



# Practical Machine Learning

## Workshop 2

### Principal Components Analysis (PCA) & Support Vector Machine (SVM)

Dr. Suyong Eum



# **Principal Component Analysis (PCA)**

# Principal Component Analysis (PCA): definition

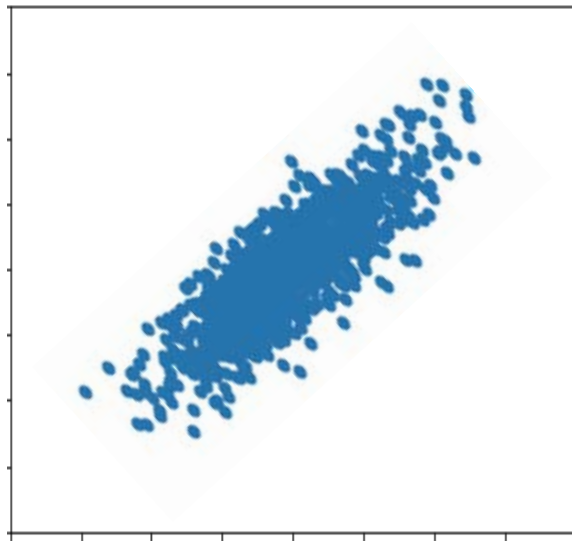
A statistical procedure that uses an orthogonal transformation to convert a set of observations of possibly correlated variables into a set of values of linearly uncorrelated variables called principal components.

*In Wikipedia:*

# Principal Component Analysis (PCA): definition

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:*

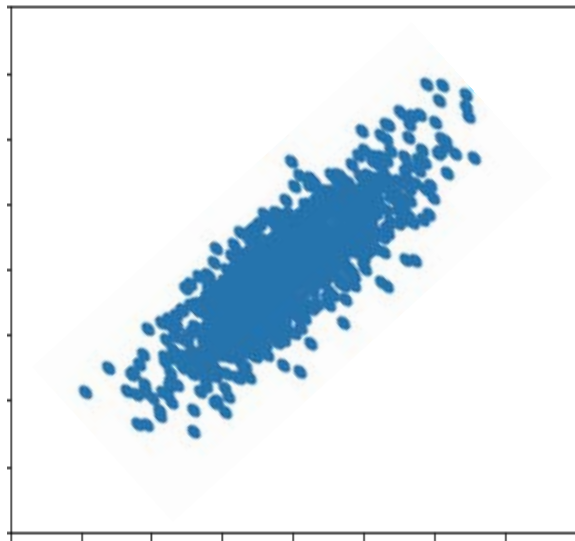


$$\begin{bmatrix} 8 & 1 \\ 1 & 8 \end{bmatrix}$$

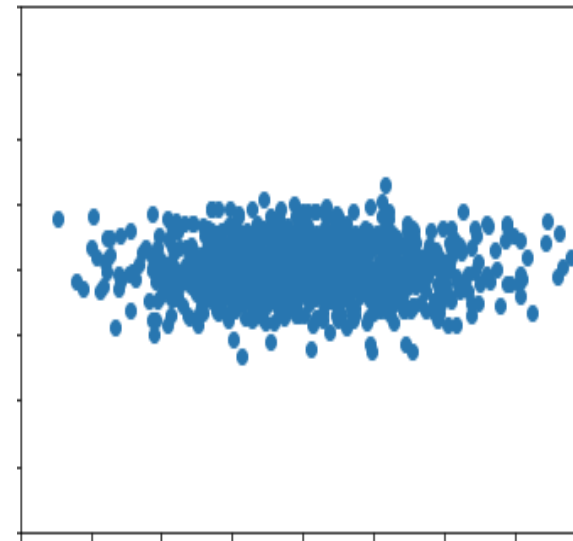
# Principal Component Analysis (PCA): definition

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:*



$$\begin{bmatrix} 8 & 1 \\ 1 & 8 \end{bmatrix}$$

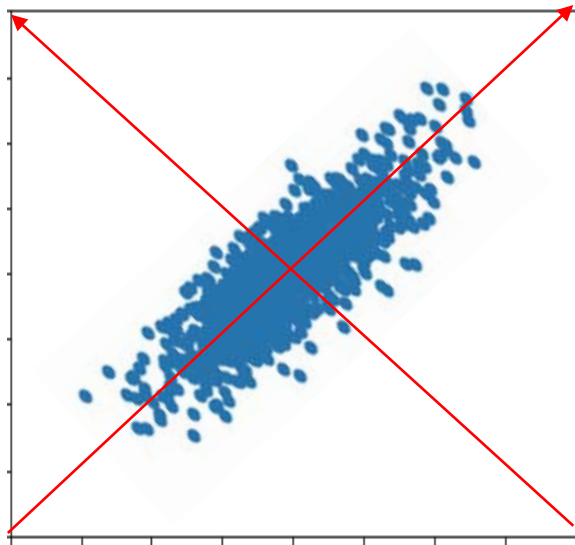


$$\begin{bmatrix} 10 & 0 \\ 0 & 1 \end{bmatrix}$$

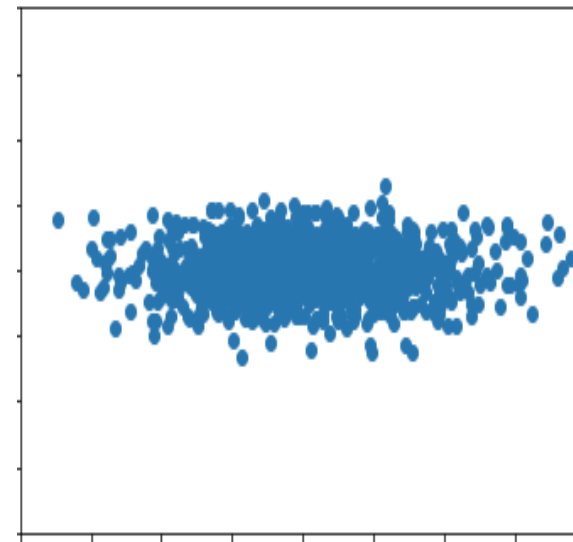
# Principal Component Analysis (PCA): definition

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:*



$$\begin{bmatrix} 8 & 1 \\ 1 & 8 \end{bmatrix}$$

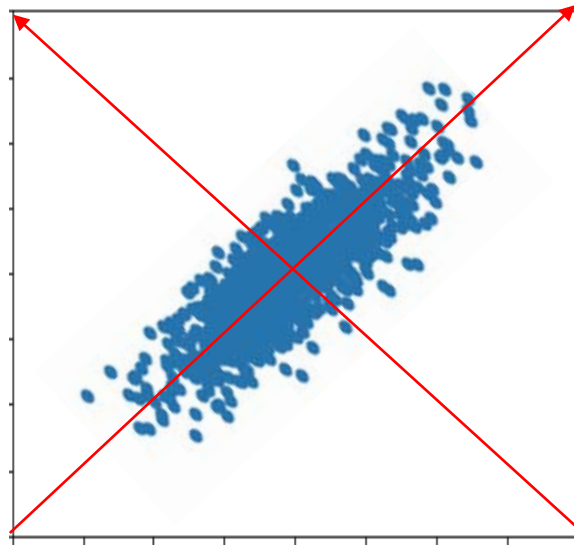


$$\begin{bmatrix} 10 & 0 \\ 0 & 1 \end{bmatrix}$$

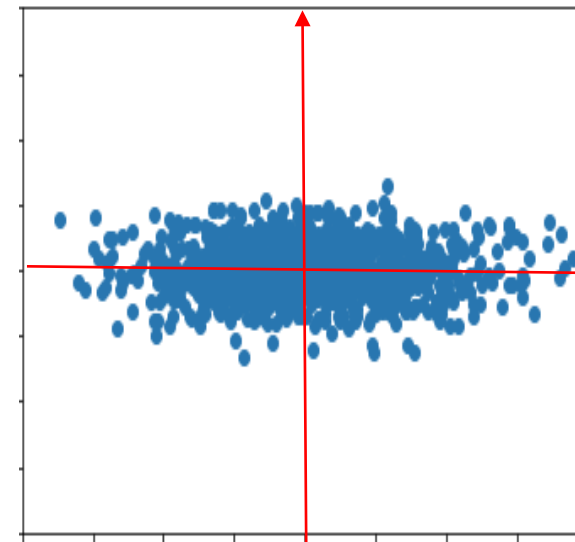
# Principal Component Analysis (PCA): definition

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:*

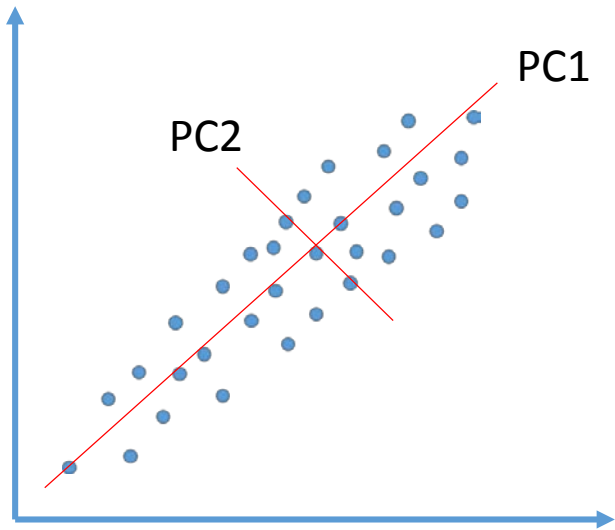


$$\begin{bmatrix} 8 & 1 \\ 1 & 8 \end{bmatrix}$$

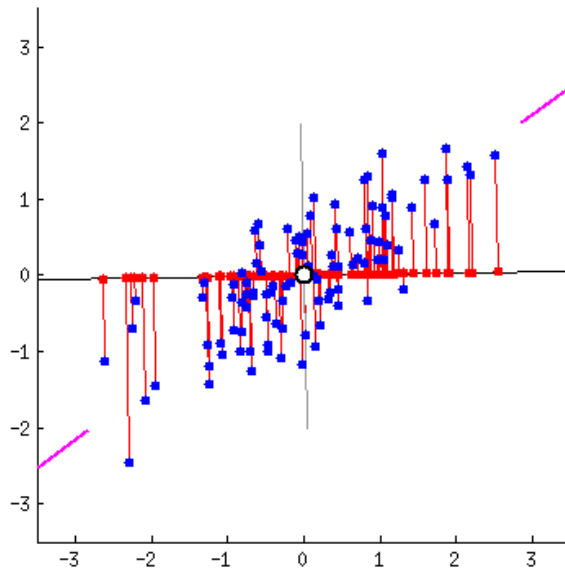


$$\begin{bmatrix} 10 & 0 \\ 0 & 1 \end{bmatrix}$$

# Principal Component Analysis (PCA): intuition

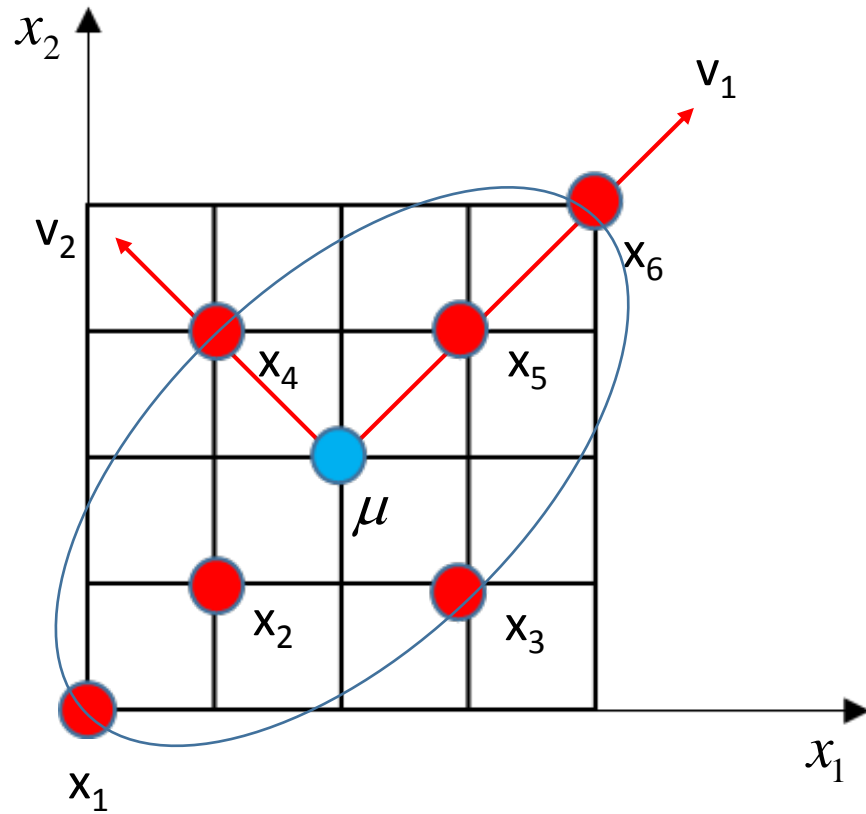


- ❑ How to select a principal component?
  - One that captures the largest variance of the data points.
- ❑ Why?
  - Because we want to clearly see how each data point is related (close) each other.
  - Then, which one (PC1 or PC2) is better?





# How to find the principal components showing the largest variance?

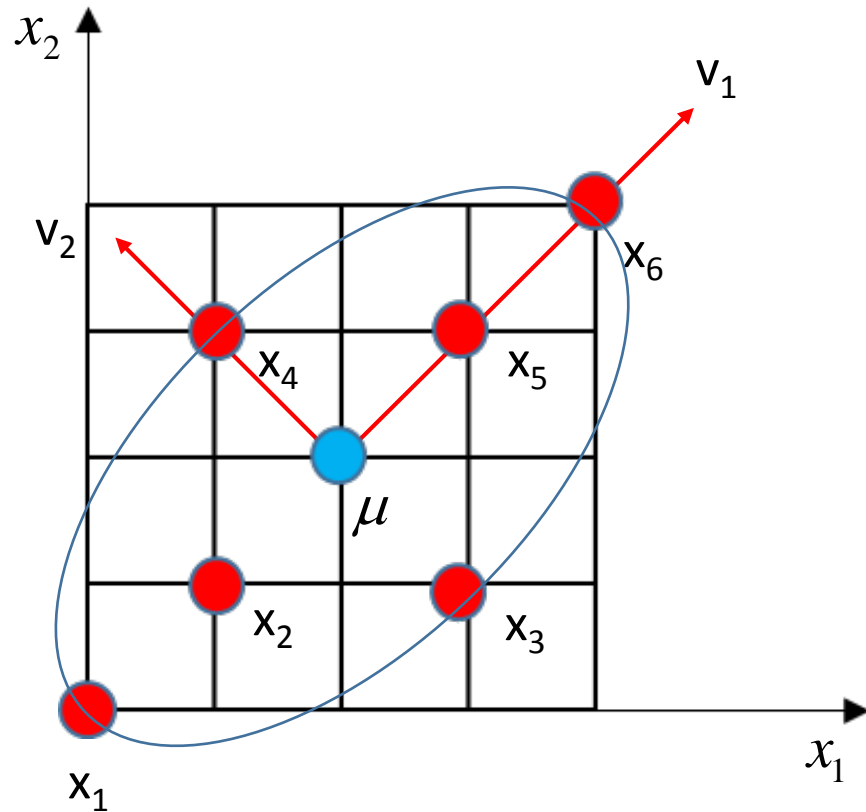


- Distance to data points from the mean along the axis of " $v_1$ "  
 $= [-2\sqrt{2}, -\sqrt{2}, 0, 0, \sqrt{2}, 2\sqrt{2}]$       variance = 4

$$\text{var} = \frac{\sum (x - \mu)^2}{N - 1}$$

$$X = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \\ x_6 \end{bmatrix} = \begin{bmatrix} -2 & -2 \\ -1 & -1 \\ 1 & -1 \\ -1 & 1 \\ 1 & 1 \\ 2 & 2 \end{bmatrix}$$

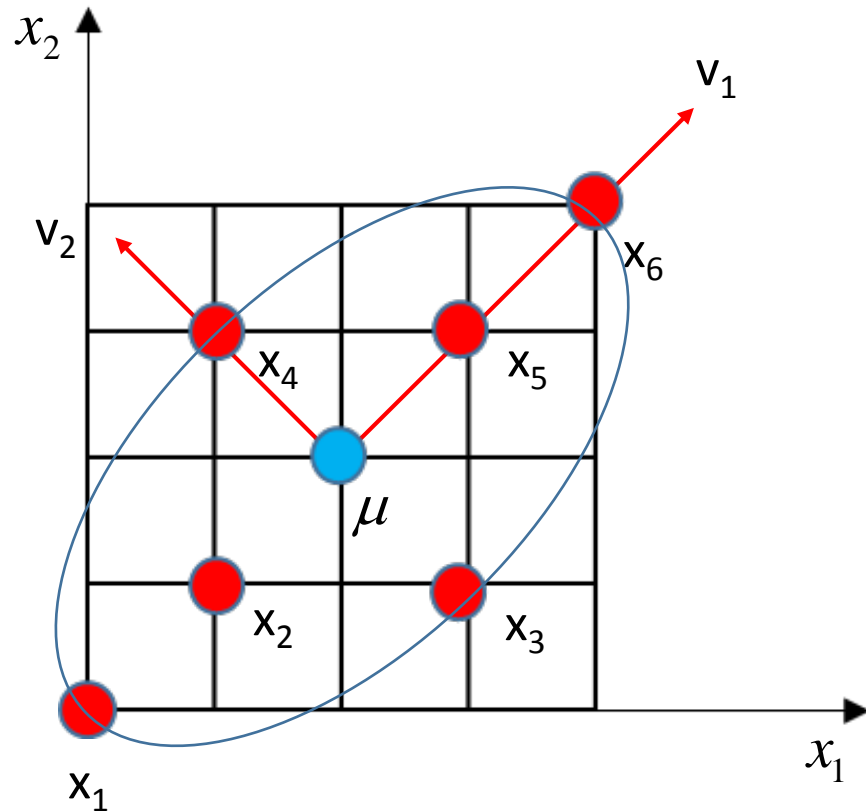
# How to find the principal components showing the largest variance?



$$\mathbf{X} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \\ x_6 \end{bmatrix} = \begin{bmatrix} -2 & -2 \\ -1 & -1 \\ 1 & -1 \\ -1 & 1 \\ 1 & 1 \\ 2 & 2 \end{bmatrix}$$

- Distance to data points from the mean along the axis of "v<sub>1</sub>"  
 $= [-2\sqrt{2}, -\sqrt{2}, 0, 0, \sqrt{2}, 2\sqrt{2}]$  variance = 4
- Distance to data points from the mean along the axis of "v<sub>2</sub>"  
 $= [0, 0, -\sqrt{2}, \sqrt{2}, 0, 0]$  variance = 0.8

# How to find the principal components showing the largest variance?



$$X = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \\ x_6 \end{bmatrix} = \begin{bmatrix} -2 & -2 \\ -1 & -1 \\ 1 & -1 \\ -1 & 1 \\ 1 & 1 \\ 2 & 2 \end{bmatrix}$$

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

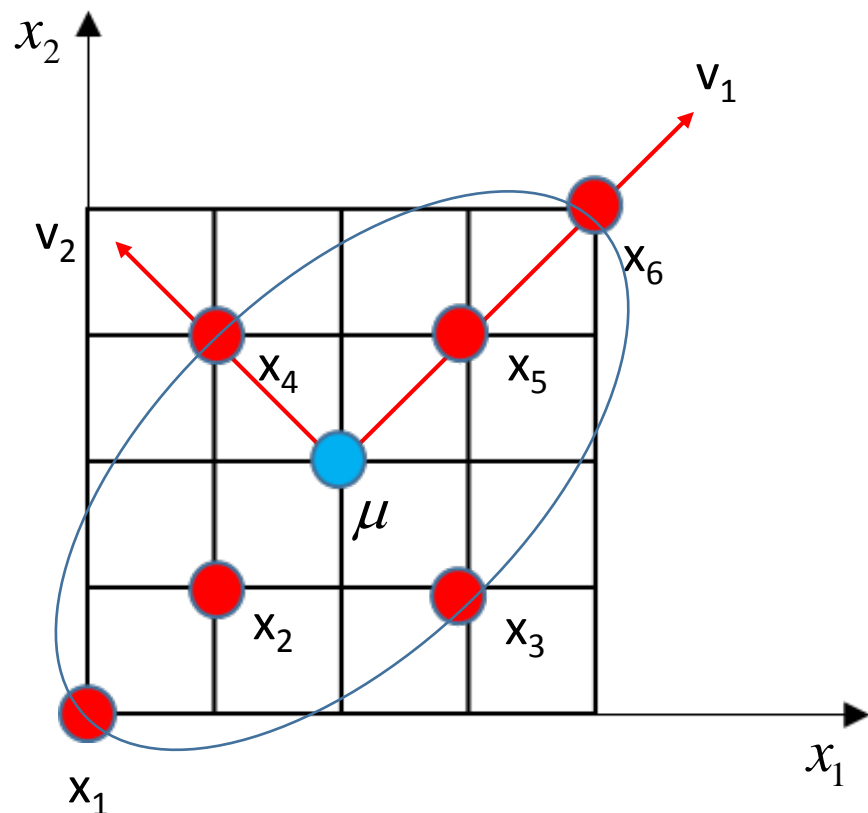
Distance to data points from the mean along the axis of " $v_2$ "  
 $= [0, 0, -\sqrt{2}, \sqrt{2}, 0, 0]$       variance = 0.8

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





# How to find the principal components showing the largest variance?



$$X = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \\ x_6 \end{bmatrix} = \begin{bmatrix} -2 & -2 \\ -1 & -1 \\ 1 & -1 \\ -1 & 1 \\ 1 & 1 \\ 2 & 2 \end{bmatrix}$$

Distance to data points from the mean along the axis of "v<sub>1</sub>"  
 $= [-2\sqrt{2}, -\sqrt{2}, 0, 0, \sqrt{2}, 2\sqrt{2}]$  variance = 4

Distance to data points from the mean along the axis of "v<sub>2</sub>"  
 $= [0, 0, -\sqrt{2}, \sqrt{2}, 0, 0]$  variance = 0.8

$\text{cov}(X) = \begin{bmatrix} 2.4 & 1.6 \\ 1.6 & 2.4 \end{bmatrix}$  variance along the axis of "x<sub>1</sub>"  
 variance along the axis of "x<sub>2</sub>"

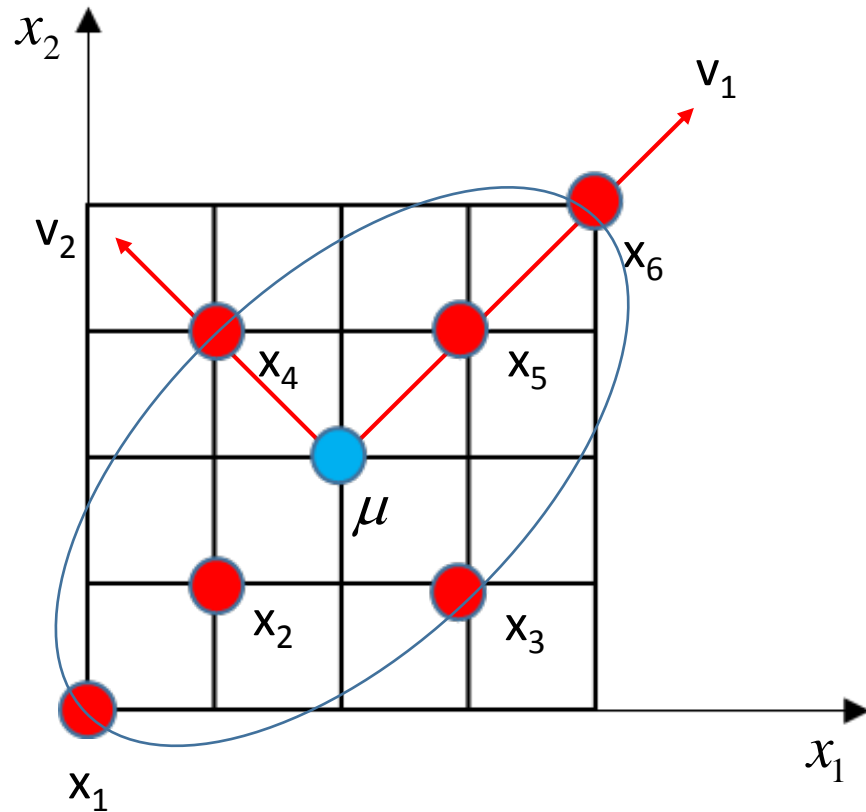
$\text{cov}(x) = V\Lambda V^T$

```
>> [vec, val] = eig(cov(x))
vec =
    -0.70711    0.70711
     0.70711    0.70711
val =
Diagonal Matrix
    0.80000    0
         0    4.00000
```

V

Λ

# How to find the principal components showing the largest variance?



$$X = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \\ x_6 \end{bmatrix} = \begin{bmatrix} -2 & -2 \\ -1 & -1 \\ 1 & -1 \\ -1 & 1 \\ 1 & 1 \\ 2 & 2 \end{bmatrix}$$

Distance to data points from the mean along the axis of "v<sub>1</sub>"  
 $= [-2\sqrt{2}, -\sqrt{2}, 0, 0, \sqrt{2}, 2\sqrt{2}]$  variance = 4

Distance to data points from the mean along the axis of "v<sub>2</sub>"  
 $= [0, 0, -\sqrt{2}, \sqrt{2}, 0, 0]$  variance = 0.8

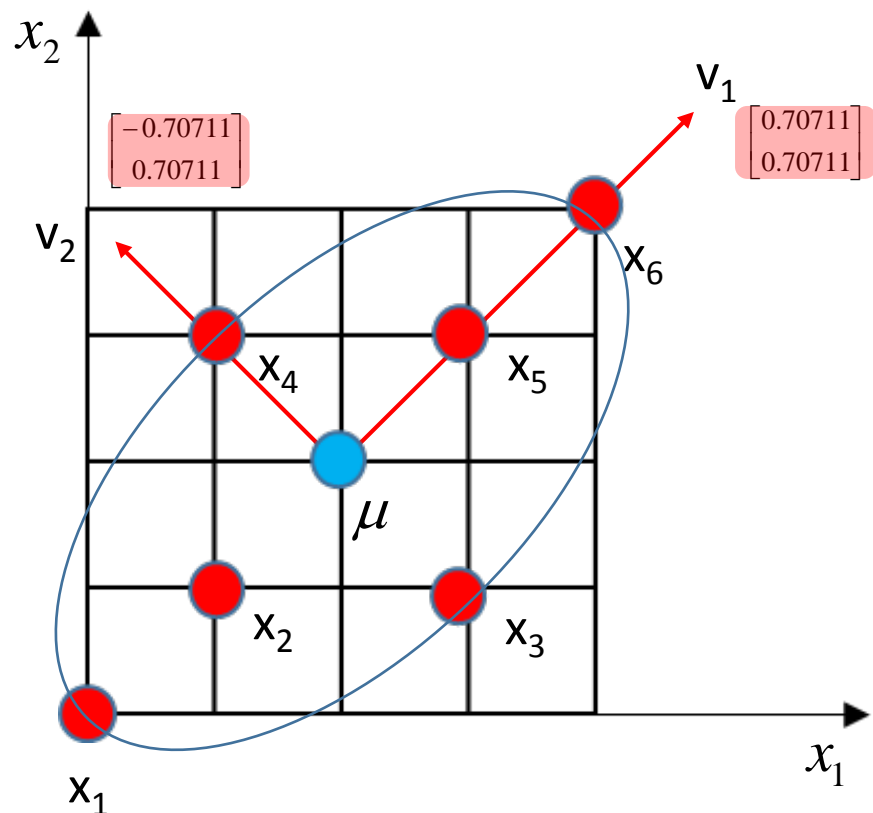
$\text{cov}(X) = \begin{bmatrix} 2.4 & 1.6 \\ 1.6 & 2.4 \end{bmatrix}$  variance along the axis of "x<sub>1</sub>"  
 variance along the axis of "x<sub>2</sub>"

$\text{cov}(x) = V\Lambda V^T$

```
>> [vec, val] = eig(cov(x))
vec =
    -0.70711    0.70711
     0.70711    0.70711
val =
Diagonal Matrix
    0.80000    0
         0    4.00000
```

V  
 Λ

# How to find the principal components showing the largest variance?



$$X = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \\ x_6 \end{bmatrix} = \begin{bmatrix} -2 & -2 \\ -1 & -1 \\ 1 & -1 \\ -1 & 1 \\ 1 & 1 \\ 2 & 2 \end{bmatrix}$$

Distance to data points from the mean along the axis of "v<sub>1</sub>"  
 $= [-2\sqrt{2}, -\sqrt{2}, 0, 0, \sqrt{2}, 2\sqrt{2}]$  variance = 4

Distance to data points from the mean along the axis of "v<sub>2</sub>"  
 $= [0, 0, -\sqrt{2}, \sqrt{2}, 0, 0]$  variance = 0.8

$\text{cov}(X) = \begin{bmatrix} 2.4 & 1.6 \\ 1.6 & 2.4 \end{bmatrix}$   
 variance along the axis of "x<sub>1</sub>"  
 variance along the axis of "x<sub>2</sub>"

$\text{cov}(x) = V\Lambda V^T$

```
>> [vec, val] = eig(cov(x))
vec =
    -0.70711    0.70711
     0.70711    0.70711
val =
Diagonal Matrix
    0.80000    0
           0    4.00000
```

V  
 Λ



# How to find the principal components showing the largest variance?

1) Find the covariance matrix of data points.

```
>> 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
```

# How to find the principal components showing the largest variance?

- 1) Find the covariance matrix of data points.
- 2) Obtain the eigen values and vectors of the covariance matrix: **eigen decomposition**.

```
>> 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

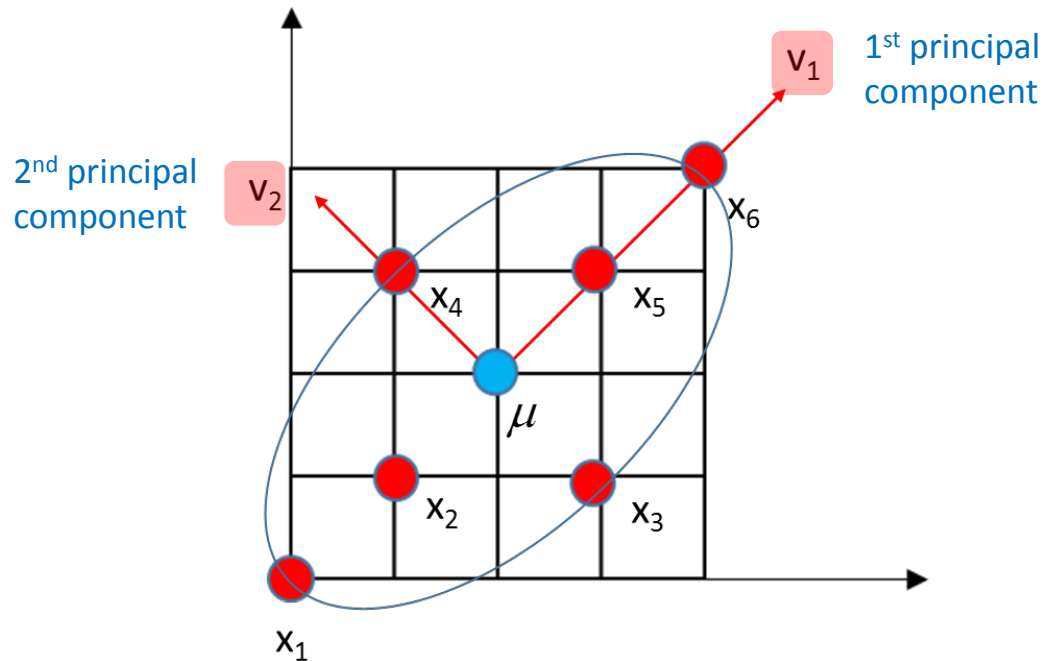
val =

Diagonal Matrix

    0.80000    0
         0    4.00000
```

# How to find the principal components showing the largest variance?

- 1) Find the covariance matrix of data points.
- 2) Obtain the eigen values and vectors of the covariance matrix: **eigen 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.



```
>> 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 =

    0.80000    0
           0    4.00000

Diagonal Matrix
```

2<sup>nd</sup> principal component      1<sup>st</sup> principal component

# How to find the principal components showing the largest variance?

- ❑ Actually, there is a more convenient way of doing it.
- ❑ It is called “Singular Value Decomposition” or **SVD**.

Eigen Value decomposition

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

```
>> 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
```

# How to find the principal components showing the largest variance?

- ❑ Actually, there is a more convenient way of doing it.
- ❑ It is called “Singular Value Decomposition” or **SVD**.

Eigen Value 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  ?
```

# How to find the principal components showing the largest variance?

- ❑ Actually, there is a more convenient way of doing it.
- ❑ It is called “Singular Value Decomposition” or **SVD**.

Eigen Value decomposition

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

```
>> 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
```

Singular Value Decomposition (SVD)

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

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

# How to find the principal components showing the largest variance?

- ❑ Actually, there is a more convenient way of doing it.
- ❑ It is called “Singular Value Decomposition” or **SVD**.

Eigen Value decomposition

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

```
>> 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
```

Singular Value Decomposition (SVD)

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

$$\begin{aligned} X^T X &= (U \Sigma V^T)^T (U \Sigma V^T) \\ &= V \Sigma^T U^T U \Sigma V^T \\ &= V \Sigma^2 V^T \end{aligned}$$

# How to find the principal components showing the largest variance?

- ❑ Actually, there is a more convenient way of doing it.
- ❑ It is called “Singular Value Decomposition” or **SVD**.

Eigen decomposition

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

```
>> 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
```

Singular Value Decomposition (SVD)

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

$$\begin{aligned} X^T X &= (U \Sigma V^T)^T (U \Sigma V^T) \\ &= V \Sigma^T U^T U \Sigma V^T \\ &= V \Sigma^2 V^T \end{aligned} \quad \Lambda = \Sigma^2$$

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

Eigen value

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

Singular value



# How to find the principal components showing the largest variance?

- ❑ Actually, there is a more convenient way of doing it.
- ❑ It 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 =

    0.80000    0
         0    4.00000

Diagonal Matrix
```

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

   -0.70711    0.70711
    0.70711    0.70711

val =

    4.0000    0
         0   20.0000

Diagonal Matrix
```

$$\Lambda = \Sigma^2$$

$$4.4721^2 = 20$$

```
>> [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 =

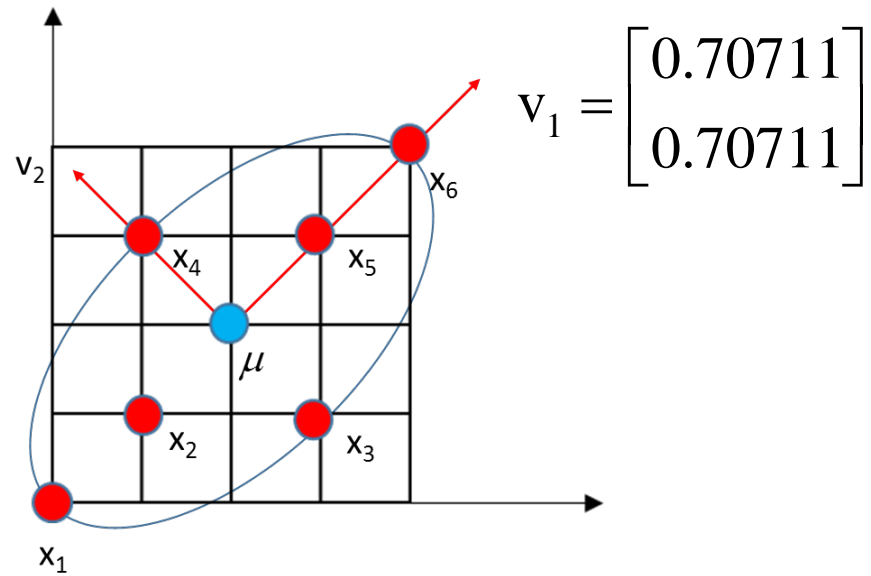
    4.4721    0
         0    2.0000
         0    0
         0    0
         0    0
         0    0

Diagonal Matrix

V =

    0.70711   -0.70711
    0.70711    0.70711
```

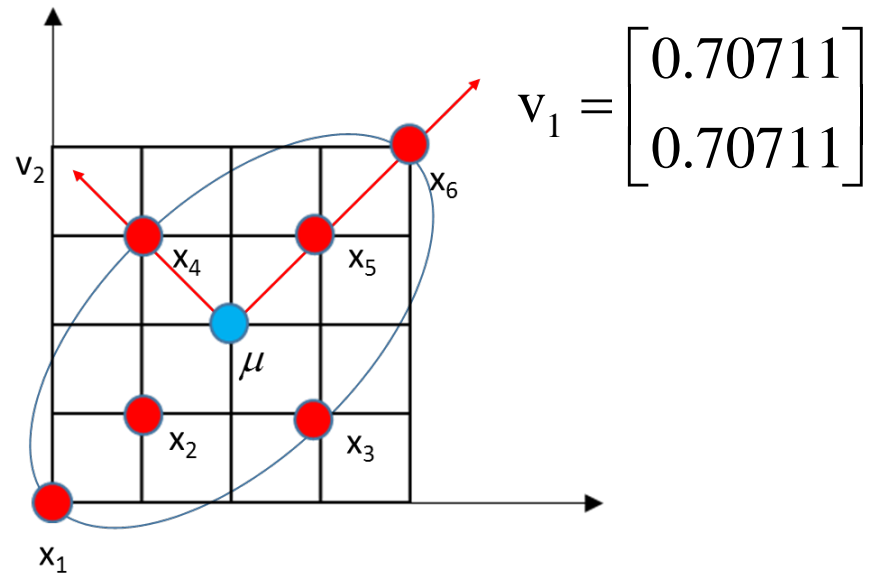
# Now we know how to find the principal components



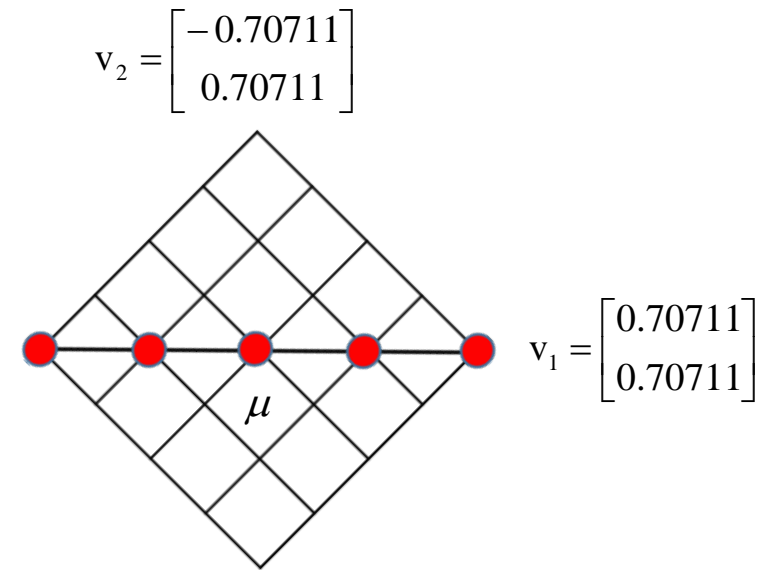
$$\mathbf{X} = \begin{bmatrix} -2 & -2 \\ -1 & -1 \\ 1 & -1 \\ -1 & 1 \\ 1 & 1 \\ 2 & 2 \end{bmatrix}$$

2 dimension data points can be represented into one dimension space ( $v_1$ )

# Dimension reduction



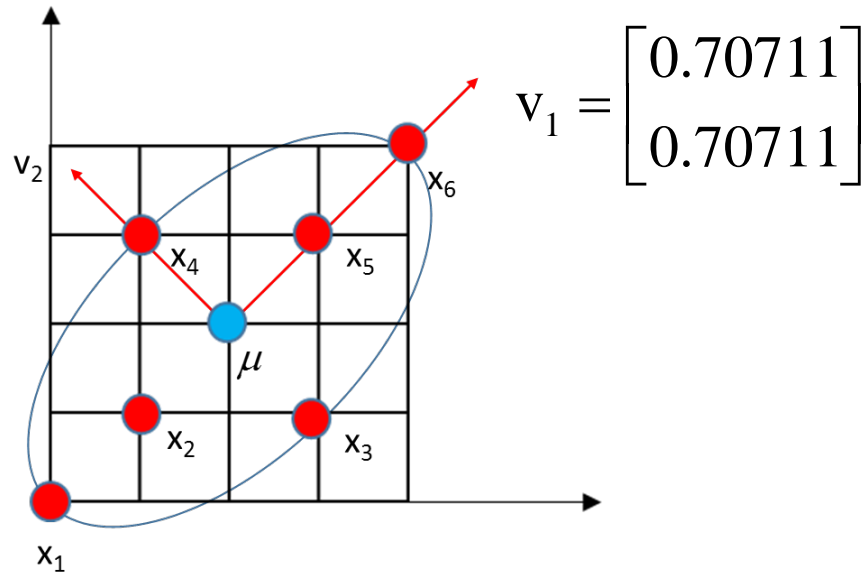
$$X = \begin{bmatrix} -2 & -2 \\ -1 & -1 \\ 1 & -1 \\ -1 & 1 \\ 1 & 1 \\ 2 & 2 \end{bmatrix}$$



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

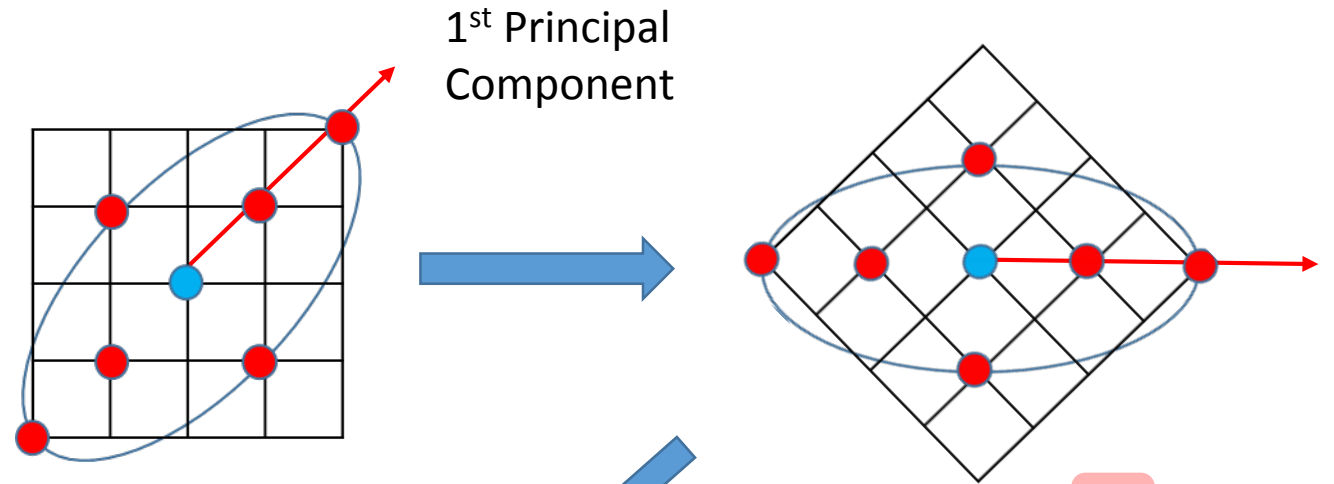
2 dimension data points can be represented into one dimension space ( $v_1$ )

# Dimension reduction

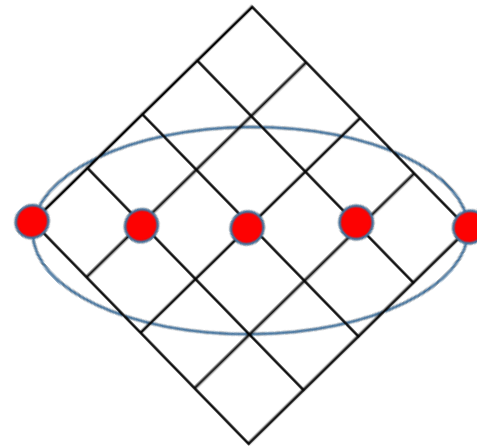


$$X = \begin{bmatrix} -2 & -2 \\ -1 & -1 \\ 1 & -1 \\ -1 & 1 \\ 1 & 1 \\ 2 & 2 \end{bmatrix}$$

2 dimension data points can be represented into one dimension space ( $v_1$ )

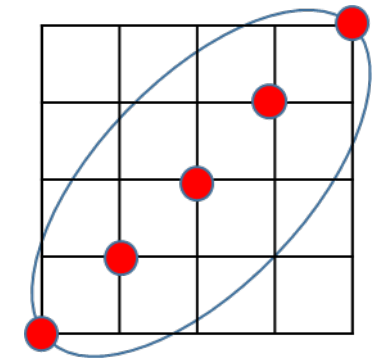


$$X_{rot} = X \cdot V$$



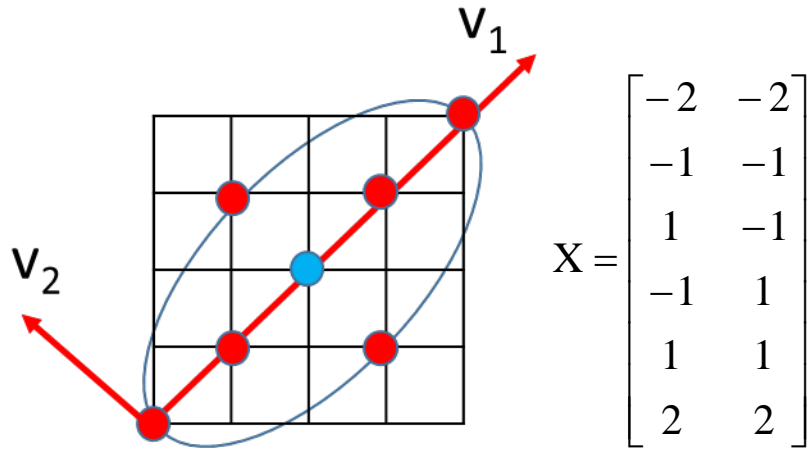
Set the " $v_2$ " into zero

$$X_{rot\_zero} = X_{rot} \cdot V^{-1}$$



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

# Example



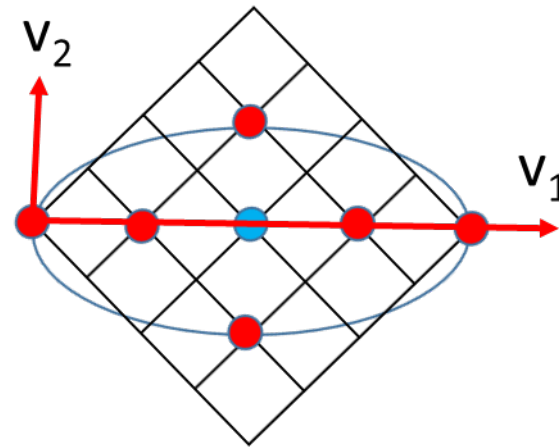
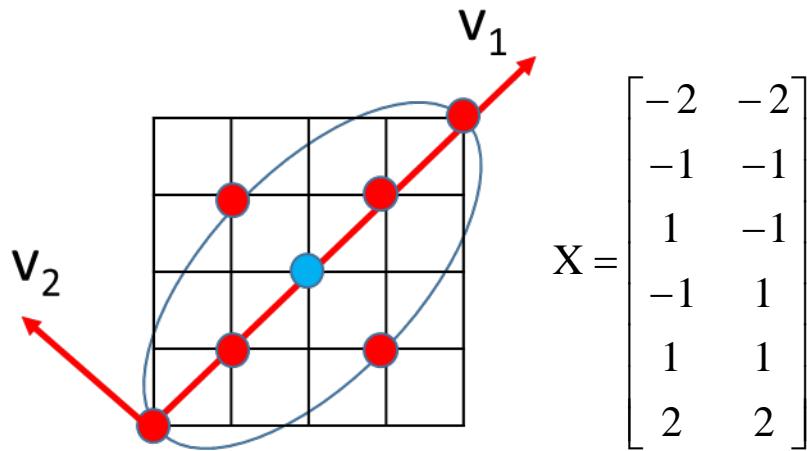
$$X = U\Sigma V^T$$

```
>> [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
```

# Example



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

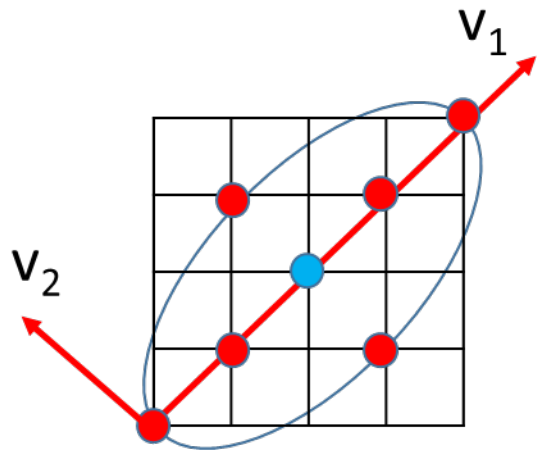
$$X = U\Sigma V^T$$

```
>> [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

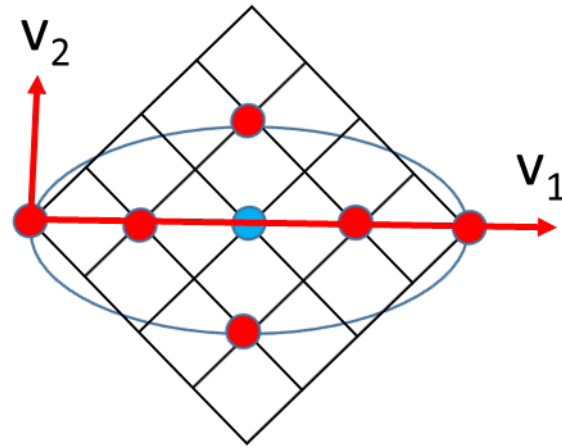
S =
Diagonal Matrix
4.4721    0
0    2.0000
0    0
0    0
0    0
0    0

V =
0.70711   -0.70711
0.70711    0.70711
```

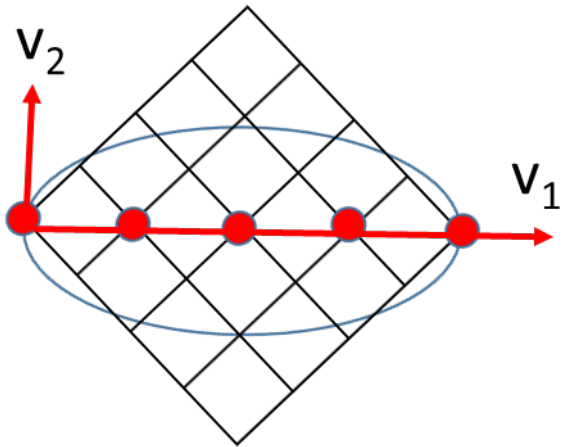
# Example



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



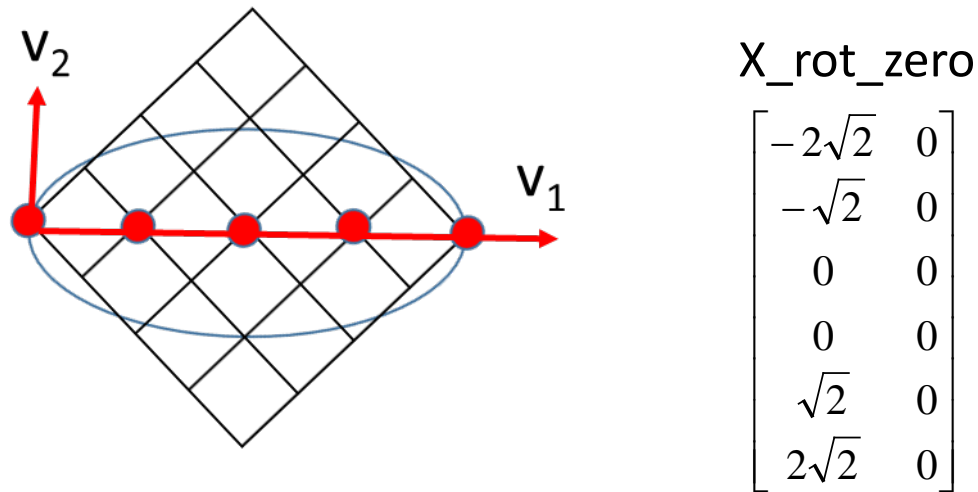
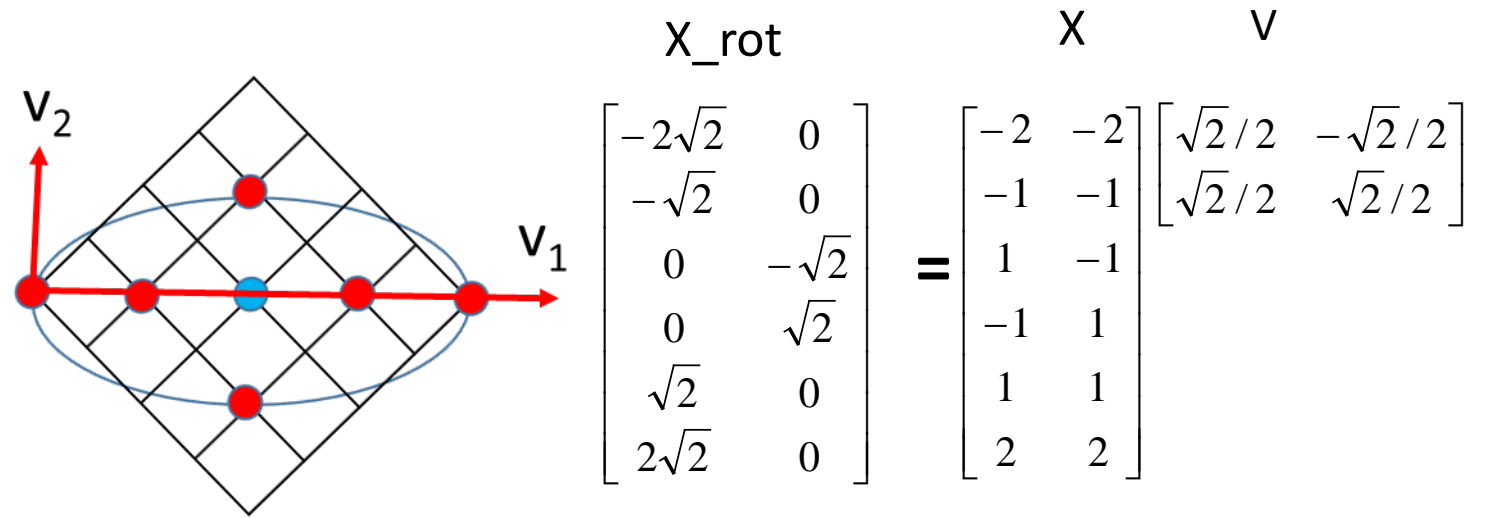
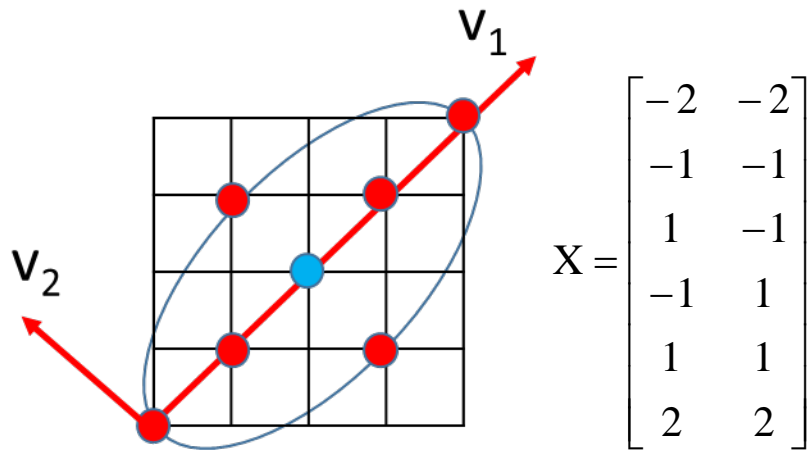
$$\begin{matrix} X_{\text{rot}} & X & V \\ \begin{bmatrix} -2\sqrt{2} & 0 \\ -\sqrt{2} & 0 \\ 0 & -\sqrt{2} \\ 0 & \sqrt{2} \\ \sqrt{2} & 0 \\ 2\sqrt{2} & 0 \end{bmatrix} & = & \begin{bmatrix} -2 & -2 \\ -1 & -1 \\ 1 & -1 \\ -1 & 1 \\ 1 & 1 \\ 2 & 2 \end{bmatrix} \begin{bmatrix} \sqrt{2}/2 & -\sqrt{2}/2 \\ \sqrt{2}/2 & \sqrt{2}/2 \end{bmatrix} \end{matrix}$$



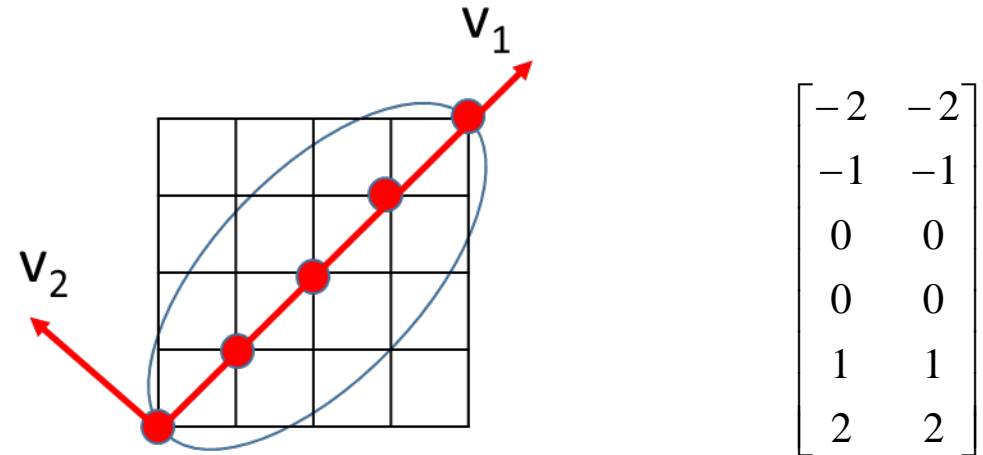
$$X_{\text{rot\_zero}} = \begin{bmatrix} -2\sqrt{2} & 0 \\ -\sqrt{2} & 0 \\ 0 & 0 \\ 0 & 0 \\ \sqrt{2} & 0 \\ 2\sqrt{2} & 0 \end{bmatrix}$$

Set the "v<sub>2</sub>" elements into zero

# Example



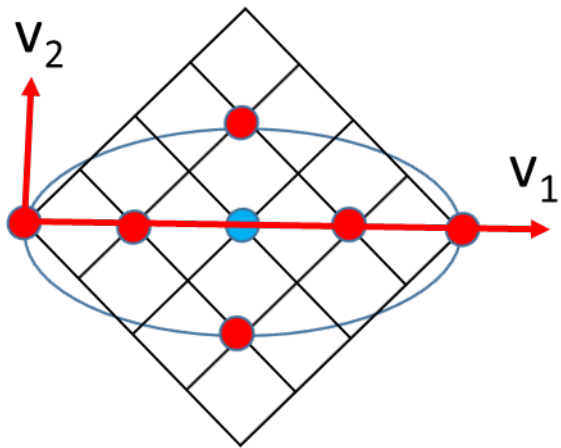
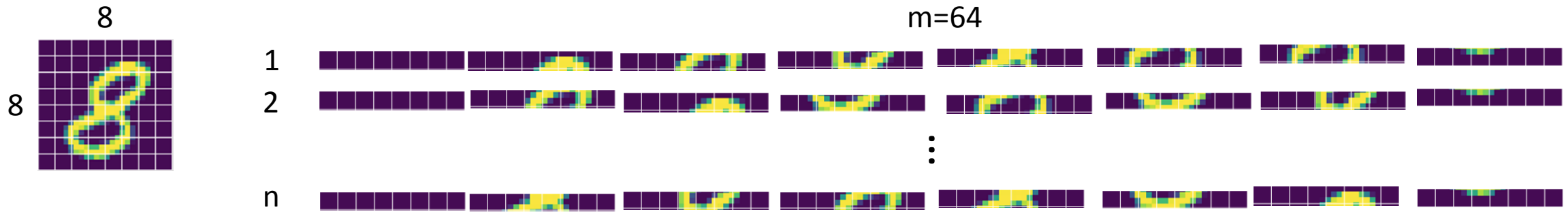
Set the " $v_2$ " elements into zero



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



# How to use PCA for machine learning?



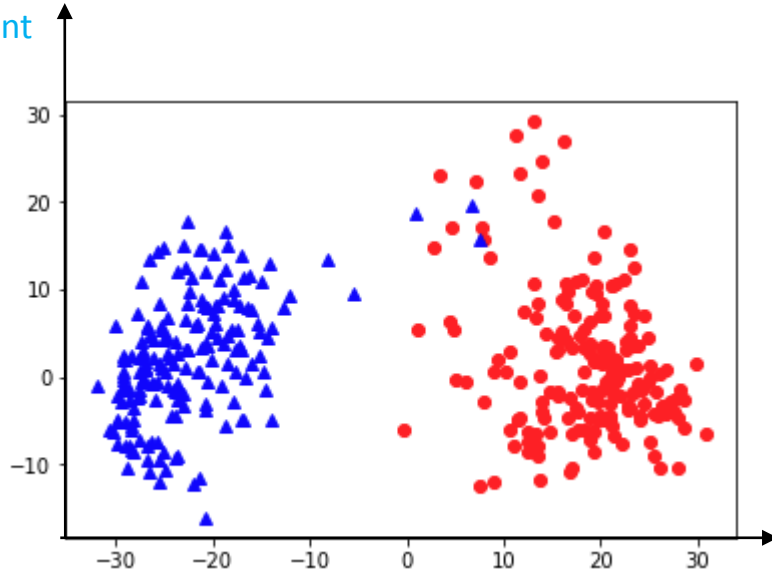
$$X_{\text{rot}} = X \cdot V$$

$$\begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1m} \\ \vdots & \vdots & \vdots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nm} \end{bmatrix}$$

A digit number with 64 dimension can be shown in 2 dimension space ( $v_1$  and  $v_2$ ).

- ❑ Each digit number has 8 by 8 = 64 dimensions.
- ❑ After SVD, the first two principal components are selected, and the data points with 64 dimension are plotted in two dimension.

2<sup>nd</sup> Principal Component

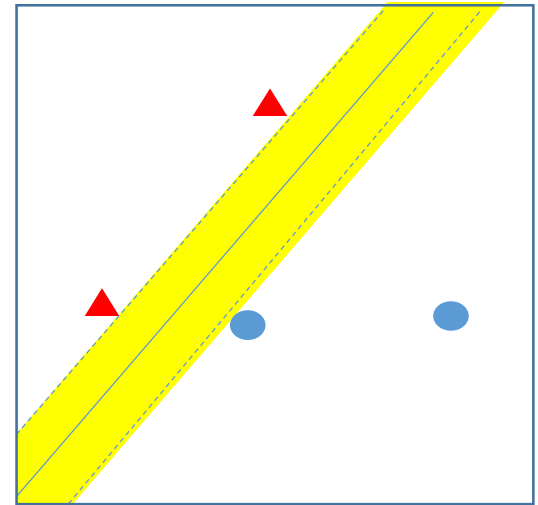
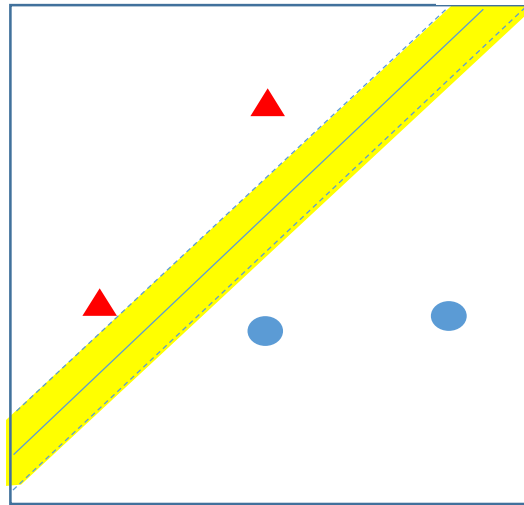
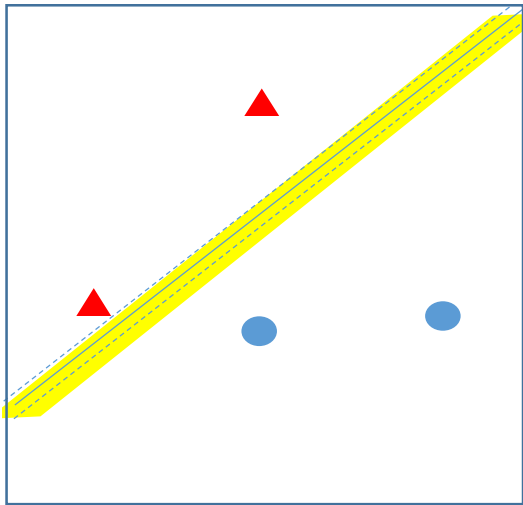


1<sup>st</sup> Principal Component

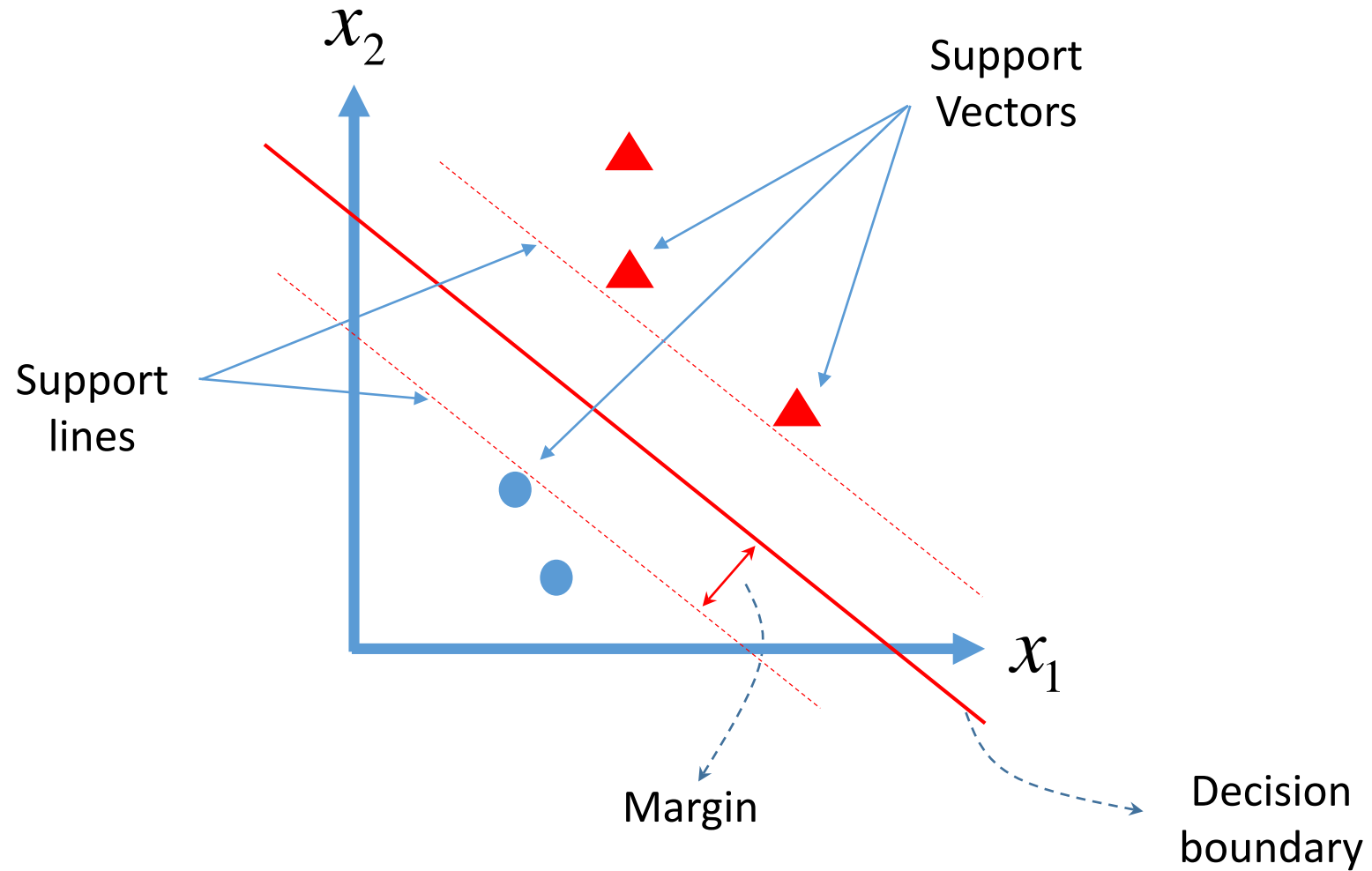


# Support Vector Machine (SVM)

36 Which one is better for classification?

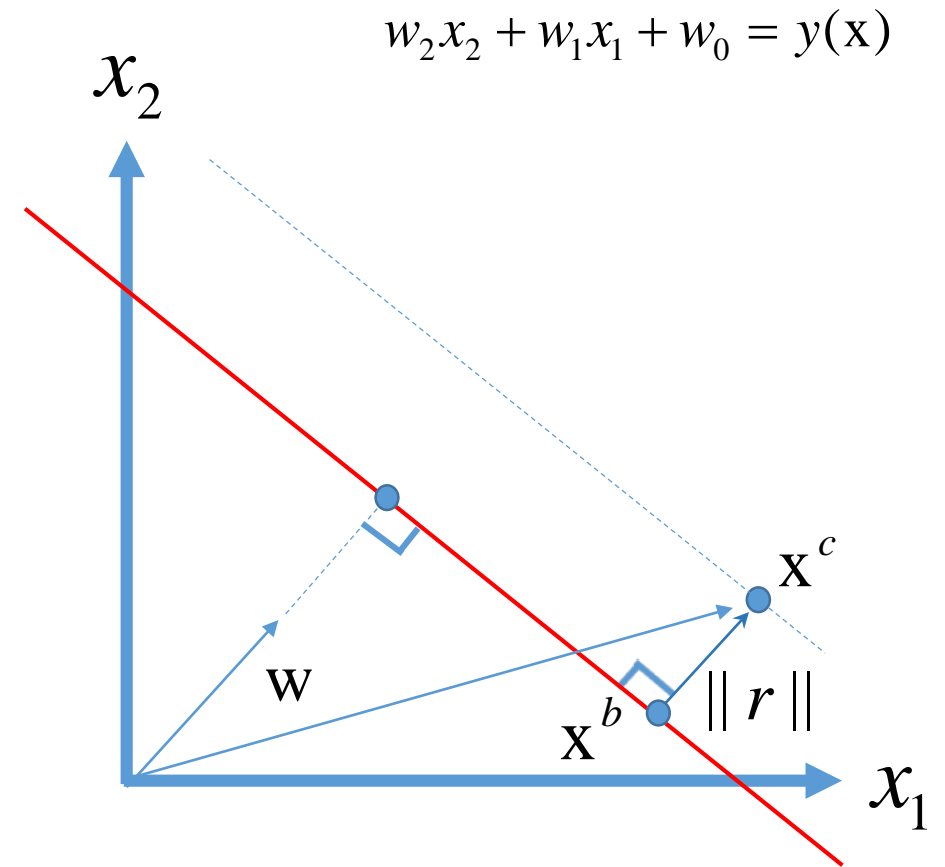


# 37 Terminology used in this lecture



# 38 Margin distance

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

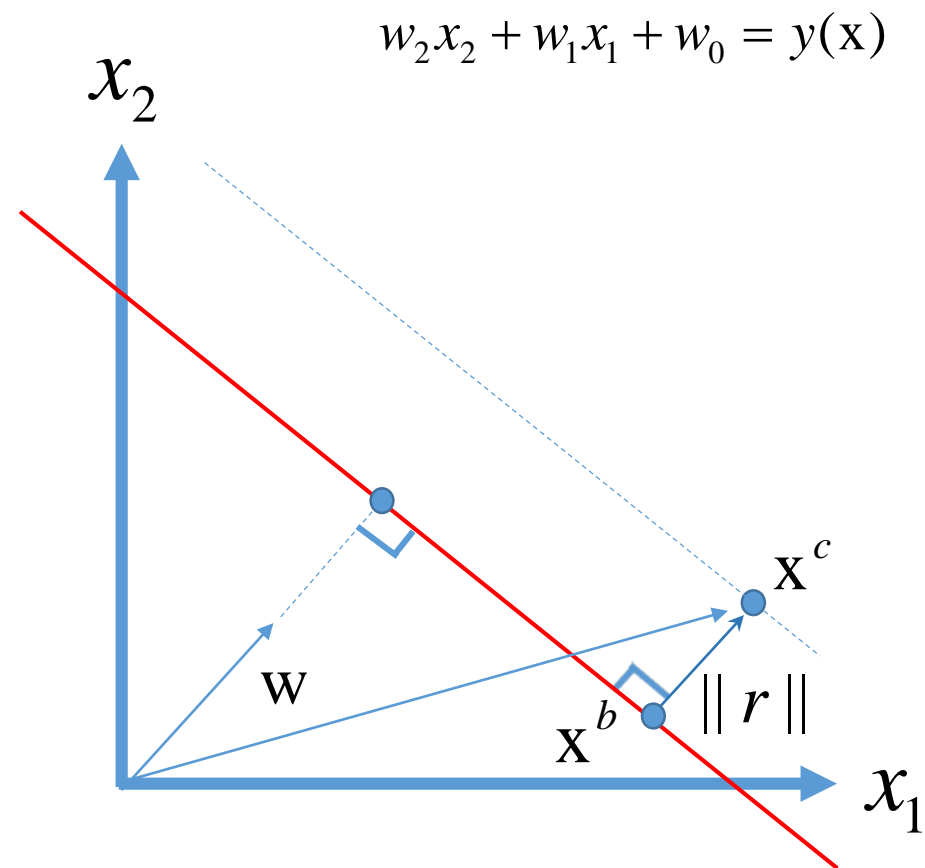


# 39 Margin distance

Size of the vector  
( $x^b \rightarrow x^c$ )

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

Unit vector showing  
the direction only



# 40 Margin distance

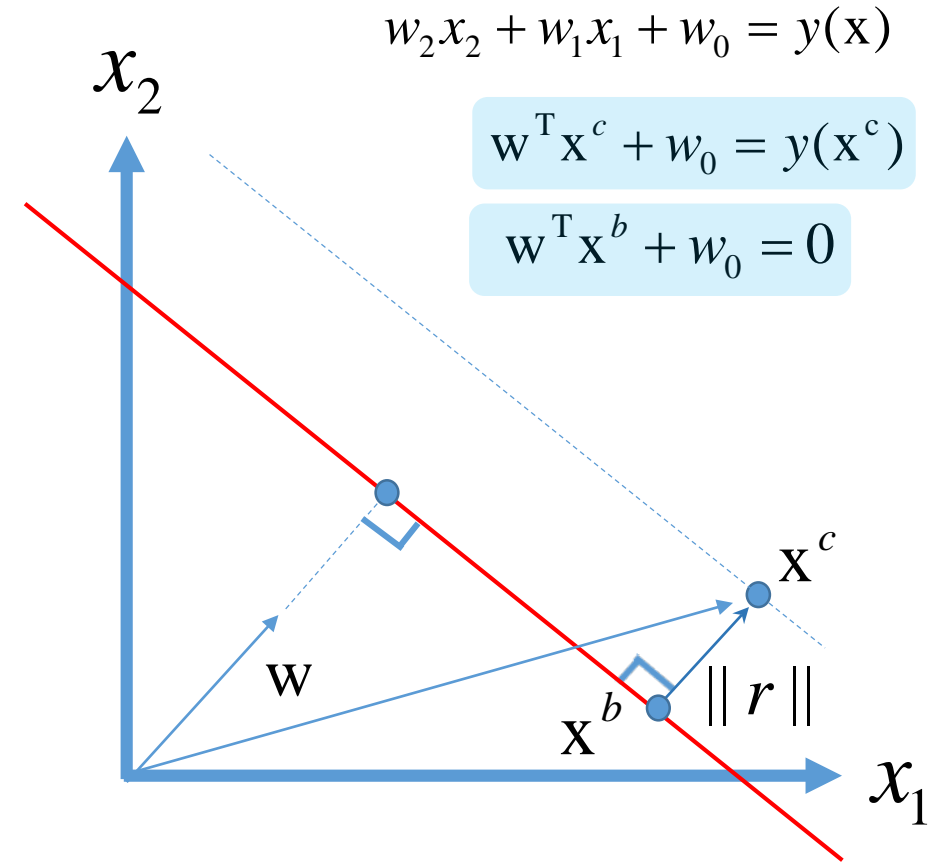
Size of the vector  
( $x^b \rightarrow x^c$ )

$$\mathbf{x}^c = \mathbf{x}^b + \underbrace{\|r\|}_{\text{Size of the vector}} \frac{\mathbf{w}}{\|\mathbf{w}\|}$$

Unit vector showing the direction only

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

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





# 41 Margin distance

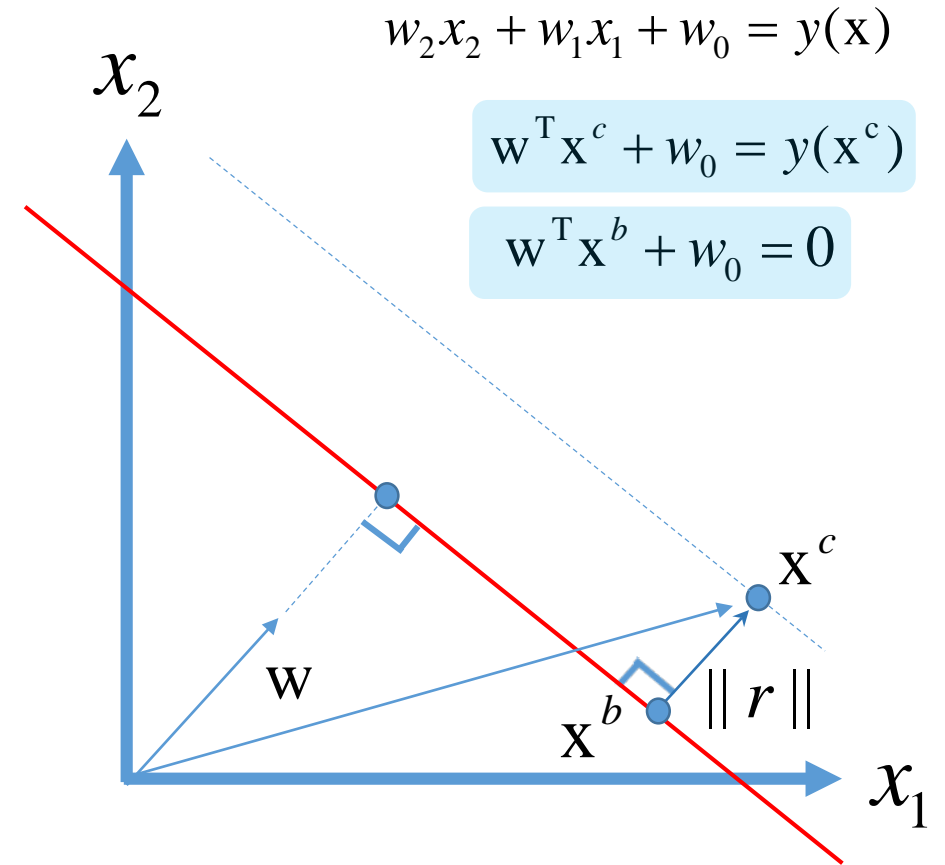
Size of the vector  
( $x^b \rightarrow x^c$ )

$$\mathbf{x}^c = \mathbf{x}^b + \underbrace{\|r\|}_{\text{Size of the vector}} \underbrace{\frac{\mathbf{w}}{\|\mathbf{w}\|}}_{\text{Unit vector showing the direction only}}$$

□ Let's multiply  $\mathbf{w}^T$  and add  $w_0$  in both sides.

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

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



# 42 Margin distance

Size of the vector  
( $x^b \rightarrow x^c$ )

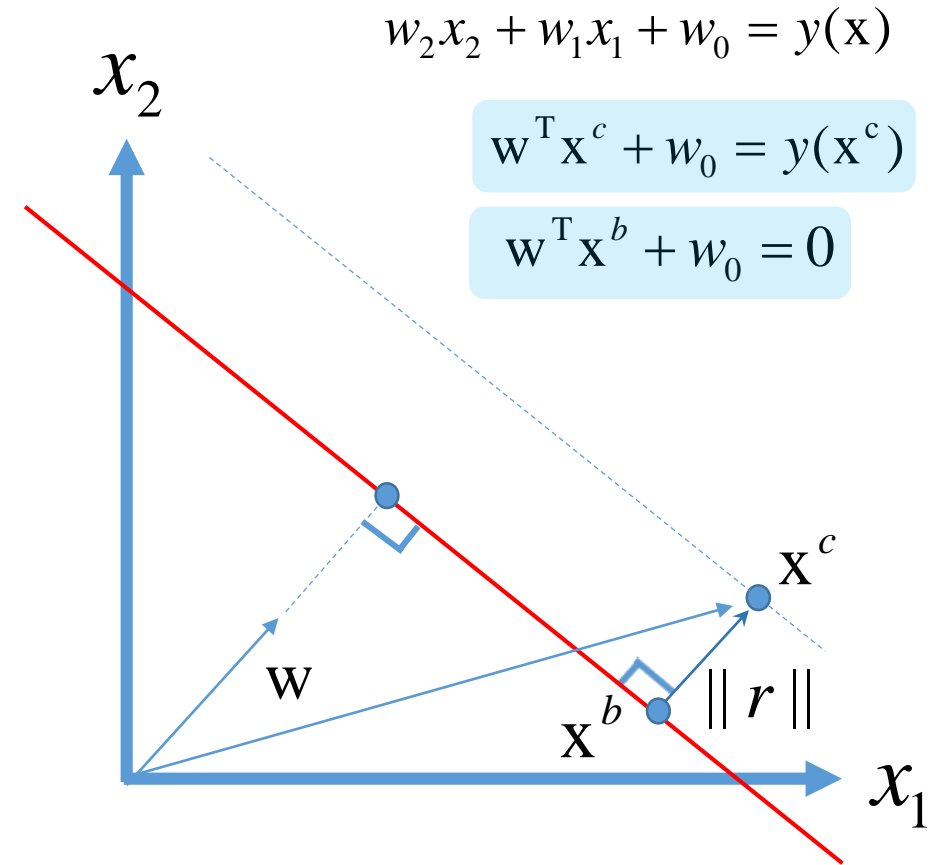
$$\mathbf{x}^c = \mathbf{x}^b + \underbrace{\|r\|}_{\text{Size of the vector}} \underbrace{\frac{\mathbf{w}}{\|\mathbf{w}\|}}_{\text{Unit vector showing the direction only}}$$

□ Let's multiply  $\mathbf{w}^T$  and add  $w_0$  in both sides.

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

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

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



# 43 Margin distance

Size of the vector  
( $x^b \rightarrow x^c$ )

$$\mathbf{x}^c = \mathbf{x}^b + \underbrace{\|r\|}_{\text{Size of the vector}} \frac{\mathbf{w}}{\|\mathbf{w}\|}$$

Unit vector showing the direction only

□ Let's multiply  $\mathbf{w}^T$  and add  $w_0$  in both sides.

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

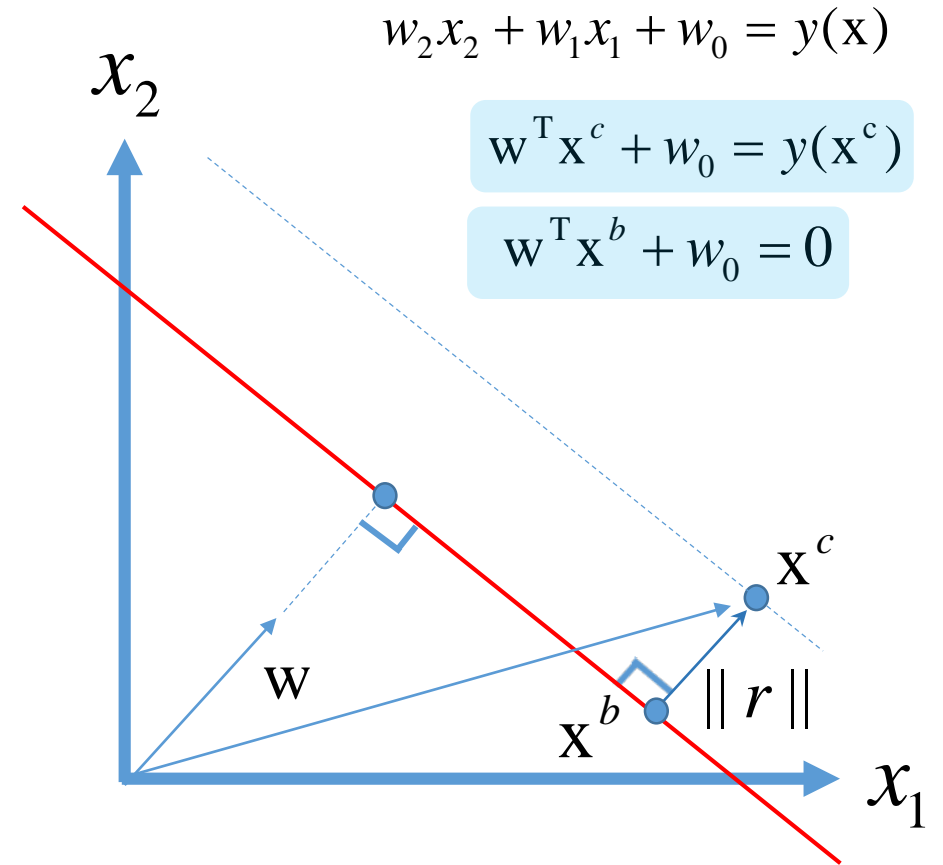
$$y(\mathbf{x}^c) = \mathbf{w}^T \|r\| \frac{\mathbf{w}}{\|\mathbf{w}\|}$$

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

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

Let's say

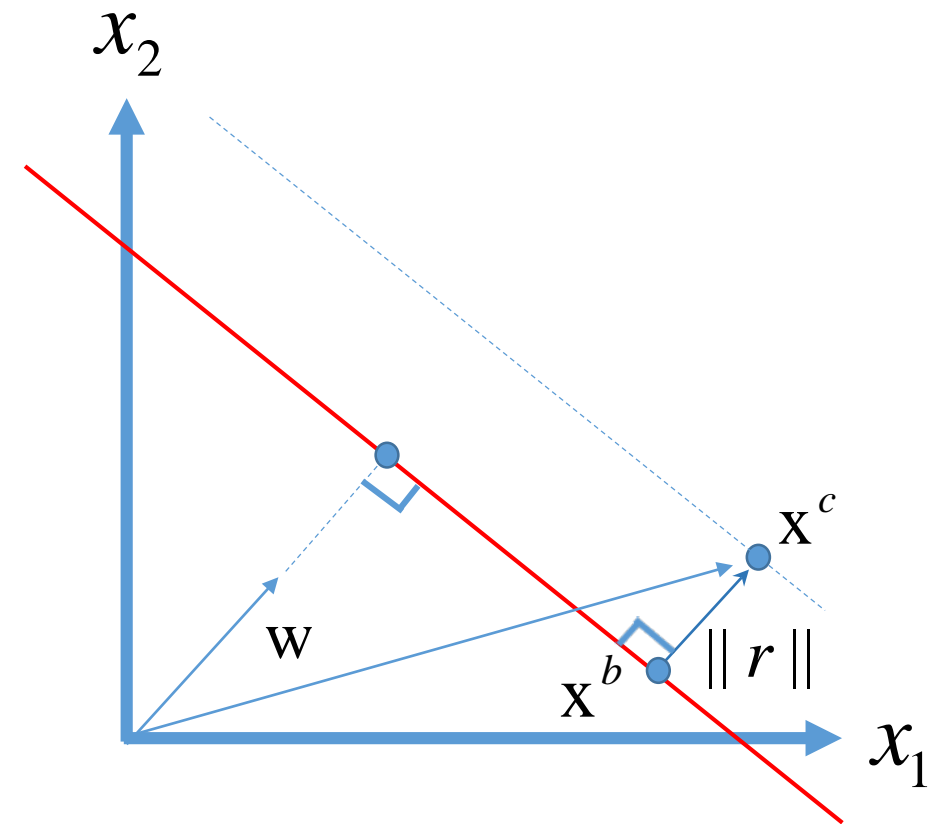
$$|y(\mathbf{x}^c)| = 1$$



# 44 Problem formulation

- Finding a decision boundary which maximizes the margin.

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



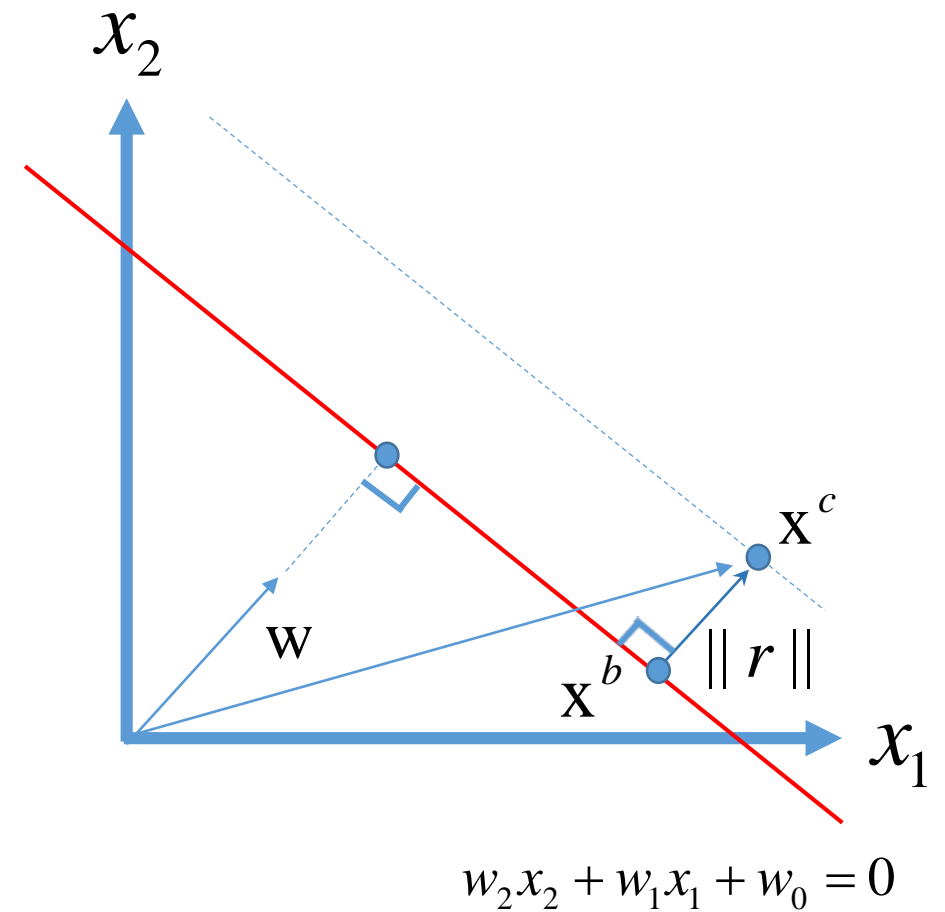
- Finding a decision boundary which maximizes the margin.

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

s.t.

$$t_n y(\mathbf{x}_n) > 0 \quad \longrightarrow \quad \text{Every data points are classified correctly.}$$

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



□ Let's make it a quadratic programming problem.

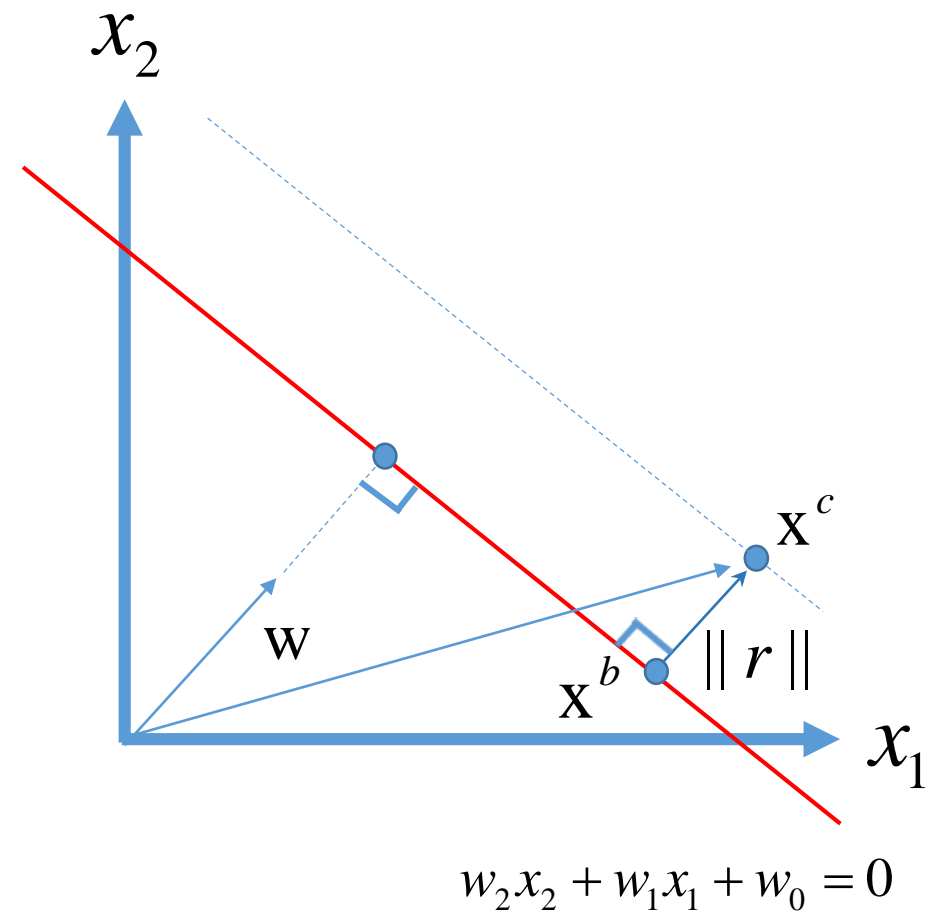
$$\max \frac{1}{\|w\|}$$

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

□ Finally

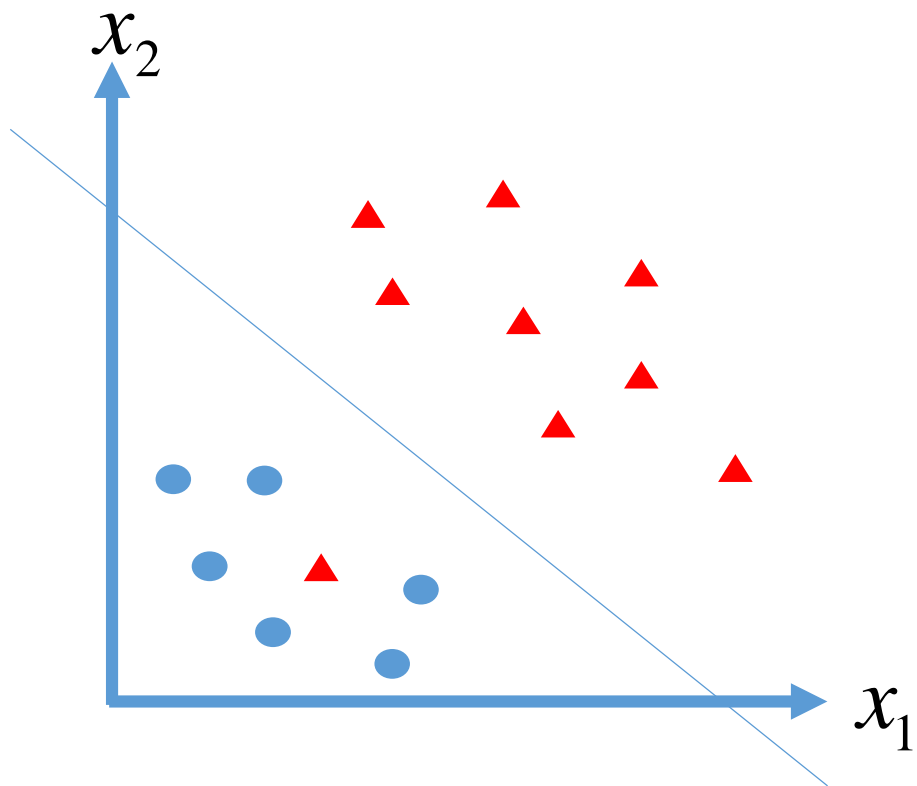
$$\min \frac{1}{2} \|w\|^2$$

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



# 47 How about non-linearly separable case?

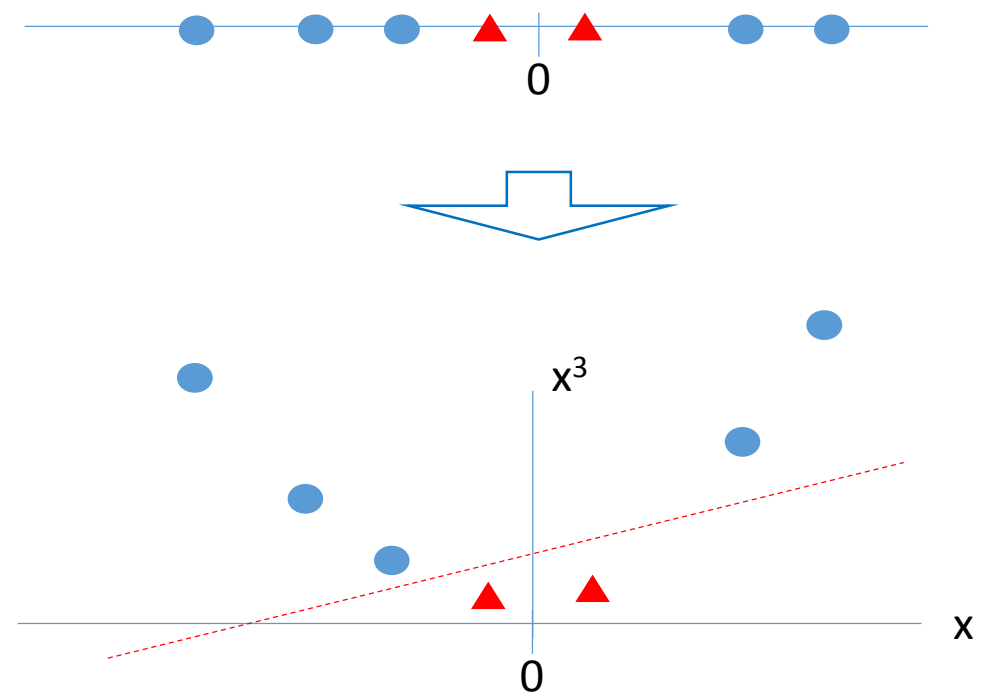
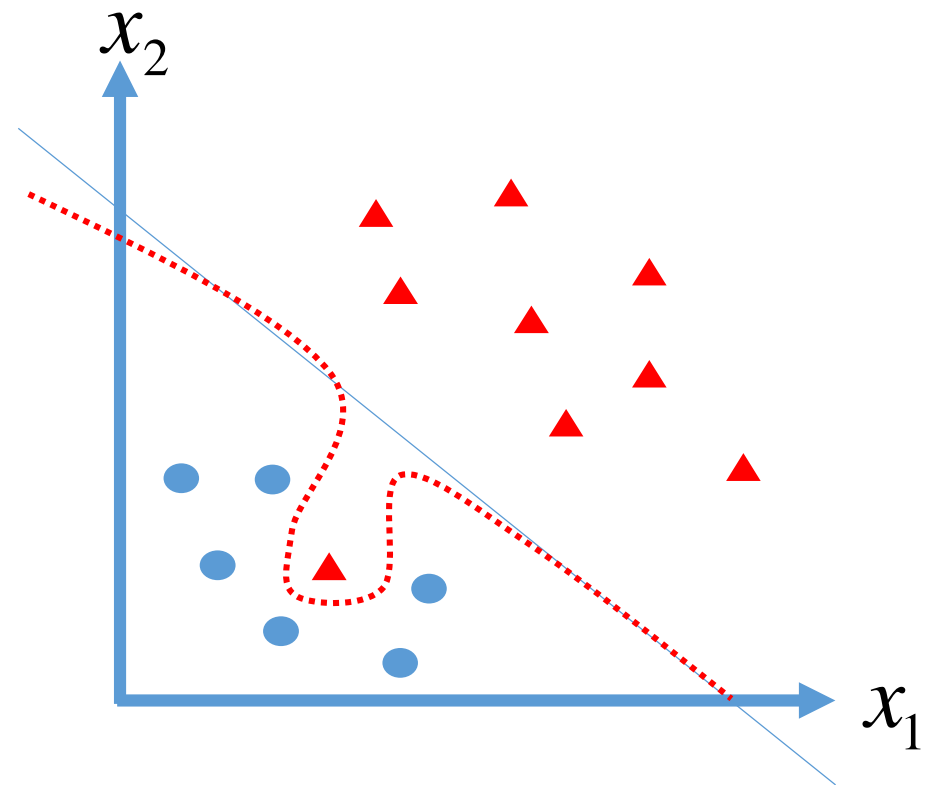
$$\begin{aligned} \min & \frac{1}{2} \| \mathbf{w} \|^2 \\ \text{s.t.} & t_n (\mathbf{w}^T \mathbf{x}_n + w_0) \geq 1, \quad \forall n \end{aligned}$$



# 48 How about non-linearly separable case?

$$\begin{aligned} \min & \frac{1}{2} \| \mathbf{w} \|^2 \\ \text{s.t.} & t_n (\mathbf{w}^T \mathbf{x}_n + w_0) \geq 1, \quad \forall n \end{aligned}$$

	Approaches
Option 1	Soft margin SVM
Option 2	Kernel trick

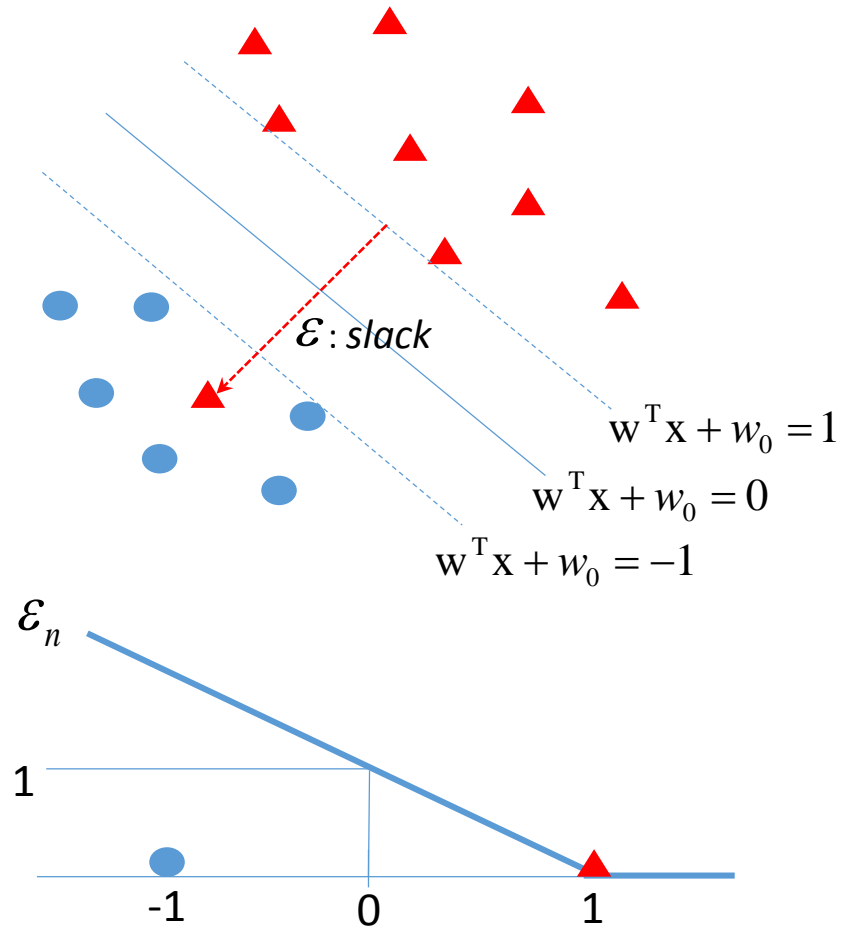


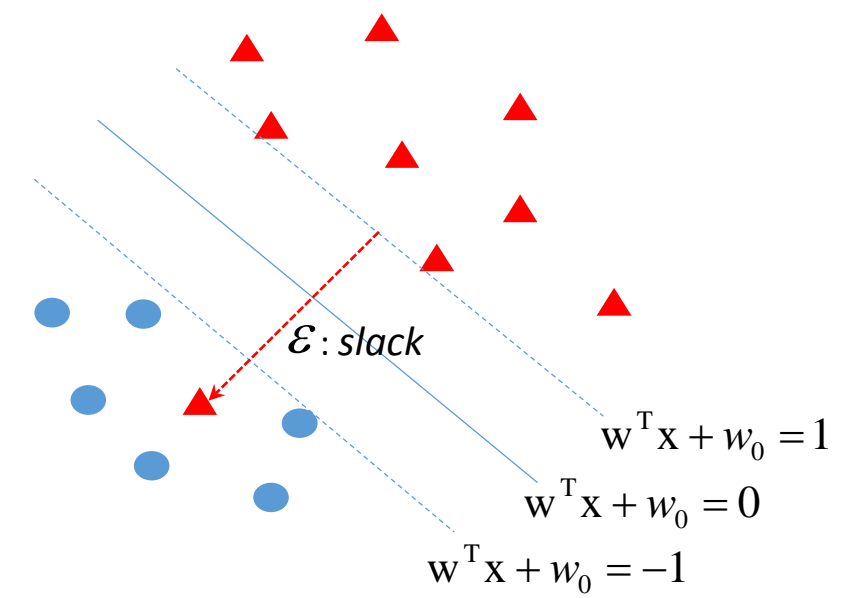


Soft margin SVM

Remember the constraint below?

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



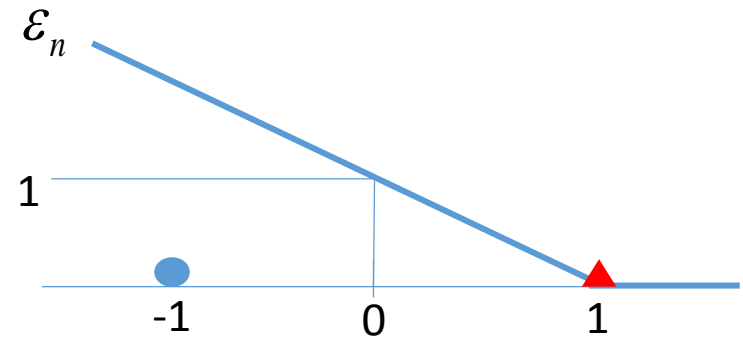


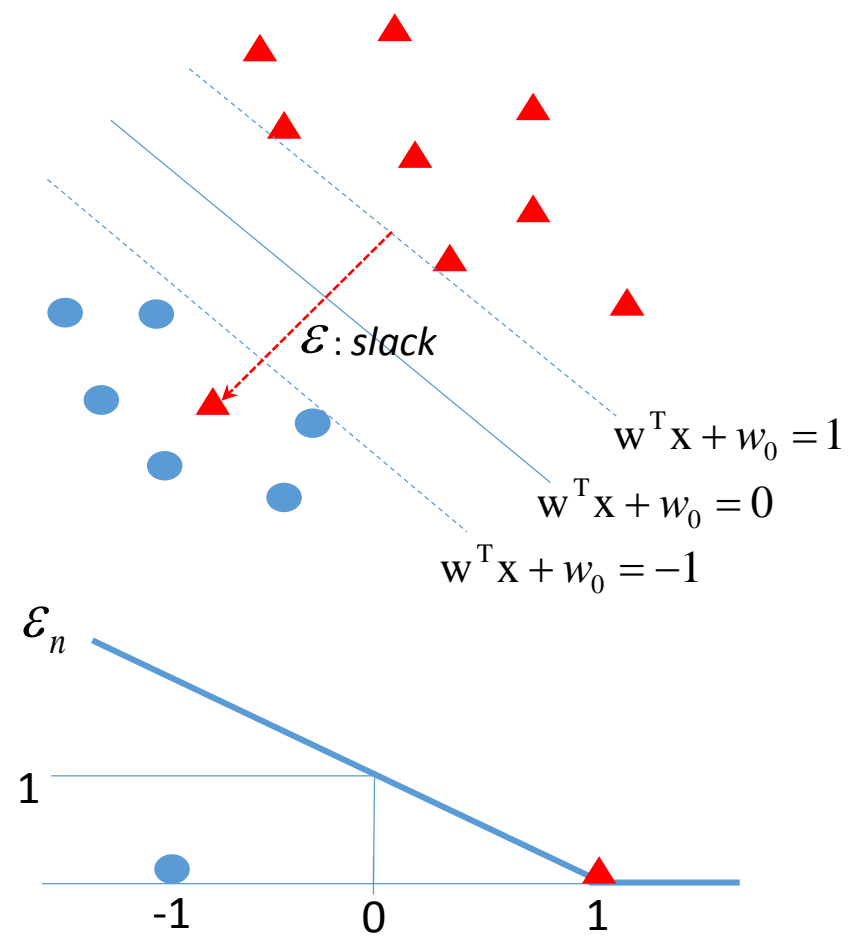
Remember the constraint below?

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

For the data points which are non-separable, we **relax** the constraint:

$$t_n (w^T x_n + w_0) \geq 1 - \varepsilon_n, \quad \forall n \quad \varepsilon_n \geq 0$$





- Remember the constraint below?

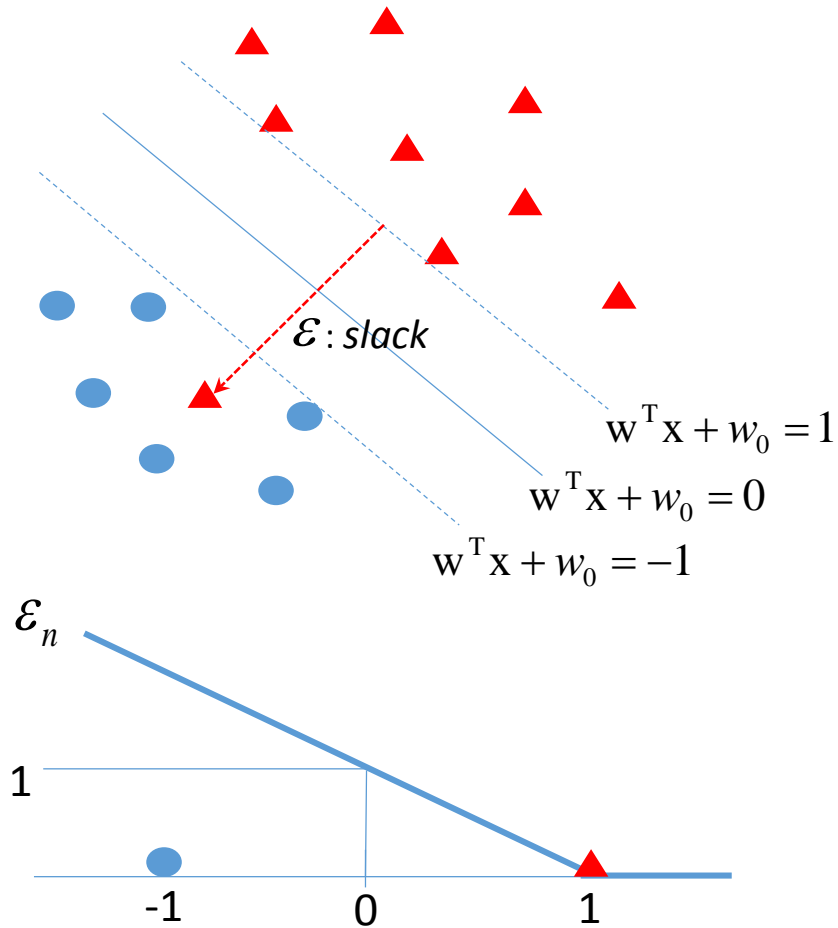
$$t_n (w^T x_n + w_0) \geq 1, \quad \forall n$$

- For the data points which are non-separable, we **relax** the constraint:

$$t_n (w^T x_n + w_0) \geq 1 - \epsilon_n, \quad \forall n \quad \epsilon_n \geq 0$$

- It says that **the distance between a data point and the decision boundary is allowed to be less than 1.**

# 53 Option 1: soft margin SVM



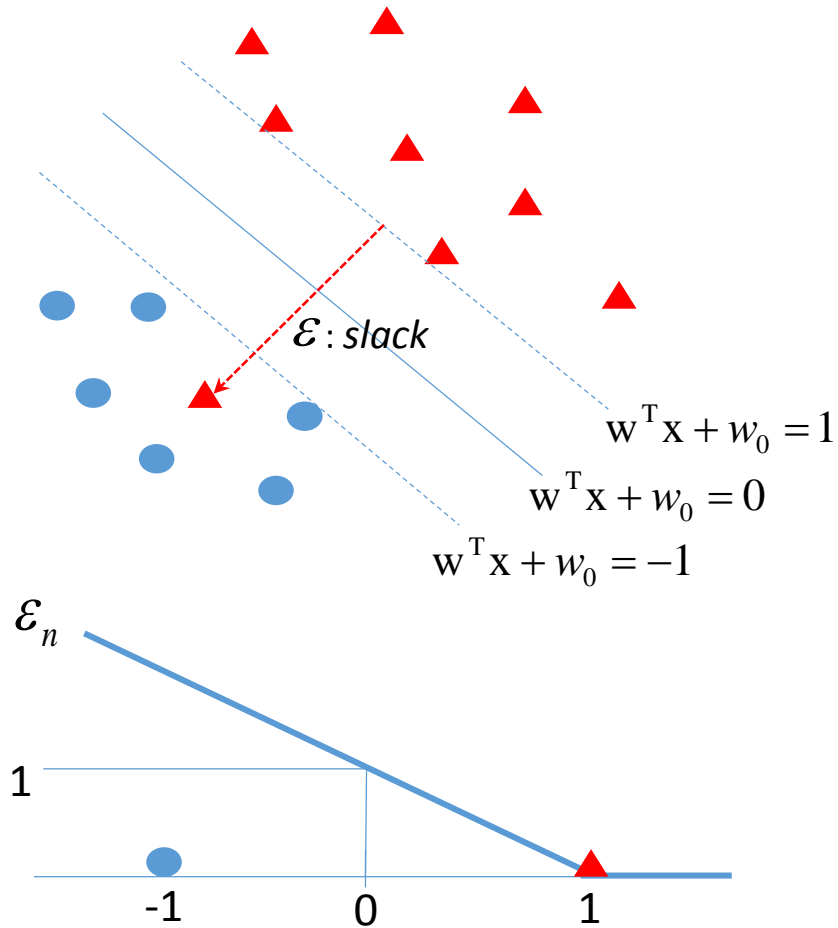
- Remember the constraint below?

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

- For the data points which are non-separable, we **relax** the constraint:

$$t_n (w^T x_n + w_0) \geq 1 - \varepsilon_n, \quad \forall n \quad \varepsilon_n \geq 0$$

- It says that **the distance between a data point and the decision boundary is allowed to be less than 1.**
- $\varepsilon_n$  is called slack variables.



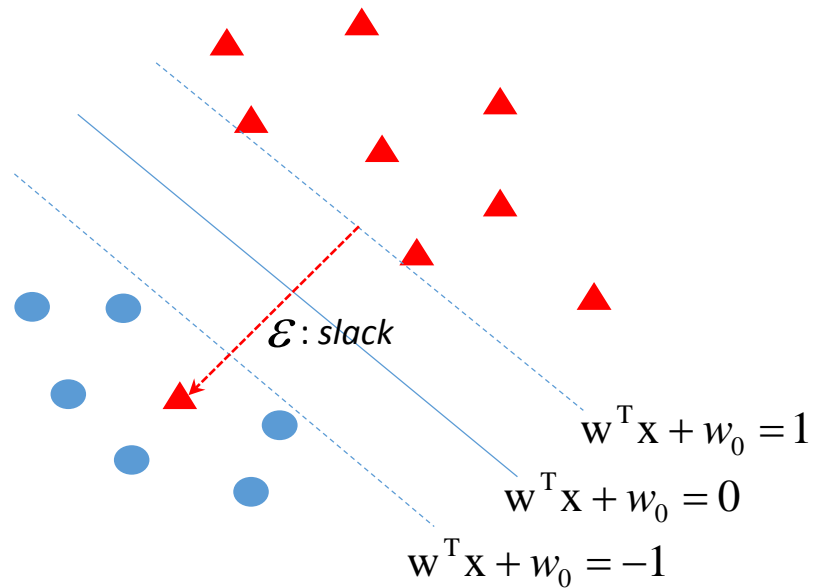
- Remember the constraint below?

$$t_n (\mathbf{w}^T \mathbf{x}_n + w_0) \geq 1, \quad \forall n$$

- For the data points which are non-separable, we **relax** the constraint:

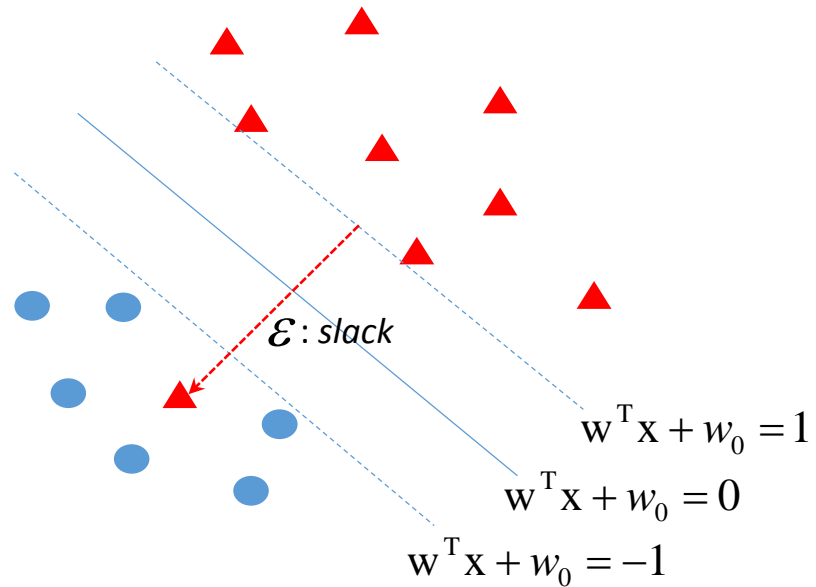
$$t_n (\mathbf{w}^T \mathbf{x}_n + w_0) \geq 1 - \varepsilon_n, \quad \forall n \quad \varepsilon_n \geq 0$$

- It says that **the distance between a data point and the decision boundary is allowed to be less than 1.**
- $\varepsilon_n$  is called slack variables.
- Question. Where is a data point when  $\varepsilon_n = 1$  ?



- So we have the constraint below. How about the objective function?

$$t_n (w^T x_n + w_0) \geq 1 - \varepsilon_n, \quad \forall n \quad \varepsilon_n \geq 0$$



- So we have the constraint below. How about the objective function?

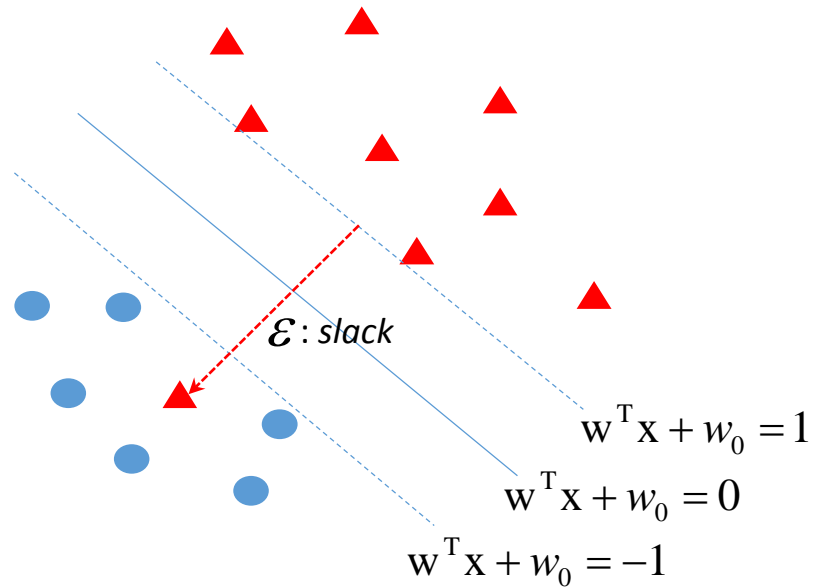
$$t_n(w^T x_n + w_0) \geq 1 - \mathcal{E}_n, \quad \forall n \quad \mathcal{E}_n \geq 0$$

- We want to minimize the slack.

$$\min \frac{1}{2} \|w\|^2 + C \sum_n \mathcal{E}_n$$



# 57 Option 1: soft margin SVM



- So we have the constraint below. How about the objective function?

$$t_n(w^T x_n + w_0) \geq 1 - \varepsilon_n, \quad \forall n \quad \varepsilon_n \geq 0$$

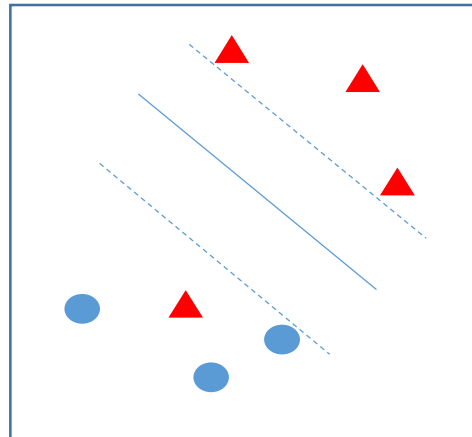
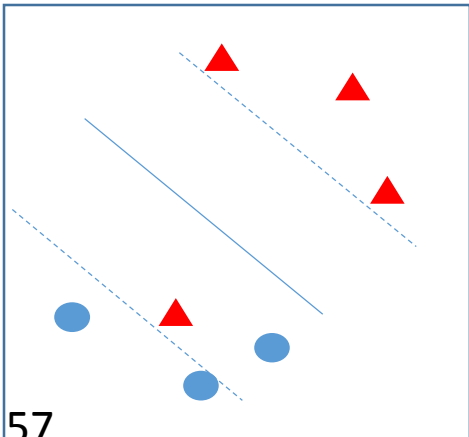
- We want to minimize the slack.

$$\min \frac{1}{2} \|w\|^2 + C \sum_n \varepsilon_n$$

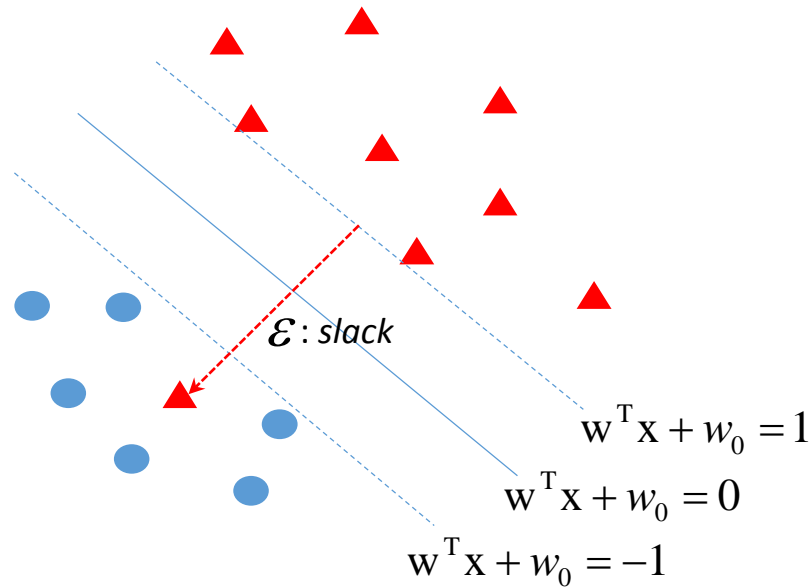
- If “C” is small, the dominant factor is  $\|w\|^2 / 2$ 
  - 1) Prefer large margin
  - 2) May cause large # of misclassified data points.

“C” is small

“C” is large



# 58 Option 1: soft margin SVM



- So we have the constraint below. How about the objective function?

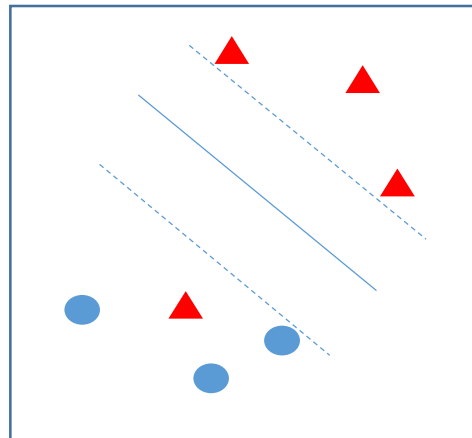
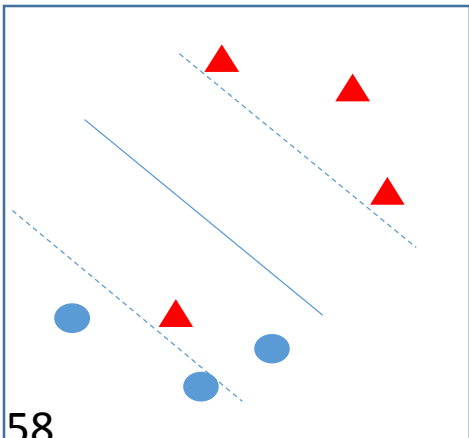
$$t_n (\mathbf{w}^T \mathbf{x}_n + w_0) \geq 1 - \varepsilon_n, \quad \forall n \quad \varepsilon_n \geq 0$$

- We want to minimize the slack.

$$\min \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_n \varepsilon_n$$

"C" is small

"C" is large



- If "C" is small, the slack contributes more
  - 1) Prefer large margin
  - 2) May cause large # of misclassified data points.
- If "C" is large, the slack contributes less
  - 1) Prefer less # of misclassified data points.
  - 2) May cause small margin.

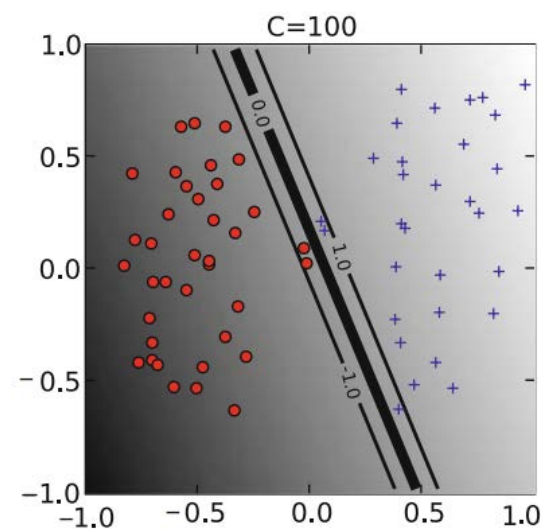
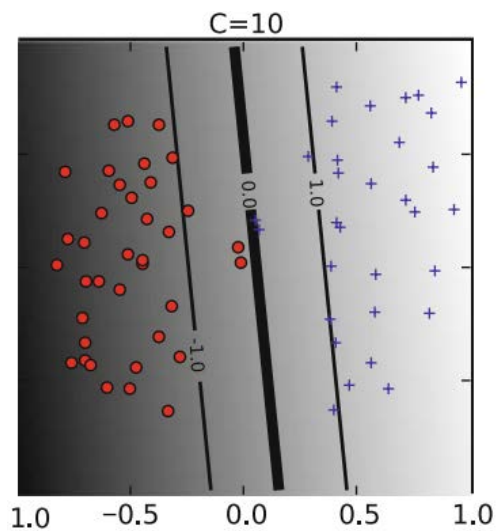
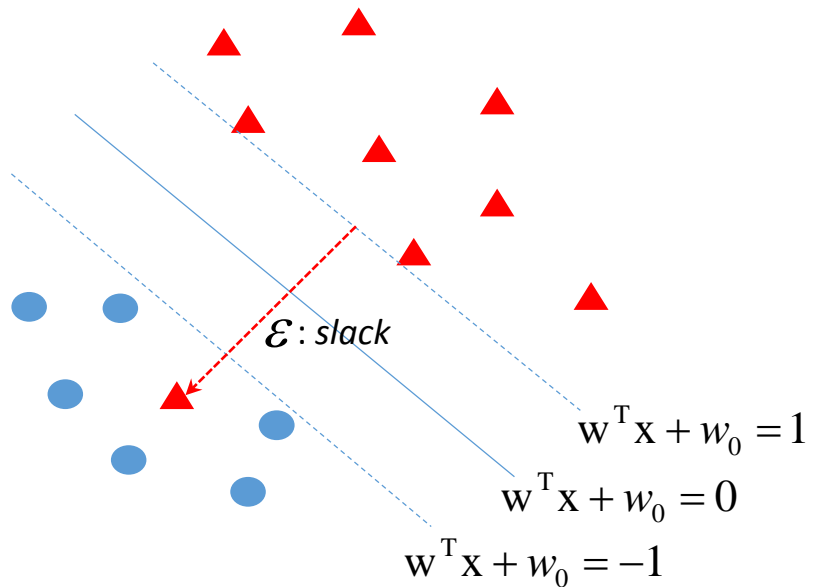
□ The formulation finally becomes

$$\min \frac{1}{2} \|w\|^2 + C \sum_n \varepsilon_n$$

*s.t.*

$$t_n (w^T x_n + w_0) \geq 1 - \varepsilon_n, \forall n$$

$$\varepsilon_n \geq 0$$

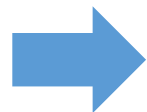


Kernel trick

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

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

Primal problem



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

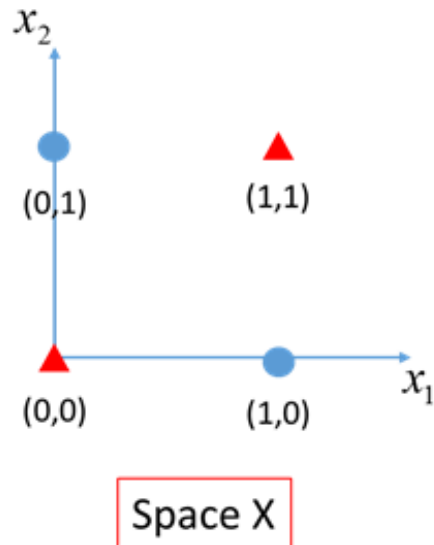
$$s.t. \quad \lambda \geq 0, \quad t^T \lambda = 0$$

Dual problem

- They are the same problem.
- $\lambda$ : Lagrange multipliers which corresponding to data points.
- $t$ : label (-1 or 1)
- It looks complicated why we border to use dual problem???

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

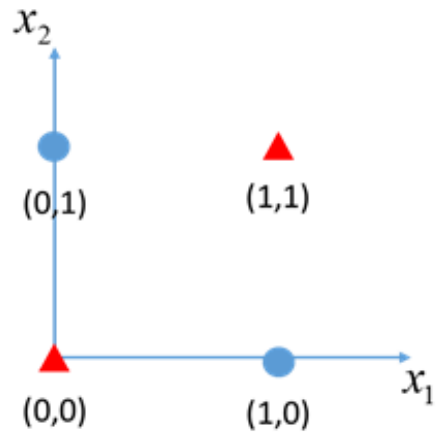
□ If data  $x_n$  are not linearly separable, what should we do?



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

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

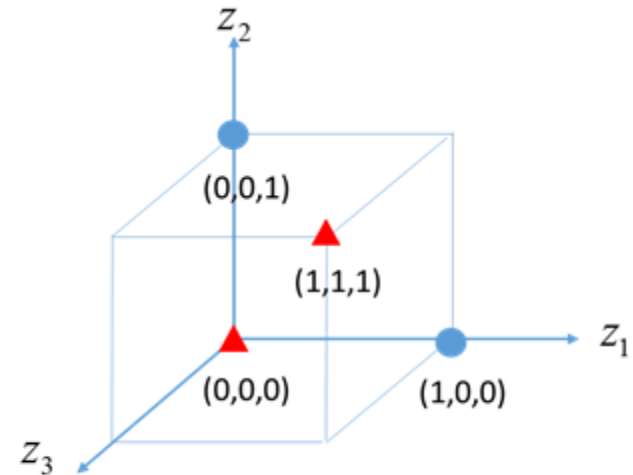
□ If data  $x_n$  are not linearly separable, what should we do?



Space X

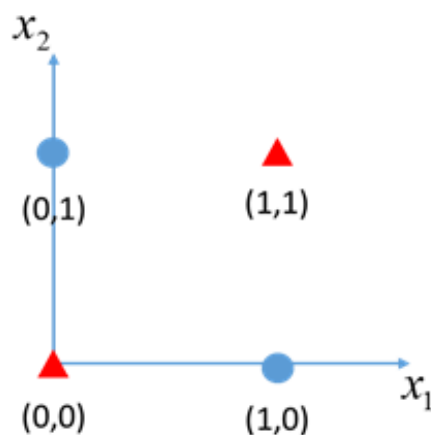
$$\phi(\mathbf{x}) = \begin{pmatrix} x_1^2 \\ x_1 x_2 \\ x_2^2 \end{pmatrix}$$

$$\phi(\mathbf{x}) \\ X \rightarrow Z$$

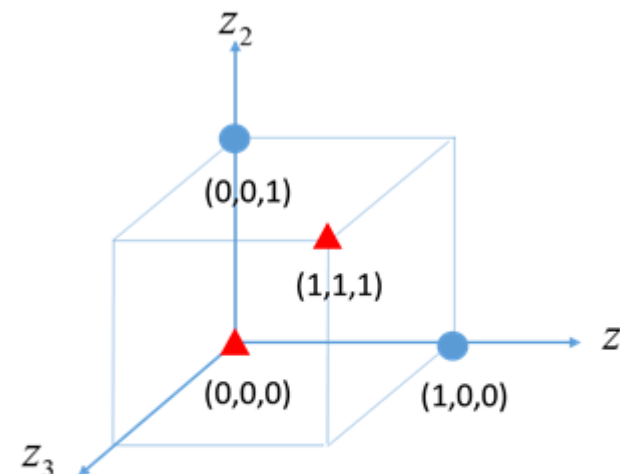
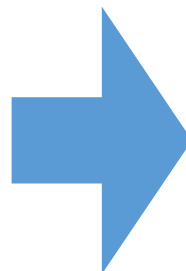


Space Z

- 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

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

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

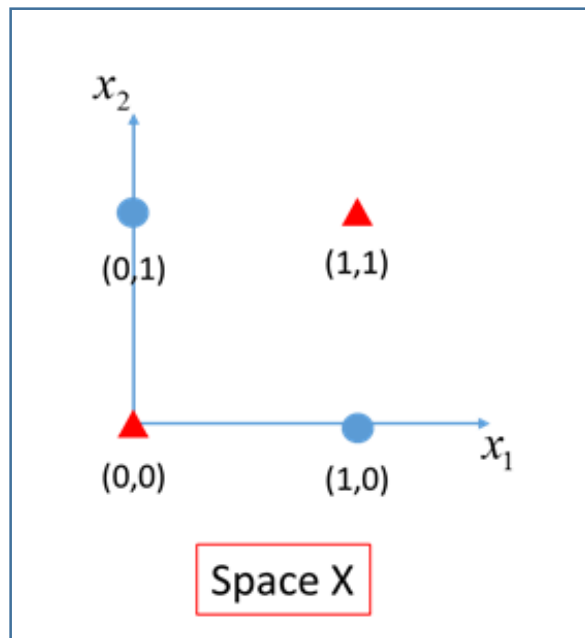
$$\min_{\lambda} L(\lambda) = \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N t_n t_m \lambda_n \lambda_m \mathbf{z}_n^T \mathbf{z}_m - \sum_{n=1}^N \lambda_n$$

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



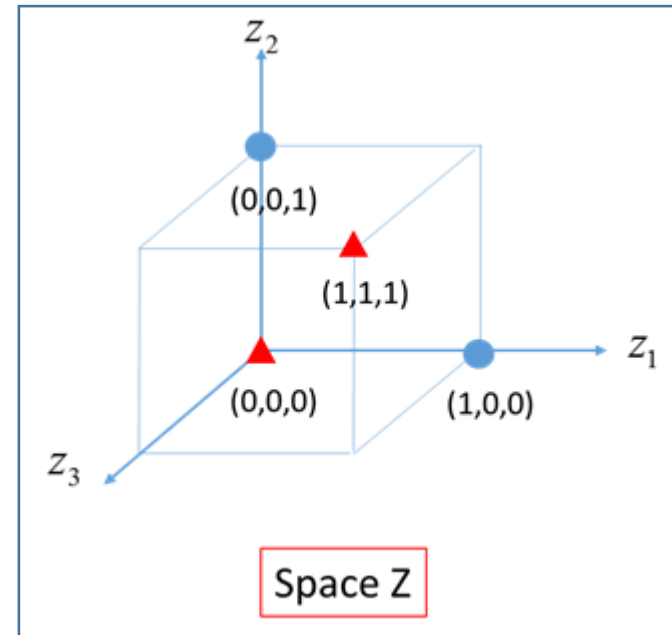
- Kernel function  $K()$  is a function which returns the scalar values (the inner product of two vectors:  $z_n$  and  $z_m$  in Z space) when the data points ( $x_n$  and  $x_m$  in X space) are given.

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



$$\mathbf{z}_n^T = \phi(\mathbf{x}_n^T)$$

$$\mathbf{z}_m = \phi(\mathbf{x}_m)$$



- 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.

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

Z space problem

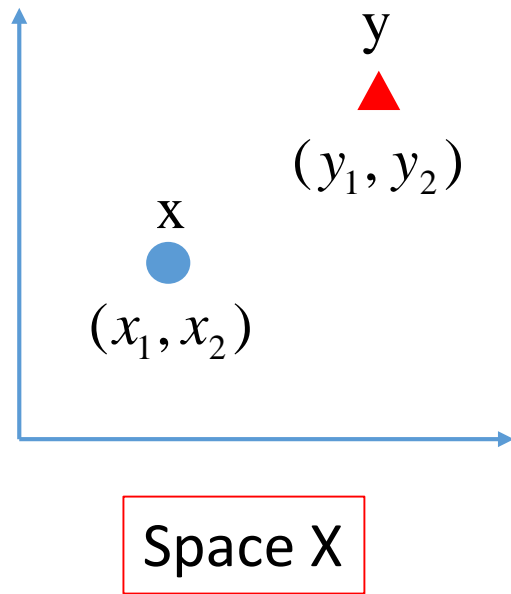


$$\min_{\lambda} L(\lambda) = \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N t_n t_m \lambda_n \lambda_m K(x_n^T x_m) - \sum_{n=1}^N \lambda_n$$
$$s.t. \quad \lambda \geq 0, \quad t^T \lambda = 0$$

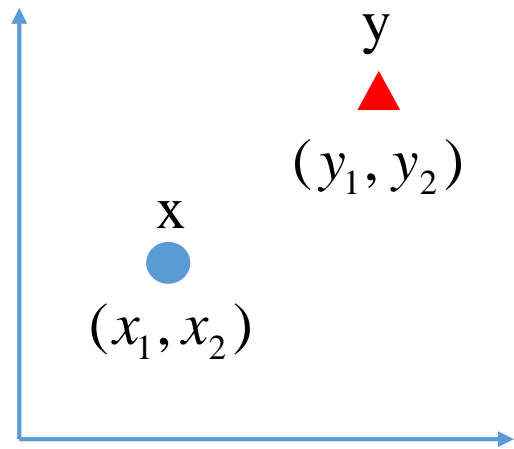
X space problem

Z space problem can be formulated with data in X space

# 67 Polynomial kernel of degree 2



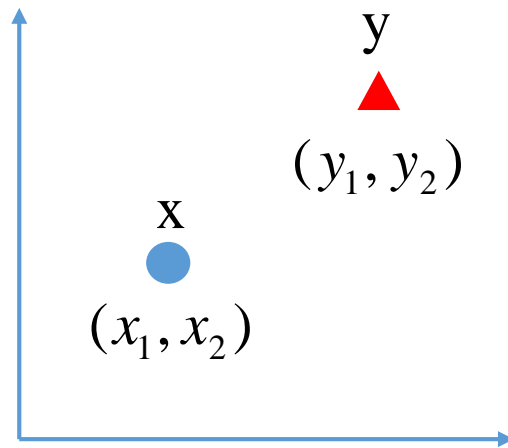
$$K(x, y) = (xy)^2$$



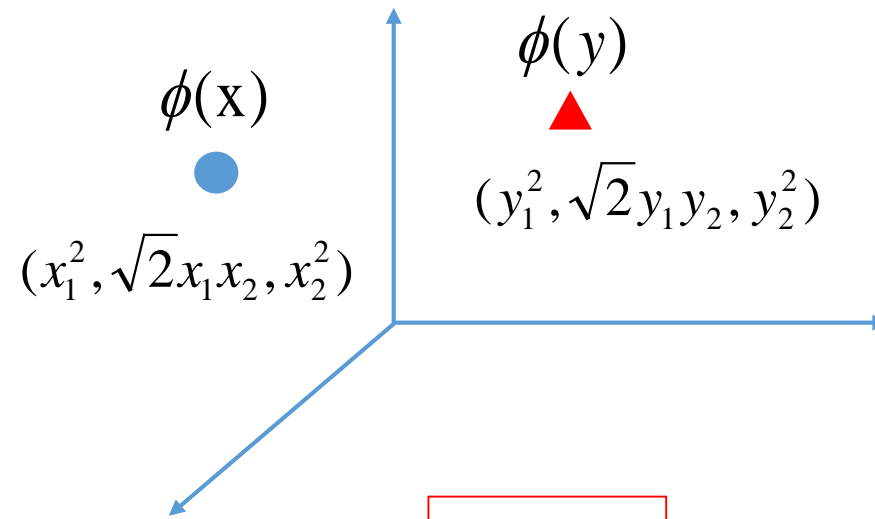
Space X

$$\begin{aligned} K(\mathbf{x}, \mathbf{y}) &= (\mathbf{x}\mathbf{y})^2 \\ &= ((x_1, x_2) \cdot (y_1, y_2))^2 \\ &= (x_1 y_1 + x_2 y_2)^2 \\ &= x_1^2 y_1^2 + 2x_1 x_2 y_1 y_2 + x_2^2 y_2^2 \end{aligned}$$

# 69 Polynomial kernel of degree 2



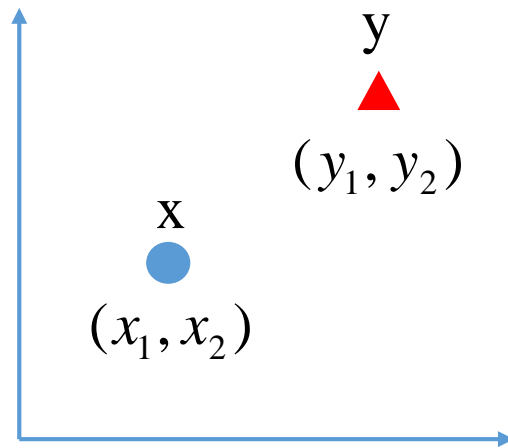
Space X



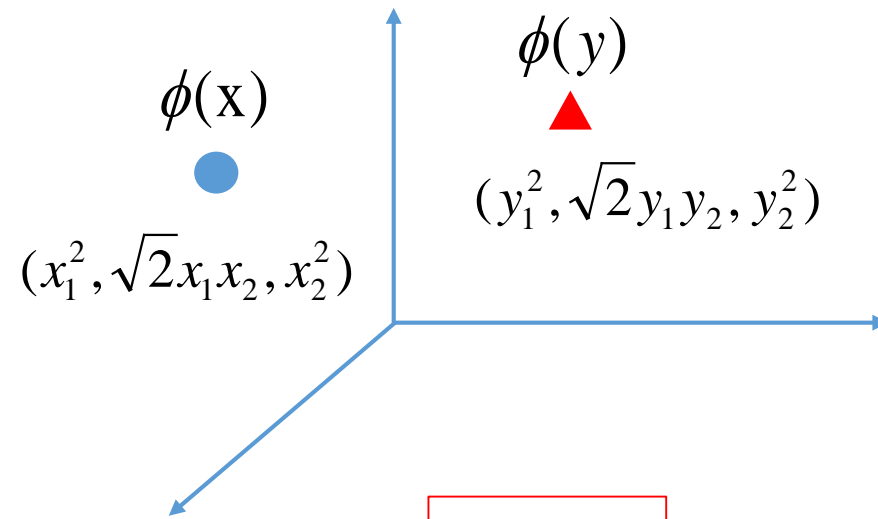
Space Z

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

# 70 Polynomial kernel of degree 2



Space X

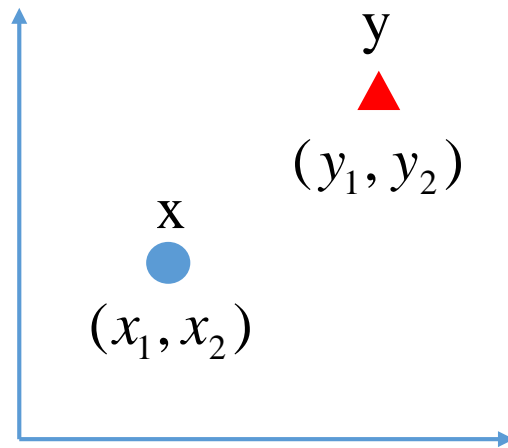


Space Z

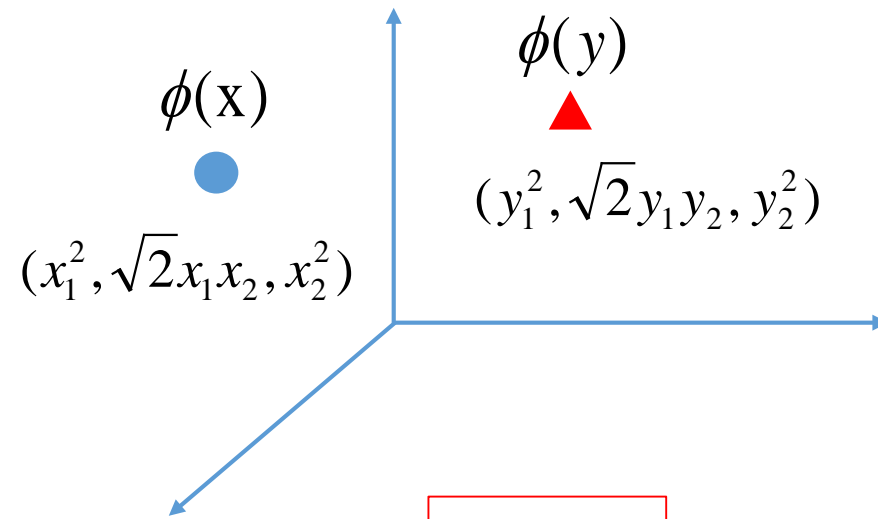
$$\begin{aligned} K(x, y) &= (xy)^2 \\ &= ((x_1, x_2) \cdot (y_1, y_2))^2 \\ &= (x_1y_1 + x_2y_2)^2 \\ &= x_1^2y_1^2 + 2x_1x_2y_1y_2 + x_2^2y_2^2 \end{aligned}$$

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

# 71 Polynomial kernel of degree 2



Space X



Space Z

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

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

Mapping to 3-dimension

$$K(\mathbf{x}_n, \mathbf{x}_m) = \exp(-\alpha \|\mathbf{x}_n - \mathbf{x}_m\|^2)$$



$$\begin{aligned} K(\mathbf{x}_n, \mathbf{x}_m) &= \exp(-\alpha \|\mathbf{x}_n - \mathbf{x}_m\|^2) \\ &= \exp(-\alpha x_n^2) \exp(-\alpha x_m^2) \exp(2\alpha x_n x_m) \end{aligned}$$

$$K(\mathbf{x}_n, \mathbf{x}_m) = \exp(-\alpha \|\mathbf{x}_n - \mathbf{x}_m\|^2)$$

$$= \exp(-\alpha \mathbf{x}_n^2) \exp(-\alpha \mathbf{x}_m^2) \exp(2\alpha \mathbf{x}_n \mathbf{x}_m)$$

$$= \exp(-\alpha \mathbf{x}_n^2) \exp(-\alpha \mathbf{x}_m^2) \sum_{k=0}^{\infty} \frac{(2\alpha)^k (\mathbf{x}_n)^k (\mathbf{x}_m)^k}{k!}$$

Taylor series expansion of  
an exponential function

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



$$K(\mathbf{x}_n, \mathbf{x}_m) = \exp(-\alpha \|\mathbf{x}_n - \mathbf{x}_m\|^2)$$

$$= \exp(-\alpha \mathbf{x}_n^2) \exp(-\alpha \mathbf{x}_m^2) \exp(2\alpha \mathbf{x}_n \mathbf{x}_m)$$

$$= \exp(-\alpha \mathbf{x}_n^2) \exp(-\alpha \mathbf{x}_m^2) \sum_{k=0}^{\infty} \frac{(2\alpha)^k (\mathbf{x}_n)^k (\mathbf{x}_m)^k}{k!}$$

$$= \sum_{k=0}^{\infty} \sqrt{\frac{(2\alpha)^k}{k!}} \exp(-\alpha \mathbf{x}_n^2) (\mathbf{x}_n)^k \sqrt{\frac{(2\alpha)^k}{k!}} \exp(-\alpha \mathbf{x}_m^2) (\mathbf{x}_m)^k$$

Taylor series expansion of  
an exponential function

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



# Gaussian Kernel: derivation (inner product in the infinite z space)

$$K(\mathbf{x}_n, \mathbf{x}_m) = \exp(-\alpha \|\mathbf{x}_n - \mathbf{x}_m\|^2)$$

Taylor series expansion of  
an exponential function

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

$$= \exp(-\alpha \mathbf{x}_n^2) \exp(-\alpha \mathbf{x}_m^2) \exp(2\alpha \mathbf{x}_n \mathbf{x}_m)$$

$$= \exp(-\alpha \mathbf{x}_n^2) \exp(-\alpha \mathbf{x}_m^2) \sum_{k=0}^{\infty} \frac{(2\alpha)^k (\mathbf{x}_n)^k (\mathbf{x}_m)^k}{k!}$$

$$= \sum_{k=0}^{\infty} \sqrt{\frac{(2\alpha)^k}{k!}} \exp(-\alpha \mathbf{x}_n^2) (\mathbf{x}_n)^k \sqrt{\frac{(2\alpha)^k}{k!}} \exp(-\alpha \mathbf{x}_m^2) (\mathbf{x}_m)^k$$



# 77 Gaussian Kernel: derivation (inner product in the infinite z space)

$$K(\mathbf{x}_n, \mathbf{x}_m) = \exp(-\alpha \|\mathbf{x}_n - \mathbf{x}_m\|^2)$$

Taylor series expansion of an exponential function

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

$$= \exp(-\alpha \mathbf{x}_n^2) \exp(-\alpha \mathbf{x}_m^2) \exp(2\alpha \mathbf{x}_n \mathbf{x}_m)$$

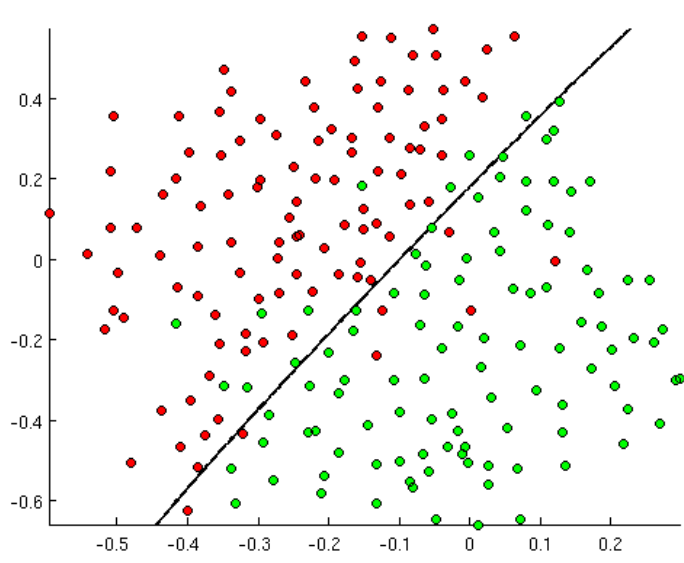
$$= \exp(-\alpha \mathbf{x}_n^2) \exp(-\alpha \mathbf{x}_m^2) \sum_{k=0}^{\infty} \frac{(2\alpha)^k (\mathbf{x}_n)^k (\mathbf{x}_m)^k}{k!}$$

$$= \sum_{k=0}^{\infty} \sqrt{\frac{(2\alpha)^k}{k!}} \exp(-\alpha \mathbf{x}_n^2) (\mathbf{x}_n)^k \sqrt{\frac{(2\alpha)^k}{k!}} \exp(-\alpha \mathbf{x}_m^2) (\mathbf{x}_m)^k$$

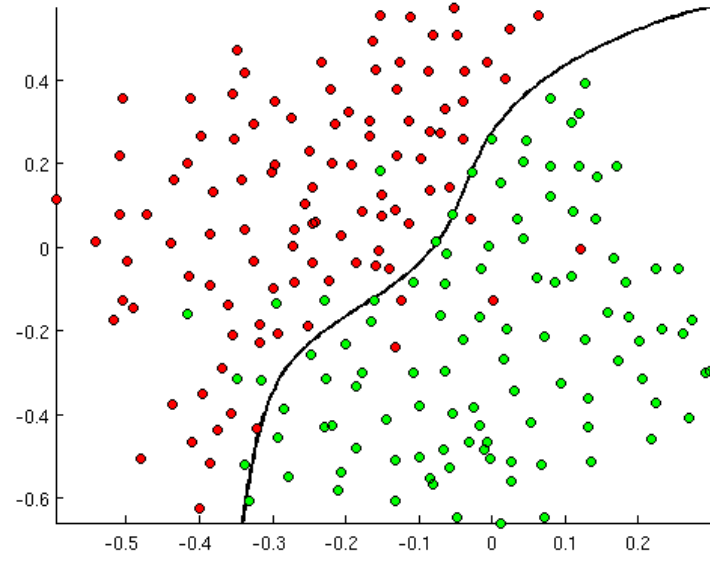
$$= \phi(\mathbf{x}_n) \phi(\mathbf{x}_m)$$

Mapping to infinite-dimension !

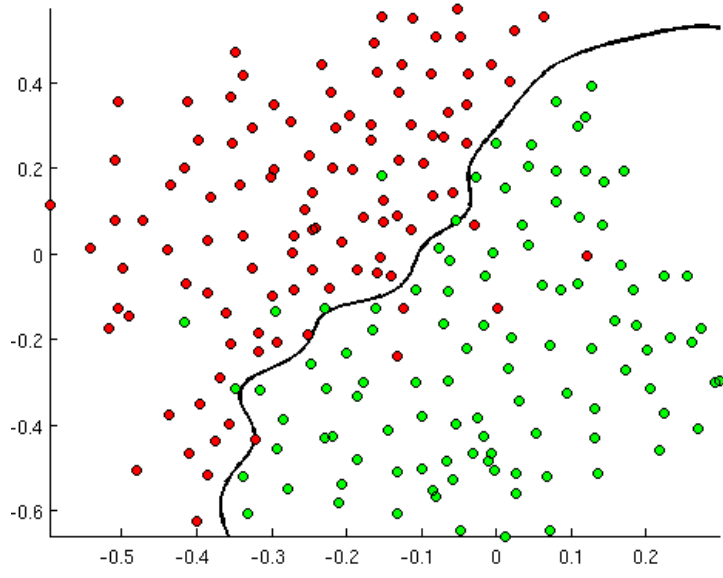
$\alpha = 1$



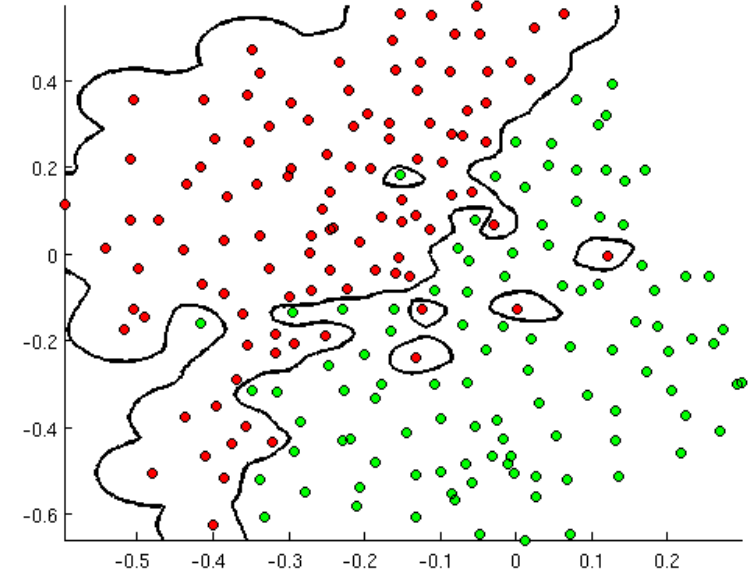
$\alpha = 10$



$\alpha = 100$



$\alpha = 1000$



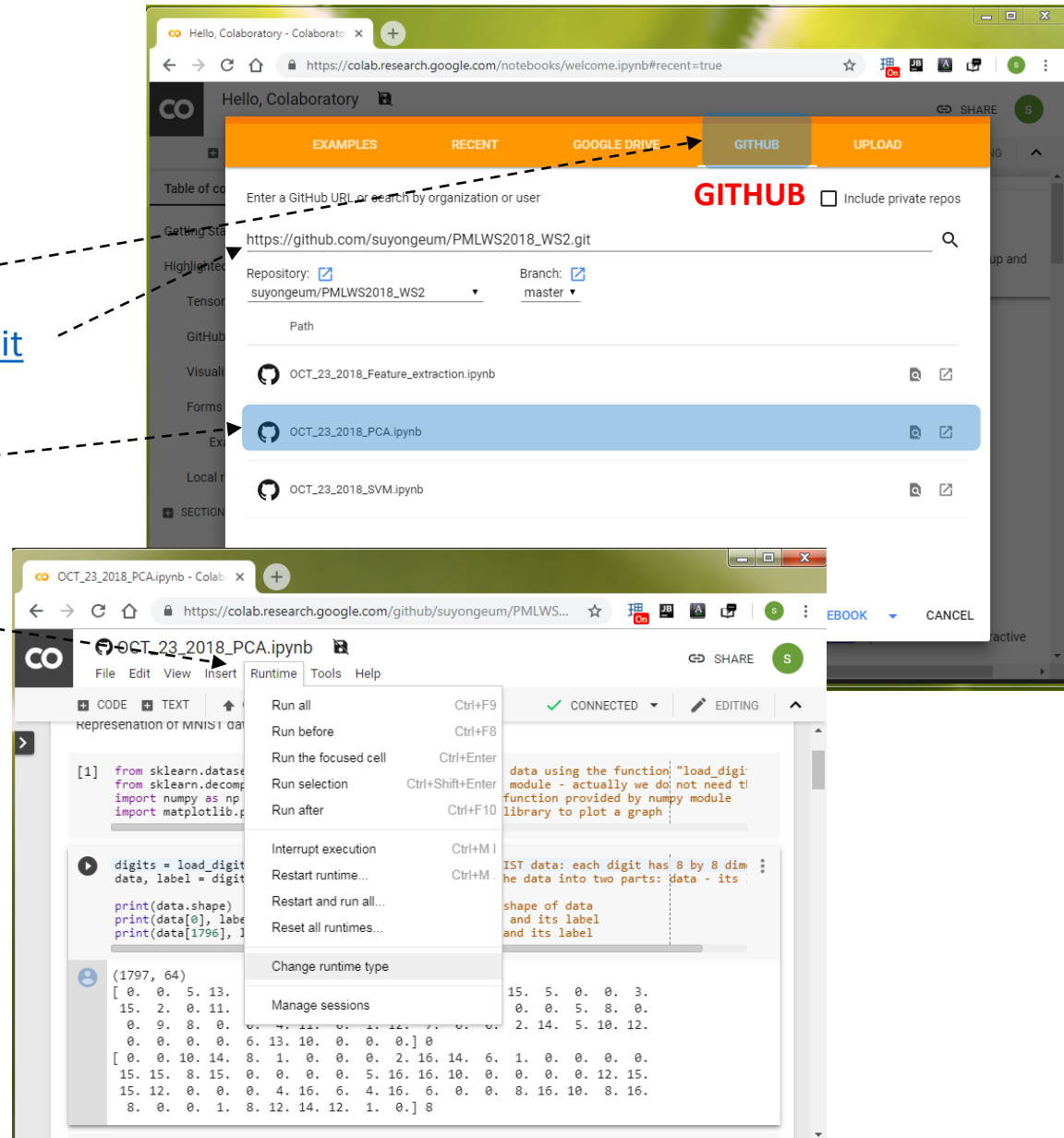
# Hand-on Experience

- ❑ A web base free google cloud service
  - Jupyter Notebook with Google Drive
  
- ❑ You can even use GPU for free!
  - Good but it provides the best effort service
    - You must save your things in your google drive or somewhere else.
  
- ❑ Resource check
  - `!cat /proc/meminfo`
  - `!cat /proc/cpuinfo`
  - `!df -h`



# 81 Colab: Principal Component Analysis (PCA)

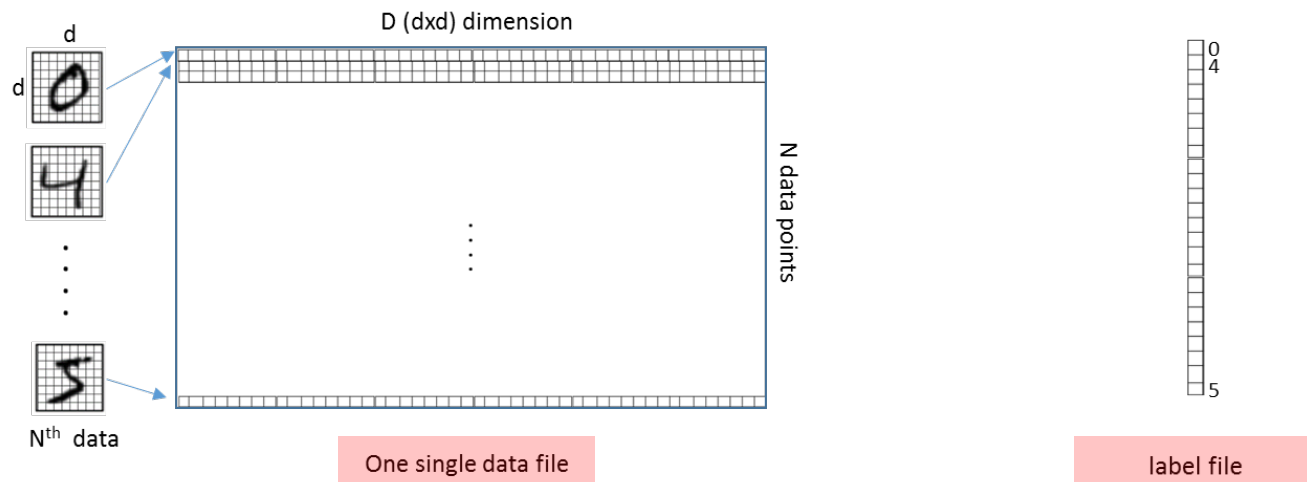
- 1) Go to the Colab
  - <https://colab.research.google.com>
- 2) Select “GITHUB” and copy the link below into
  - [https://github.com/suyongeuem/PMLWS2018\\_WS2.git](https://github.com/suyongeuem/PMLWS2018_WS2.git)
- 3) Select the notebook in the list
  - OCT\_23\_2018\_PCA.ipynb
- 4) Go to “Runtime” – “Change runtime type”
  - Python 3
  - GPU
- 5) Save it into your gdrive
  - “File” - “Save a copy in Drive ...”



# Data loading: MNIST

## □ MNIST data set

- <http://yann.lecun.com/exdb/mnist/>
- Training data
  - One single file (45M) which includes 60,000 hand digit images for training,
  - One single file (59K) which includes corresponding labels.
- Testing data
  - One single file (7.5M) which includes 10,000 hand digit images for testing,
  - One single file (9.8K) which includes corresponding labels.



# Making your gdrive as a working directory

- ❑ Defining a root directory where you mount your gdrive

```
from google.colab import drive
drive.mount('/gdrive/')
```

- /gdrive/My Drive/Colab Notebooks/

- ❑ Running time measurement

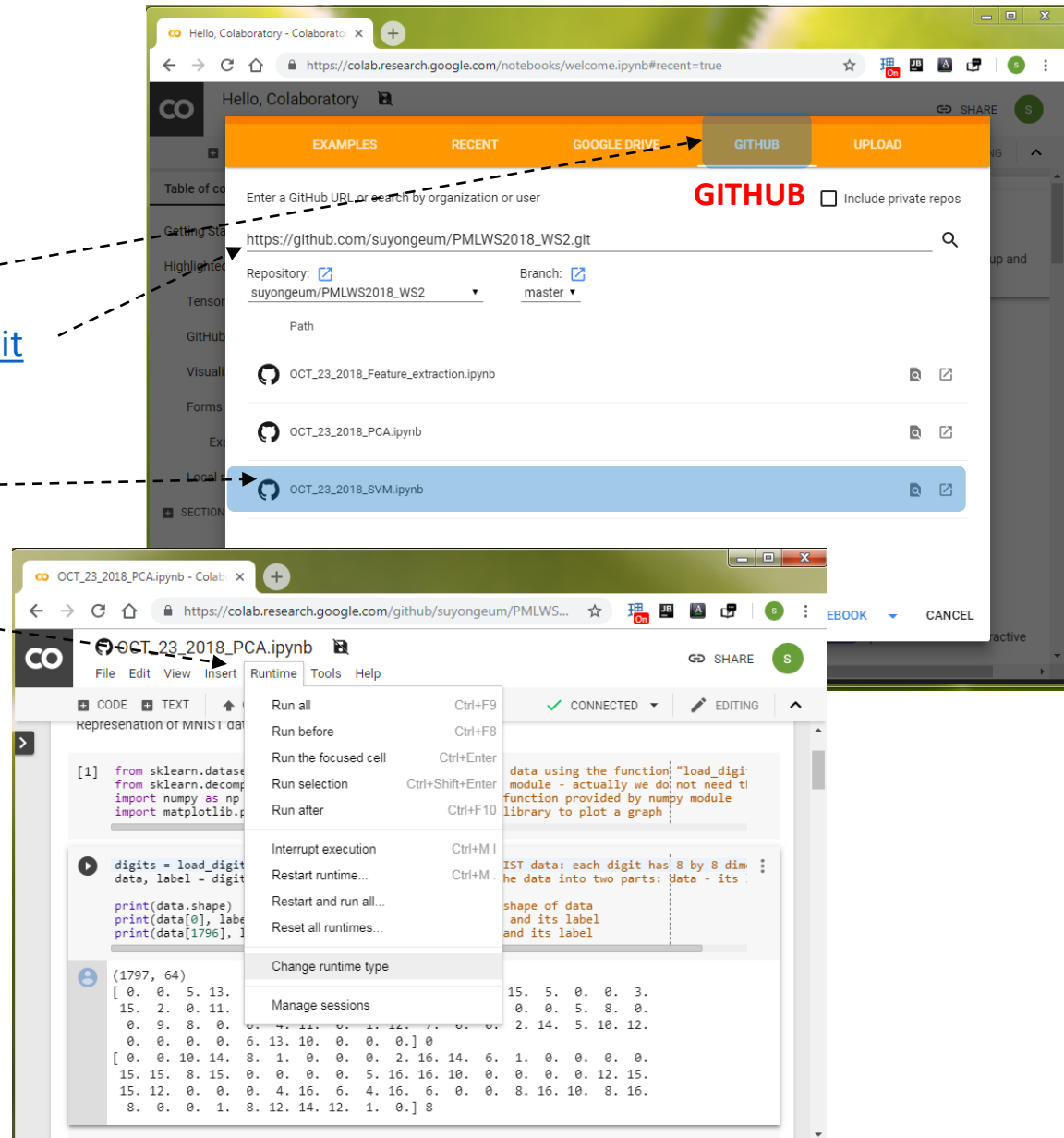
```
import datetime
before = datetime.datetime.now().timestamp()

...

after = datetime.datetime.now().timestamp()
print( "Time taken:", after - before)
```

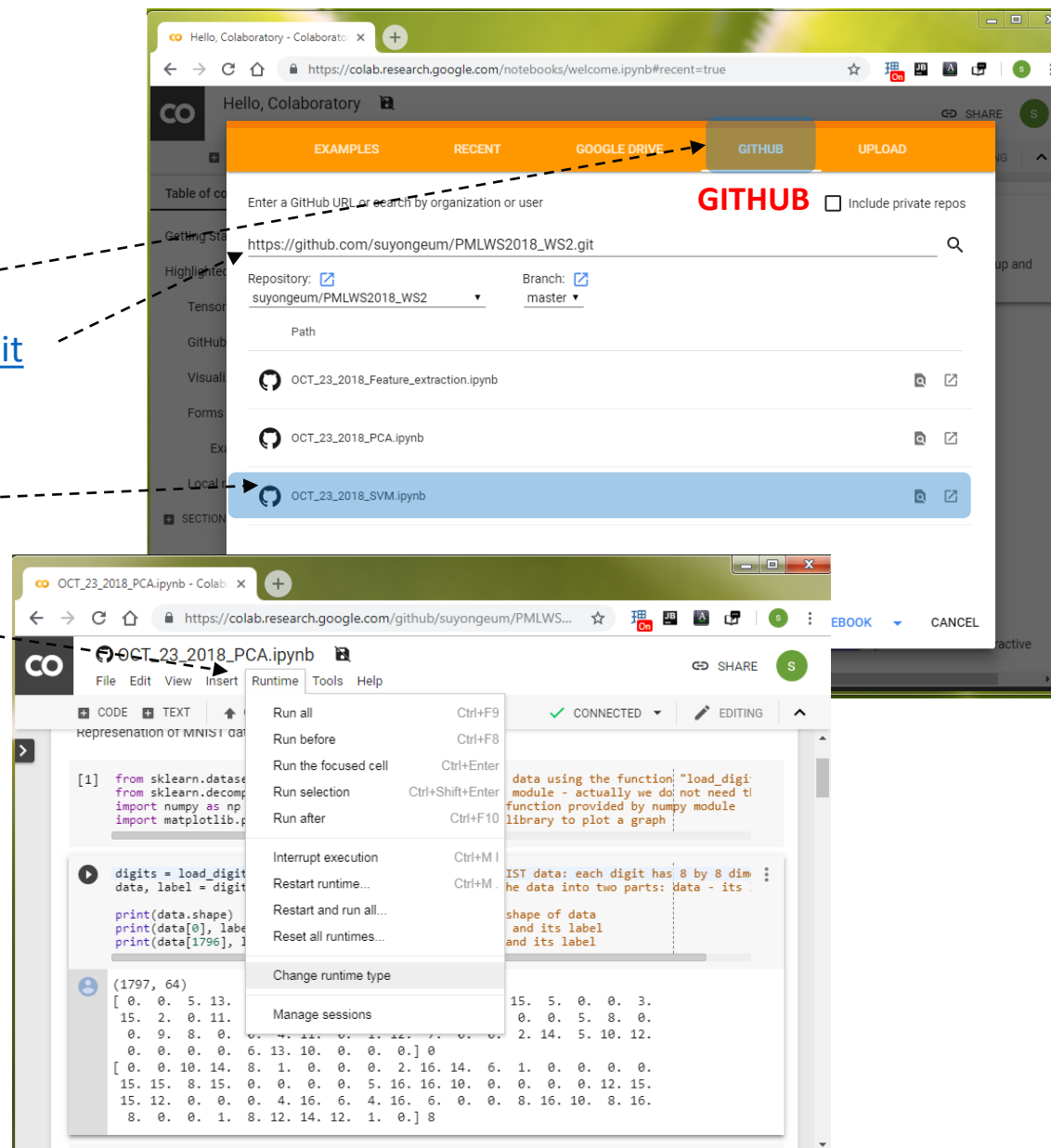
# 84 Colab: Support Vector Machine (SVM)

- 1) Go to the Colab
  - <https://colab.research.google.com>
- 2) Select “GITHUB” and copy the link below into
  - [https://github.com/suyongeuem/PMLWS2018\\_WS2.git](https://github.com/suyongeuem/PMLWS2018_WS2.git)
- 3) Select the notebook in the list
  - [OCT\\_23\\_2018\\_SVM.ipynb](#)
- 4) Go to “Runtime” – “Change runtime type”
  - Python 3
  - GPU
- 5) Save it into your gdrive
  - “File” - “Save a copy in Drive ...”



# 85 Colab: Feature extraction

- 1) Go to the Colab
  - <https://colab.research.google.com>
- 2) Select “GITHUB” and copy the link below into
  - [https://github.com/suyongeuem/PMLWS2018\\_WS2.git](https://github.com/suyongeuem/PMLWS2018_WS2.git)
- 3) Select the notebook in the list
  - [OCT\\_23\\_2018\\_Feature\\_extraction.ipynb](#)
- 4) Go to “Runtime” – “Change runtime type”
  - Python 3
  - GPU
- 5) Save it into your gdrive
  - “File” - “Save a copy in Drive ...”



Backup Slides

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

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

$$\mathbf{w} = \sum_{z_n \in SV} \lambda_n t_n \mathbf{z}_n \quad w_0 = t_n - \sum_{z_n \in SV} \lambda_n t_n z_n z_n = t_n - \sum_{z_n \in SV} \lambda_n t_n K(\mathbf{x}_n, \mathbf{x}_n)$$

$$\text{sign}(\mathbf{w}^T \mathbf{z} + w_0)$$

$$\text{sign} \left( \sum \lambda_n t_n \mathbf{z}_n \mathbf{z} + t_n - \sum_{z_n \in SV} \lambda_n t_n K(\mathbf{x}_n, \mathbf{x}_n) \right)$$

$$\text{sign} \left( \sum \lambda_n t_n K(x_n, x) + t_n - \sum \lambda_n t_n K(x_n, x_n) \right)$$

- Now you have a function, which classifies a data point in **z space** without mapping the data point to **z space** at all.
- Do you see why it is called a trick?