



Practical Machine Learning

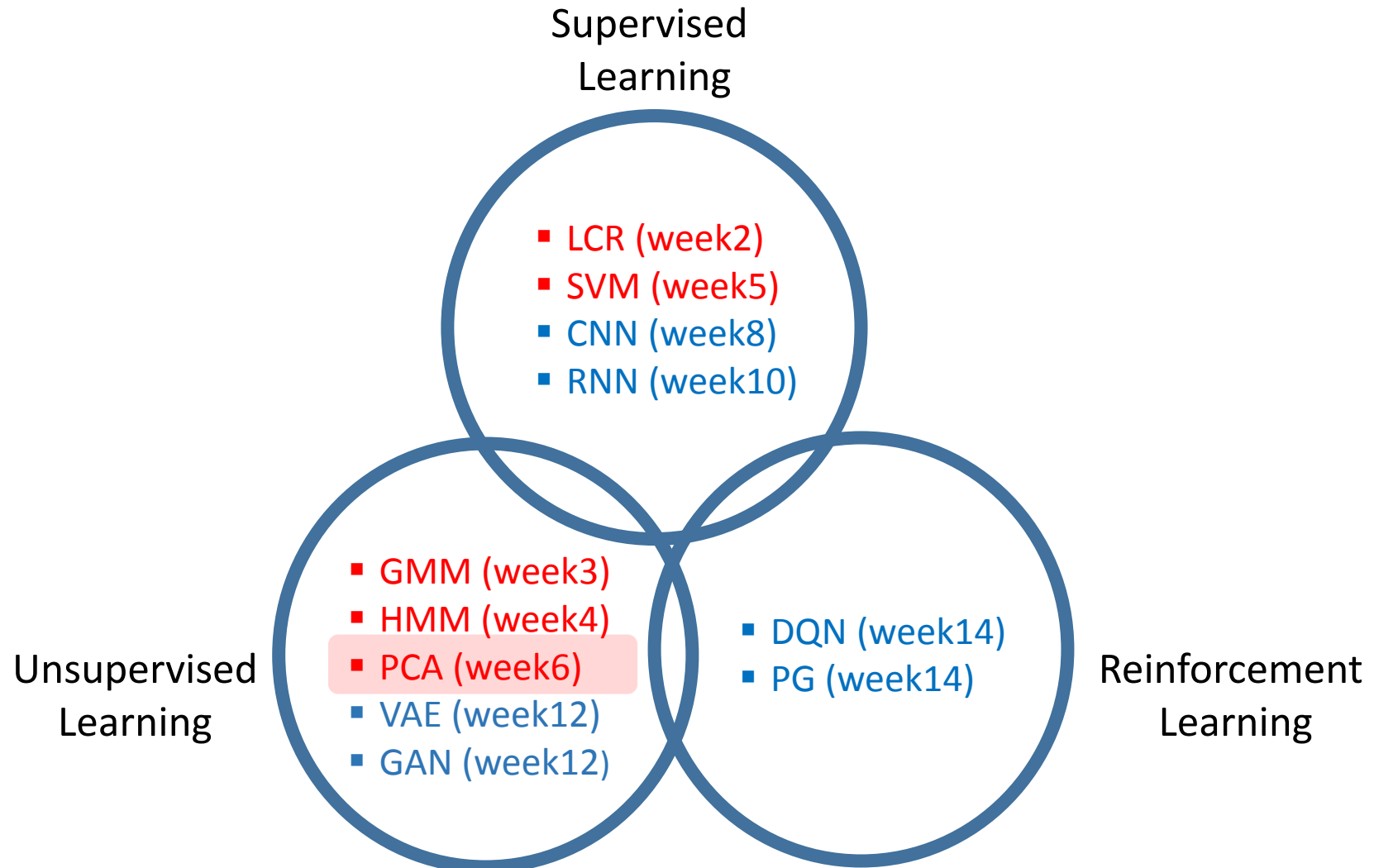
Lecture 6

Principal Components Analysis (PCA)

Dr. Suyong Eum



Where we are



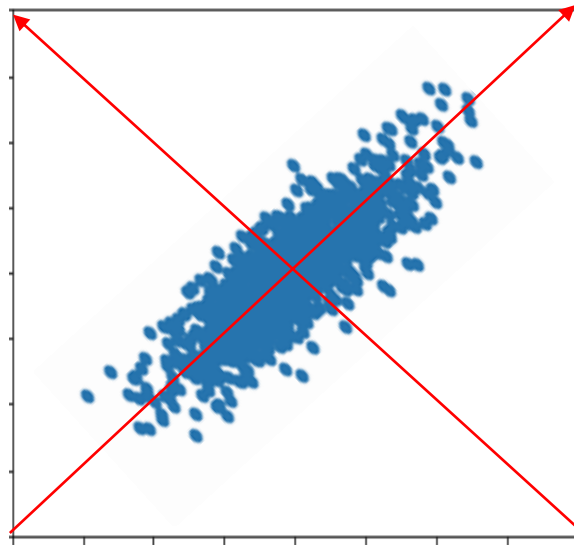
You are going to learn

- ❑ Why we need PCA
- ❑ How to obtain principal components
 - Eigen value decomposition and singular value decomposition
- ❑ SVD: data compression and visualization
- ❑ How to apply PCA for machine learning

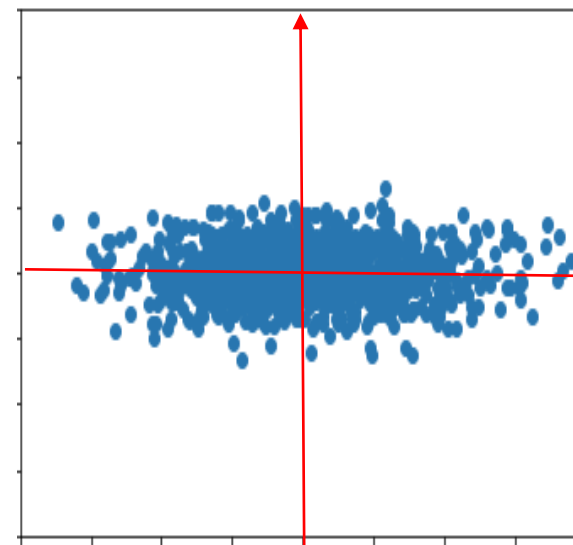
Principal Component Analysis (PCA): definition

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:

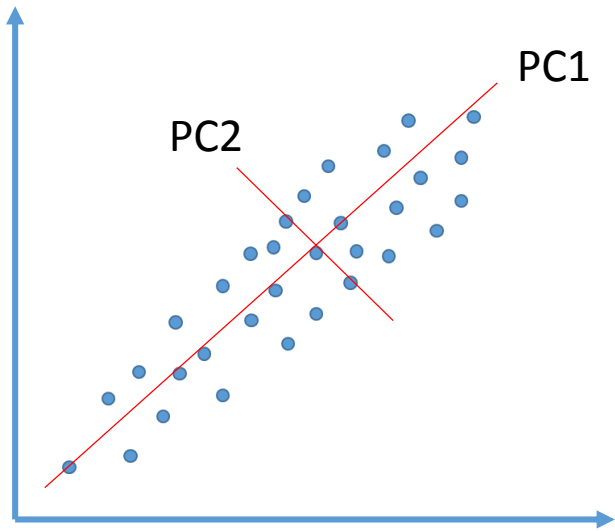


$$\begin{bmatrix} 8 & 1 \\ 1 & 8 \end{bmatrix}$$

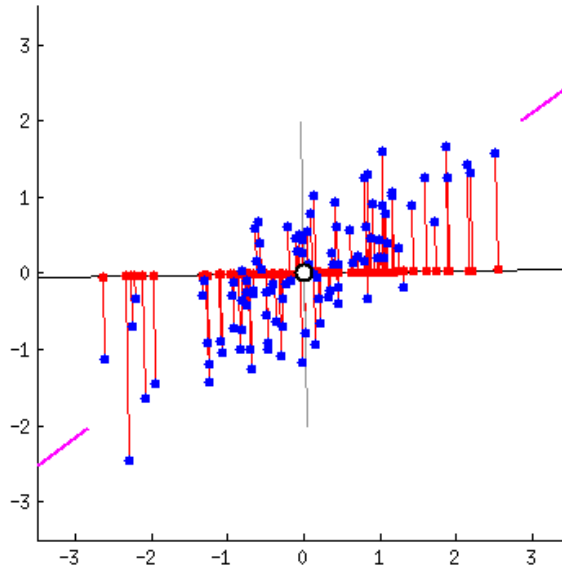


$$\begin{bmatrix} 10 & 0 \\ 0 & 1 \end{bmatrix}$$

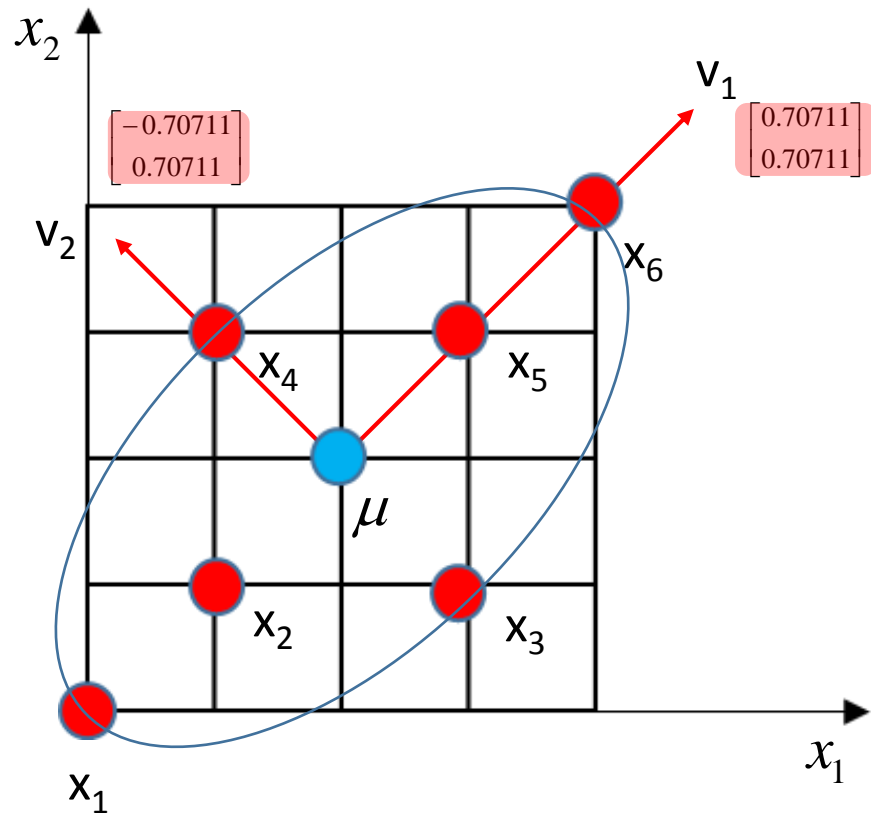
Principal Component Analysis (PCA): intuition



- ❑ 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?



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

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

$\text{cov}(X) = \begin{bmatrix} 2.4 & 1.6 \\ 1.6 & 2.4 \end{bmatrix}$
 variance along the axis of " x_1 "
 variance along the axis of " x_2 "

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

```
>> [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

Largest eigen value of the covariance matrix == largest variances in the data set?

- ❑ “ m_i ” shows the distance between 0 (mean) to the point where “ x_i ” is projected on the vector “ V ”.

$$m_i = x_i v$$

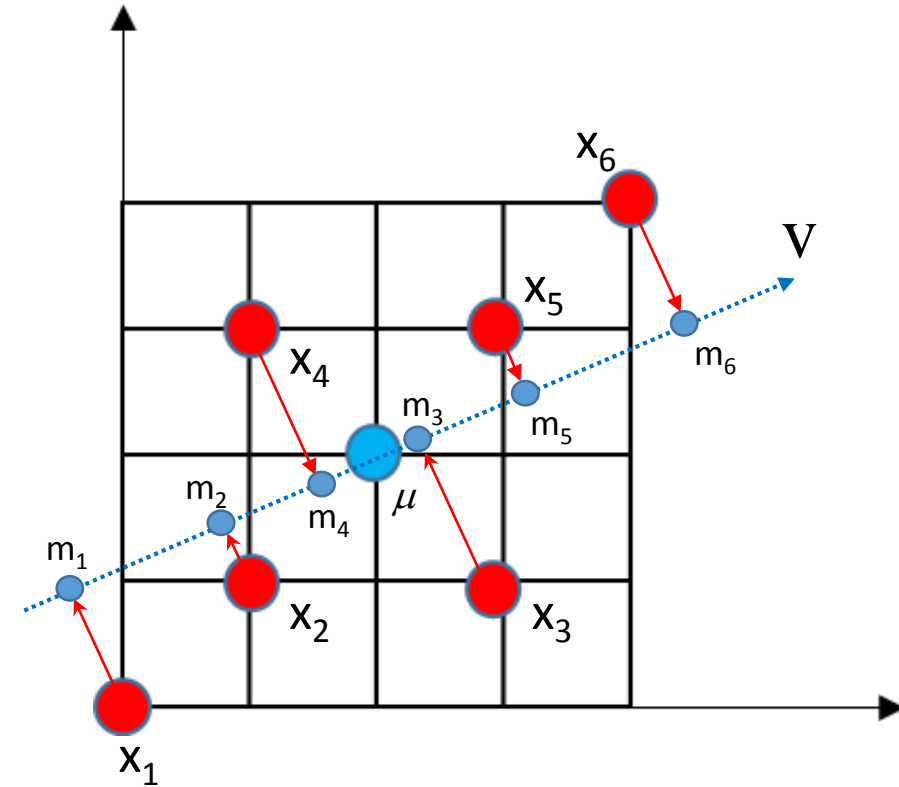
- ❑ Let's define the variance of data points “ m ”

$$\text{var}(m) = \frac{1}{N-1} \sum_{i=1}^N (m_i - \mu v)^2 = \frac{1}{N-1} \sum_{i=1}^N (x_i v)^2$$

- ❑ Let's maximize the variance with a constraint (v must be an unit vector). Then, see what it would be.

$$\max \frac{1}{N-1} \sum_{i=1}^N (x_i v)^2 \quad s.t. \quad \sum_{i=1}^d v_i^2 = 1$$

An Unit vector



Largest eigen value of the covariance matrix == largest variances in the data set?

- ❑ Let's convert the constrained problem to unconstrained problem using Lagrange method (again!).

$$L(v) = \frac{1}{N-1} \sum_{i=1}^N (x_i v)^2 - \lambda_i \left(\sum_{i=1}^d v_i^2 - 1 \right)$$

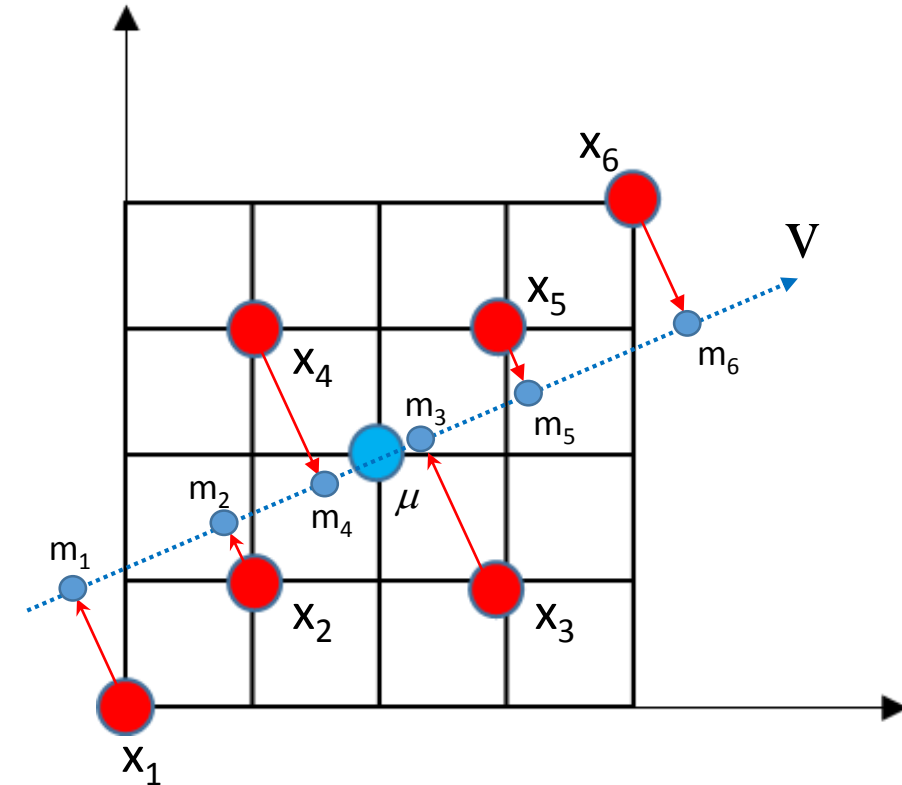
- ❑ We look for the vector “**v**” which maximizes the variance. Thus, differentiating the above with respect to “**v**”

$$\frac{\partial L(v)}{\partial v} = \frac{2}{N-1} \sum_{i=1}^N (x_i v)(x_i) - 2\lambda_i \left(\sum_{i=1}^d v_i \right) = 0$$

$$\left(\frac{1}{N-1} \sum_{i=1}^N (x_i)^2 \right) v = \lambda v$$

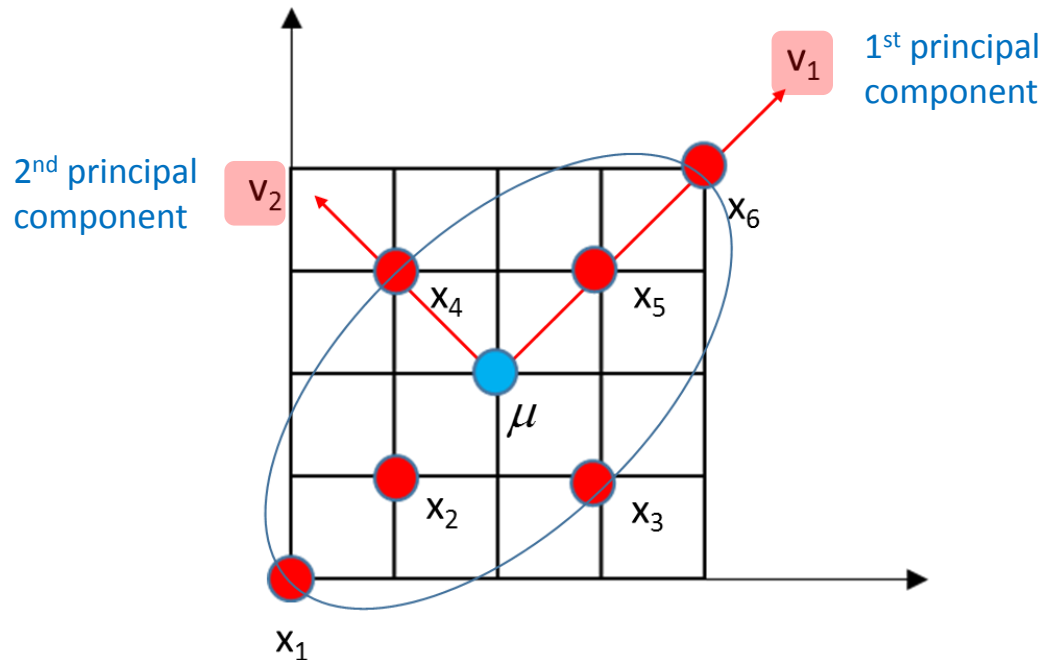
Covariance Eigen value

- When “v” is selected to maximize the variance, covariance matrix becomes equivalent to its own eigen value.
- Eigen value has diagonal elements, which represent variances along eigen vectors – no correlation.



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
```

Diagram illustrating the eigen decomposition of the covariance matrix. The eigen vectors (vec) are shown as columns of the matrix. The eigen values (val) are shown as diagonal elements of the Diagonal Matrix. The first principal component (V_1) corresponds to the largest eigen value (4.00000), and the second principal component (V_2) corresponds to the smallest eigen value (0.80000).

How to find the principal components showing the largest variance?

- ❑ Actually, there is a more convenient way of doing it (finding eigen vectors).
- ❑ 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 (finding eigen vectors).
- ❑ 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
```

$$\Lambda = \Sigma^2$$

$$4.4721^2 = 20$$

Singular Value Decomposition (SVD)

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

```
>> [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

V =
0.70711 -0.70711
0.70711 0.70711
```

Singular Value Decomposition (SVD): data compression

$$X \in \mathbf{R}^{n \times m}$$

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

$$U \in \mathbf{R}^{n \times n}$$

$$\Sigma \in \mathbf{R}^{n \times m}$$

$$V \in \mathbf{R}^{m \times m}$$

$$\begin{bmatrix} -2 & -2 \\ -1 & -1 \\ 1 & -1 \\ -1 & 1 \\ 1 & 1 \\ 2 & 2 \end{bmatrix} = \begin{bmatrix} -0.63246 & 0.00000 & 0.30819 & -0.30819 & 0.28637 & 0.57274 \\ -0.31623 & -0.00000 & -0.63635 & 0.63635 & 0.13426 & 0.26851 \\ 0.00000 & -0.70711 & 0.50000 & 0.50000 & 0.00000 & 0.00000 \\ -0.00000 & 0.70711 & 0.50000 & 0.50000 & -0.00000 & -0.00000 \\ 0.31623 & 0.00000 & -0.00399 & 0.00399 & 0.94140 & -0.11720 \\ 0.63246 & 0.00000 & -0.00799 & 0.00799 & -0.11720 & 0.76560 \end{bmatrix} \begin{bmatrix} 4.4721 & 0 \\ 0 & 2.0000 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} 0.70711 & -0.70711 \\ 0.70711 & 0.70711 \end{bmatrix}$$

X: 6x2

U: 6x6

Σ : 6x2

V: 2x2

Singular Value Decomposition (SVD): data compression

$$X \in \mathbf{R}^{n \times m}$$

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

$$U \in \mathbf{R}^{n \times n}$$

$$\Sigma \in \mathbf{R}^{n \times m}$$

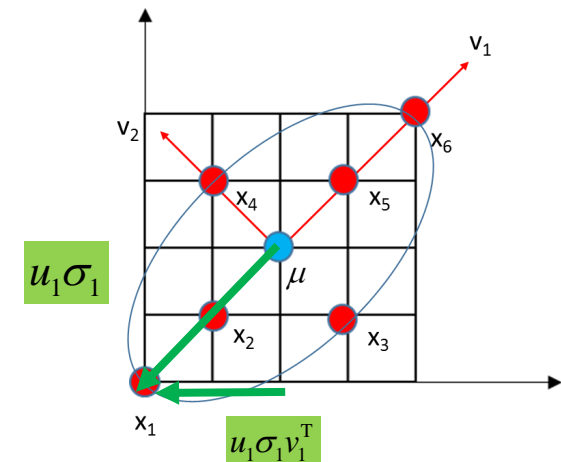
$$V \in \mathbf{R}^{m \times m}$$

$$\begin{bmatrix} -2 & -2 \\ -1 & -1 \\ 1 & -1 \\ -1 & 1 \\ 1 & 1 \\ 2 & 2 \end{bmatrix} = \begin{bmatrix} u_1 & & & & & \\ -0.63246 & 0.00000 & & & & \\ -0.31623 & -0.00000 & & & & \\ 0.00000 & -0.70711 & & & & \\ -0.00000 & 0.70711 & & & & \\ 0.31623 & 0.00000 & & & & \\ 0.63246 & 0.00000 & & & & \end{bmatrix} \begin{bmatrix} \sigma_1 & & & & & \\ 4.4721 & 0 & & & & \\ 0 & 2.0000 & & & & \end{bmatrix} \begin{bmatrix} v_1 & & & & & \\ 0.70711 & -0.70711 & & & & \\ 0.70711 & 0.70711 & & & & \end{bmatrix}$$

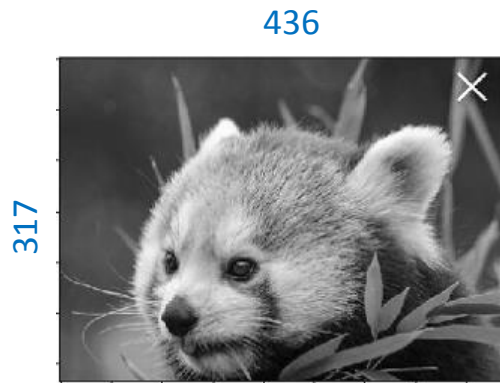
X: 6x2
U: 6x2
Σ: 2x2
V: 2x2

$$X = u_1 \sigma_1 v_1^T + u_2 \sigma_2 v_2^T$$

- New coordination system which has two basis (v1 and v2)



Singular Value Decomposition (SVD): data compression



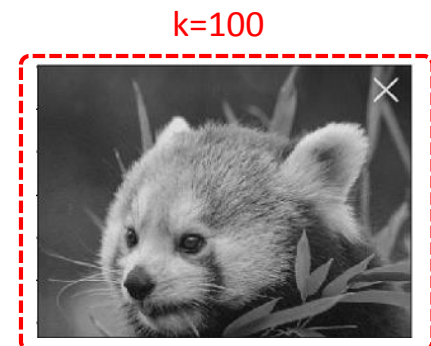
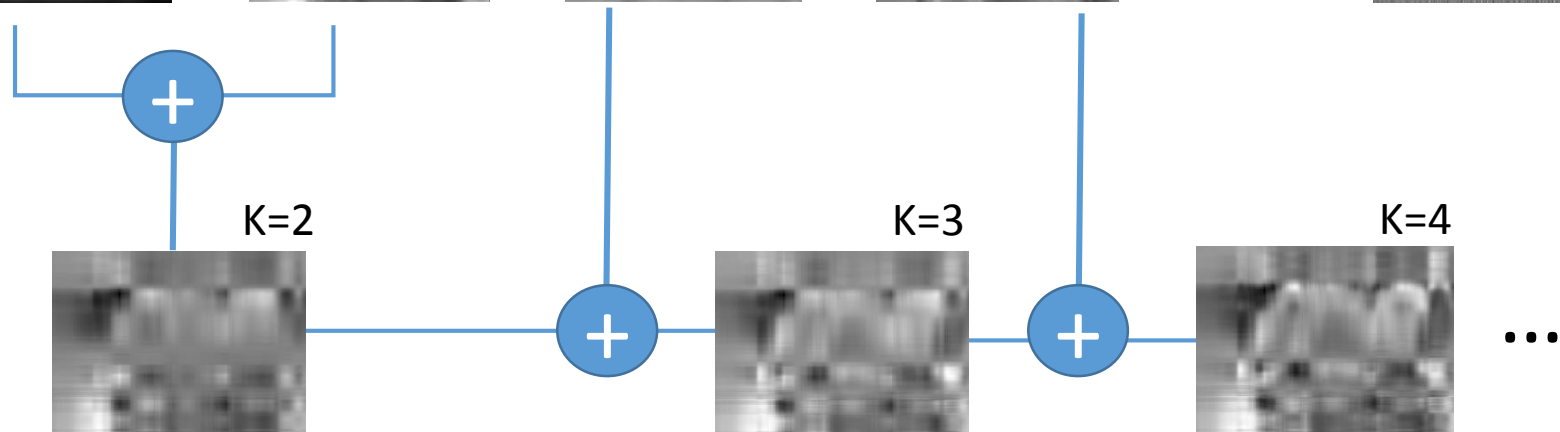
X: 317x436 = 138212

$$= u_1 \sigma_1 v_1^T + u_2 \sigma_2 v_2^T + u_3 \sigma_3 v_3^T + \dots + u_m \sigma_m v_m^T$$

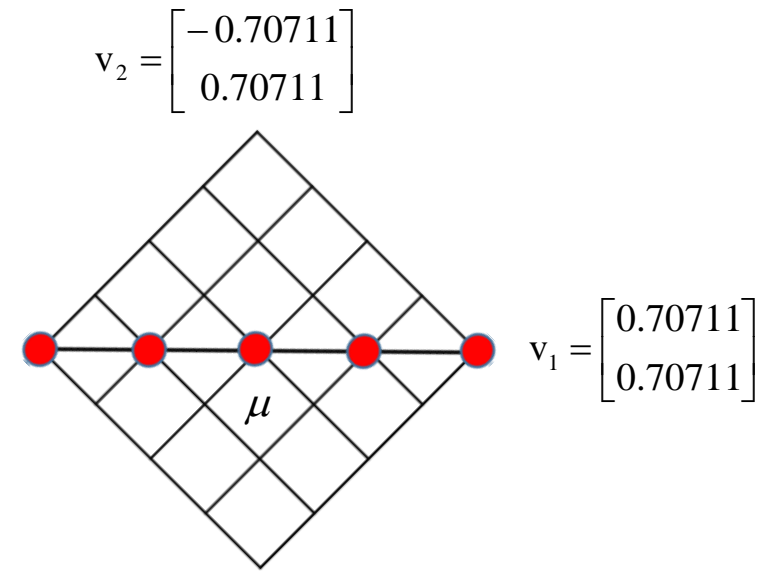
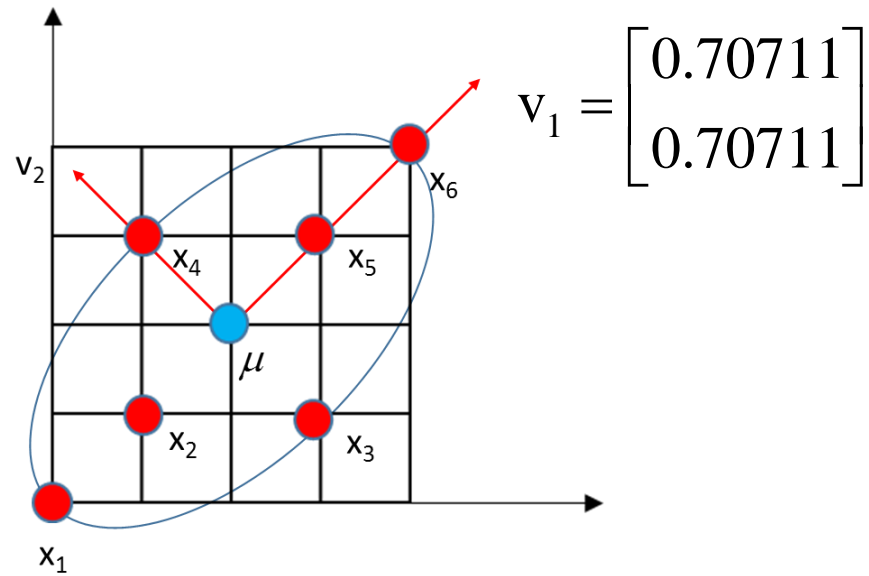
k=100

(317+1+436) x 100 = 75400

$$X = u_1 \sigma_1 v_1^T + u_2 \sigma_2 v_2^T + u_3 \sigma_3 v_3^T + u_4 \sigma_4 v_4^T + \dots + u_m \sigma_m v_m^T$$



Back to PCA: dimension reduction



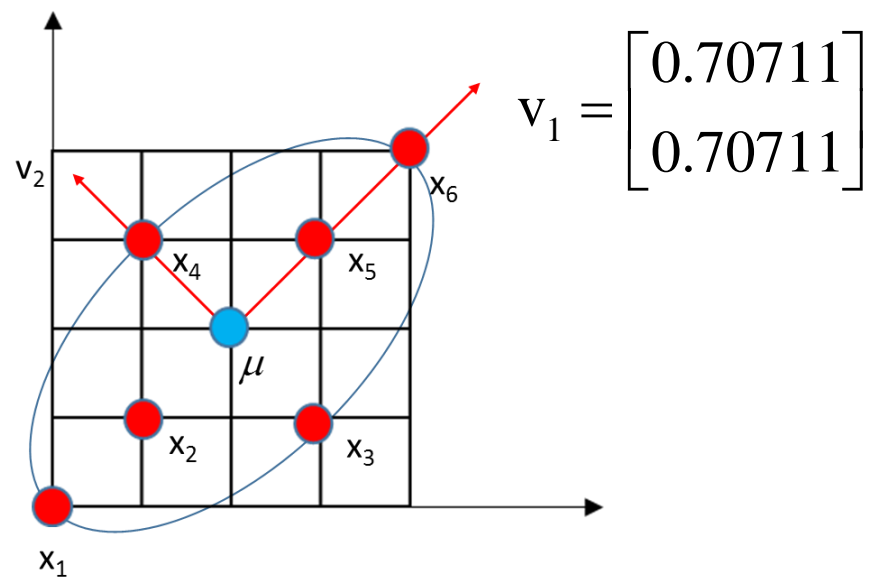
$$\begin{bmatrix} -2 & -2 \\ -1 & -1 \\ 1 & -1 \\ -1 & 1 \\ 1 & 1 \\ 2 & 2 \end{bmatrix} = \begin{bmatrix} u_1 & u_2 \end{bmatrix} \begin{bmatrix} \sigma_1 \\ \sigma_2 \end{bmatrix} \begin{bmatrix} v_1^T \\ v_2^T \end{bmatrix}$$

$$X = u_1 \sigma_1 v_1^T + u_2 \sigma_2 v_2^T$$

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

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

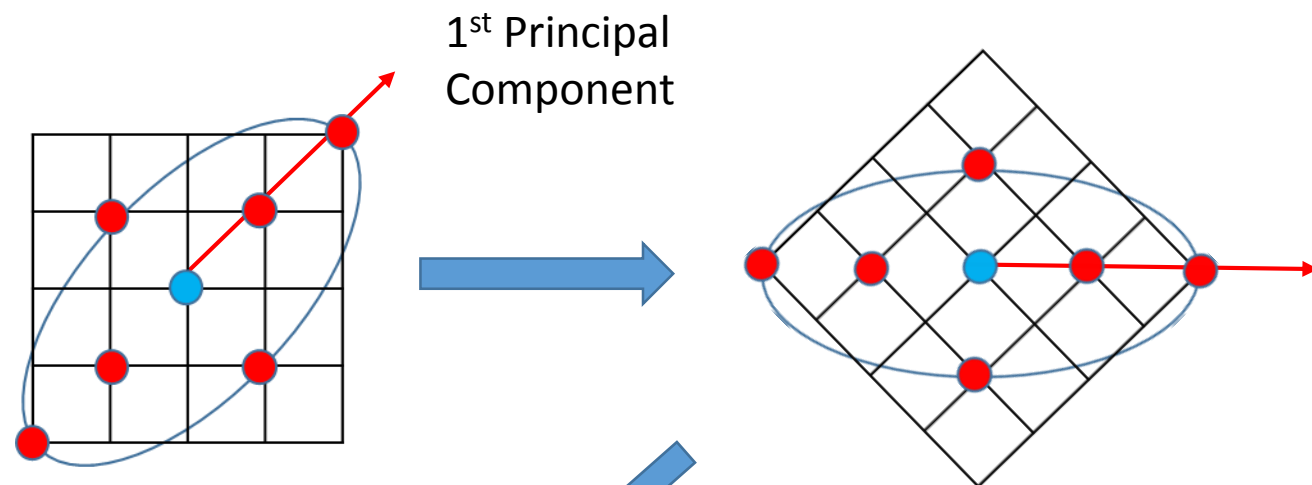
Back to PCA: dimension reduction



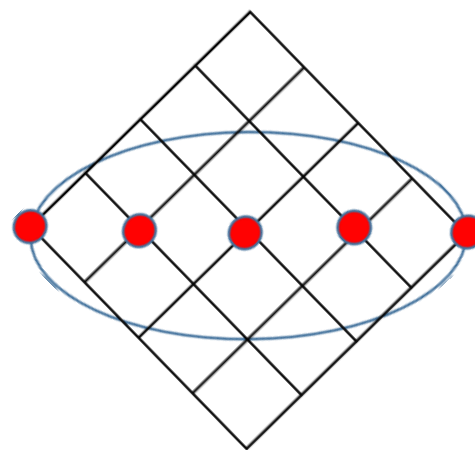
$$\begin{bmatrix} -2 & -2 \\ -1 & -1 \\ 1 & -1 \\ -1 & 1 \\ 1 & 1 \\ 2 & 2 \end{bmatrix} = \begin{bmatrix} u_1 & u_2 \end{bmatrix} \begin{bmatrix} \sigma_1 \\ \sigma_2 \end{bmatrix} \begin{bmatrix} v_1^T \\ v_2^T \end{bmatrix}$$

$$X = u_1 \sigma_1 v_1^T + u_2 \sigma_2 v_2^T$$

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

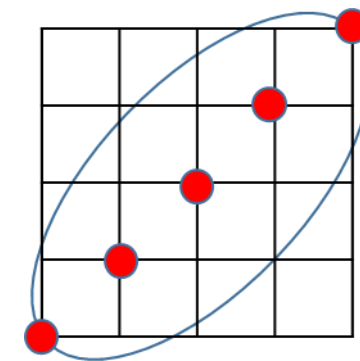


$$X_{\text{rot}} = X \cdot V$$



Set the " v_2 " into zero

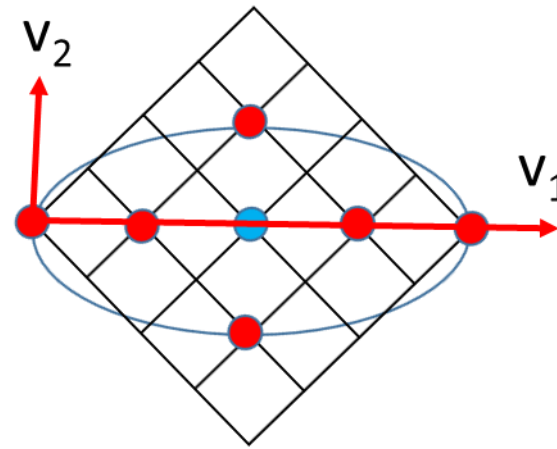
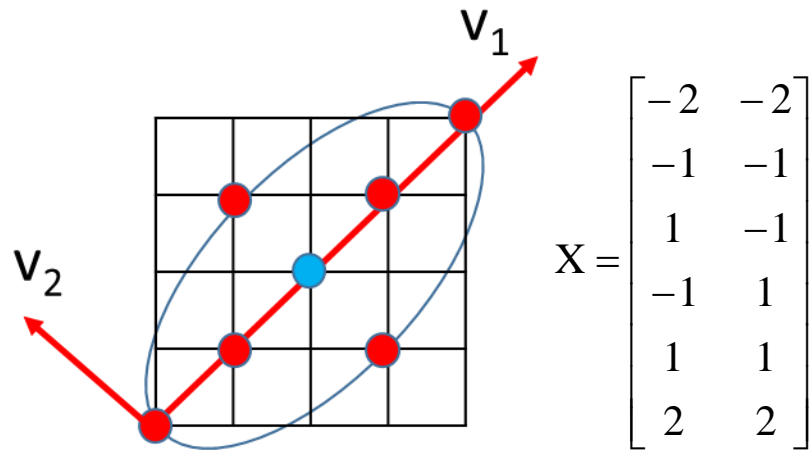
$$X_{\text{rot_zero}} = X_{\text{rot}} \cdot V^{-1}$$



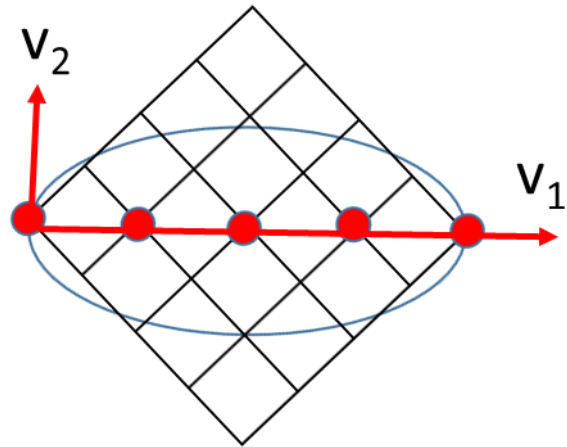
$$X' = X_{\text{rot_zero}} \cdot V^{-1}$$

$$= X_{\text{rot_zero}} \cdot V^T$$

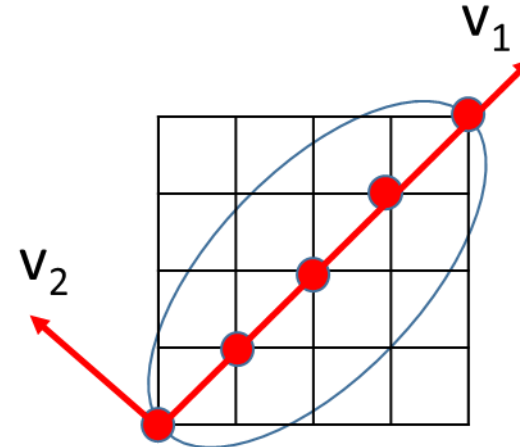
Back to PCA: example



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



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

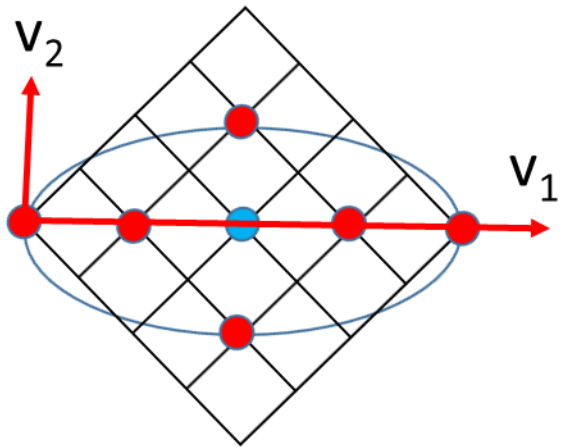
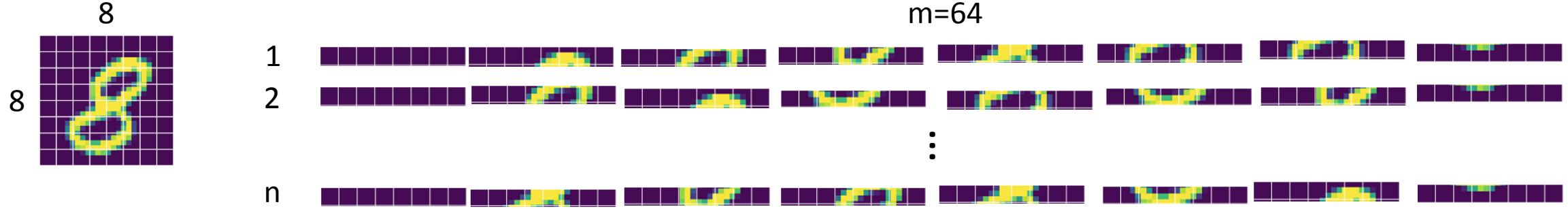


$$\begin{bmatrix} -2 & -2 \\ -1 & -1 \\ 0 & 0 \\ 0 & 0 \\ 1 & 1 \\ 2 & 2 \end{bmatrix}$$

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

Set the " v_2 " elements
into zero

How to use PCA for machine learning?



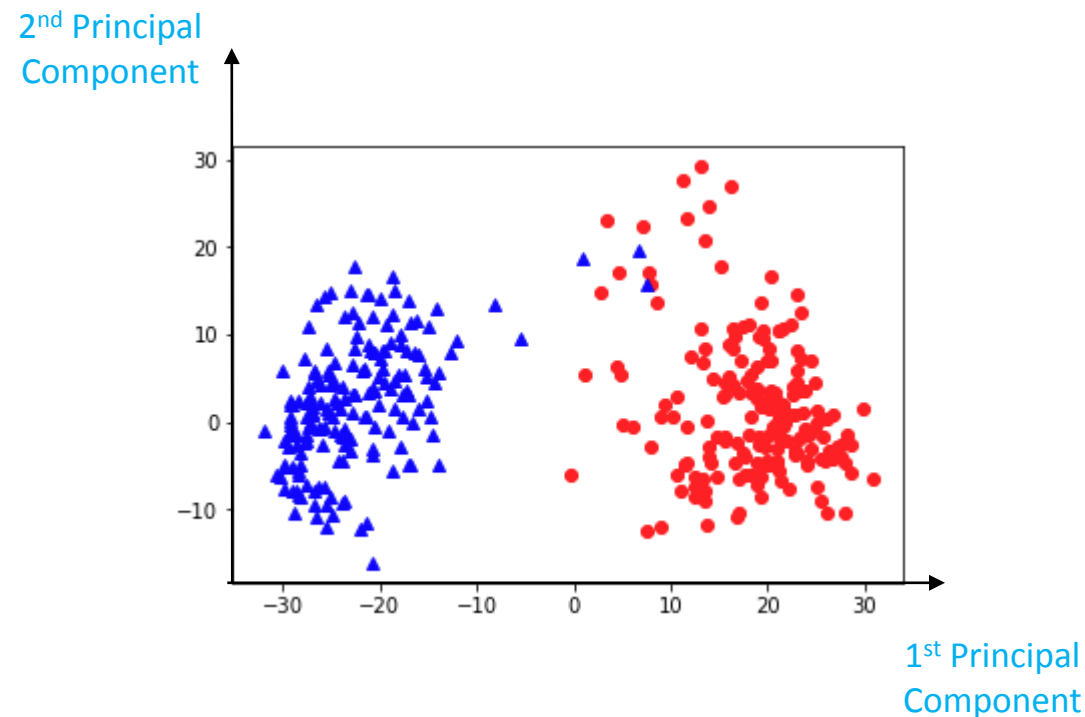
$$X_{\text{rot}} = X \cdot V$$

$$\begin{bmatrix} \overset{v_1}{a_{11}} & \overset{v_2}{a_{12}} & \cdots & a_{1m} \\ \vdots & \vdots & \vdots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nm} \end{bmatrix}$$

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

How to use PCA for machine learning?

- ❑ 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.



One that you need to be careful when carrying out PCA

- ❑ Centering the data before applying PCA.
- ❑ Normalizing or standardizing the data when features have different scale.

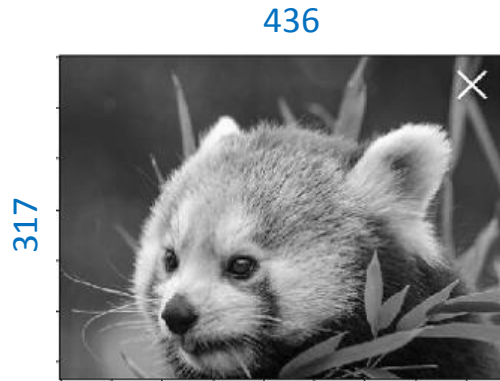
	population	area
Country 1	$5 \cdot 10^7$	92
Country 2	$2 \cdot 10^7$	74
...
Country n	$5 \cdot 10^8$	150

$$\begin{array}{l} \text{Data 1} \\ \text{Data 2} \\ \vdots \\ \text{Data n} \end{array} \begin{array}{cccc} \text{feature1} & \text{feature2} & \cdots & \text{feature m} \\ \left[\begin{array}{cccc} a_{11} & a_{12} & \cdots & a_{1n} \\ \vdots & \vdots & \vdots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{array} \right] \end{array}$$

That needs to be normalized

Backup Slides

Singular Value Decomposition (SVD): data compression



X: 317x436 = 138212

$$= u_1 \sigma_1 v_1^T + u_2 \sigma_2 v_2^T + u_3 \sigma_3 v_3^T + \cdots + u_m \sigma_m v_m^T$$

k=100

m=317

(317+1+436) x 100 = 75400



```
img = imread("sample_BW.png")[:,:]  
imshow(img)  
show()
```

Top 100

```
U,S,Vt = svd(img)
```

```
S = resize(S, [m,1])*eye(m,n)
```

```
imshow(dot(U[:,0:100], dot(S[0:100,0:100], Vt[0:100,:])))
```

```
show()
```