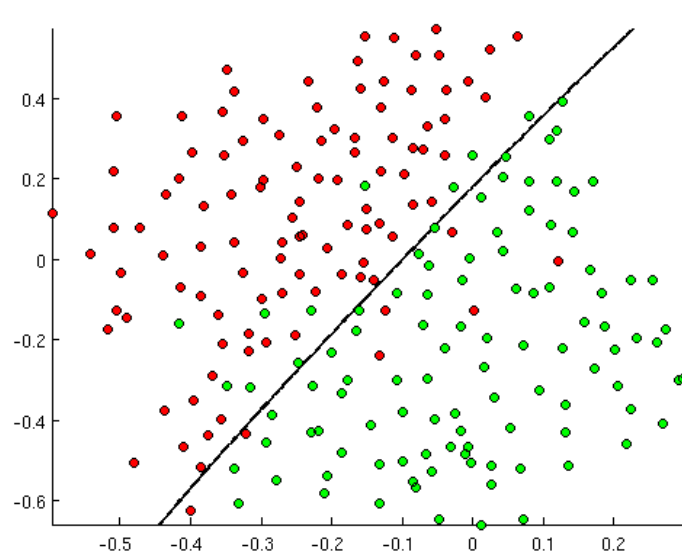
$$= \sum_{k=0}^{\infty} \sqrt{\frac{(2\alpha)^k}{k!}} \exp\left(-\alpha \mathrm{x}_n^2\right)(\mathrm{x_n})^k \sqrt{\frac{(2\alpha)^k}{k!}} \exp\left(-\alpha \mathrm{x}_m^2\right)(\mathrm{x_m})^k$$

$$= \phi(\mathrm{x_n})\phi(\mathrm{x_m})$$

Mapping to infinite-dimension !

$\alpha = 1$

$\alpha = 10$

$\alpha = 100$

$\alpha = 1000$

http://openclassroom.stanford.edu/MainFolder/DocumentPage.php?course=MachineLearning&doc=exercises/ex8/ex8.html

❑ PCA and SVM are probably the most representative conventional machine learning algorithms.

❑ PCA helps you to manipulate a set of data in a way that

- determining which features are important,

- reducing its dimension, so that the data can be processed or visualized more efficiently.

❑ SVM is a classification method founded on well defined mathematical framework, which can handle linear or nonlinear classification problems.