

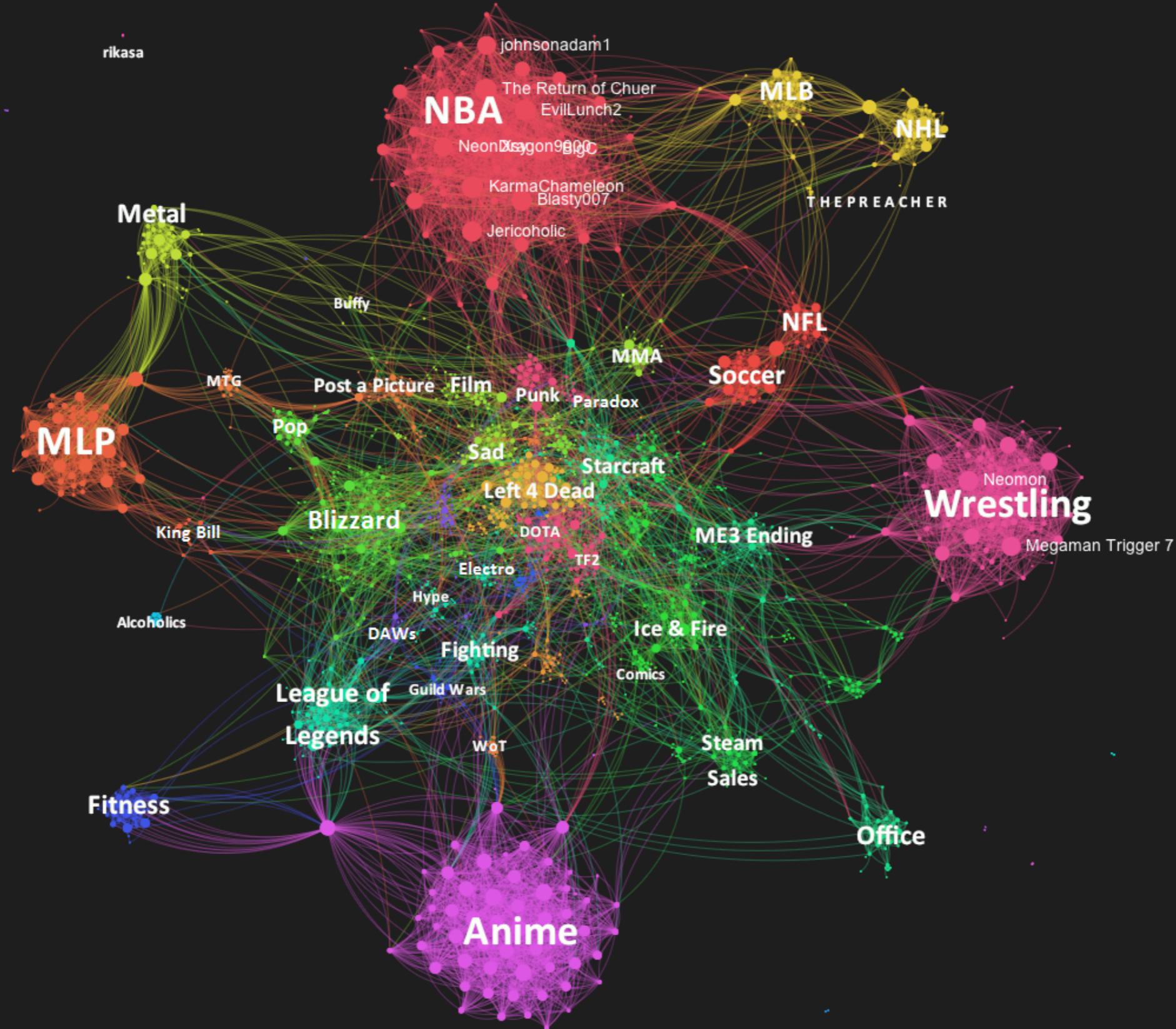
# CSE 6740: Computational Data Analysis

## Spring 2026

**Spectral clustering,  
dimensionality reduction &  
principle component analysis**

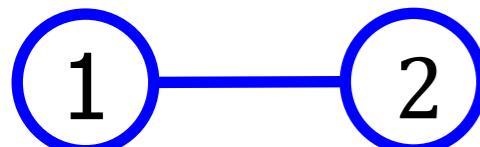
Anqi Wu  
01/22

# How about clustering nodes in social networks



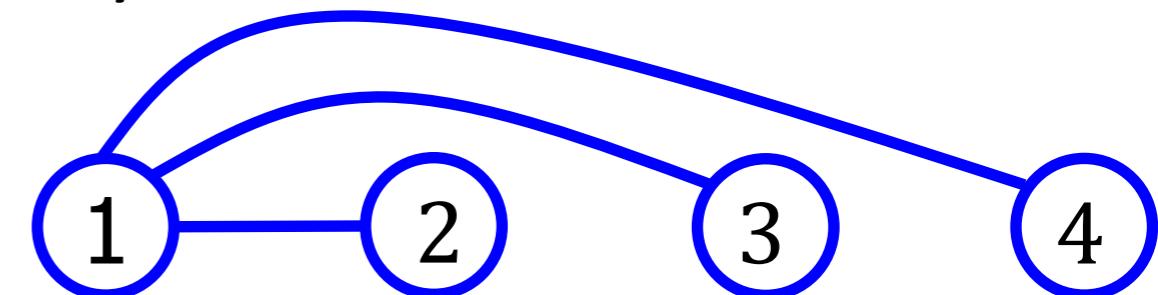
# Spectral clustering algorithm

- Step 1: represent graph as adjacency matrix  $A \in R^{m \times m}$



$$A = \begin{pmatrix} 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \end{pmatrix}$$

$$D = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}$$



$$A = \begin{pmatrix} 0 & 1 & 1 & 1 \\ 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \end{pmatrix}$$

$$D = \begin{pmatrix} 3 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}$$

- Step 2: form a special matrix  $L = D - A$ , the graph Laplacian

$$L = \begin{pmatrix} 1 & -1 & 0 & 0 \\ -1 & 1 & 0 & 0 \\ 0 & 0 & 1 & -1 \\ 0 & 0 & -1 & 1 \end{pmatrix}$$

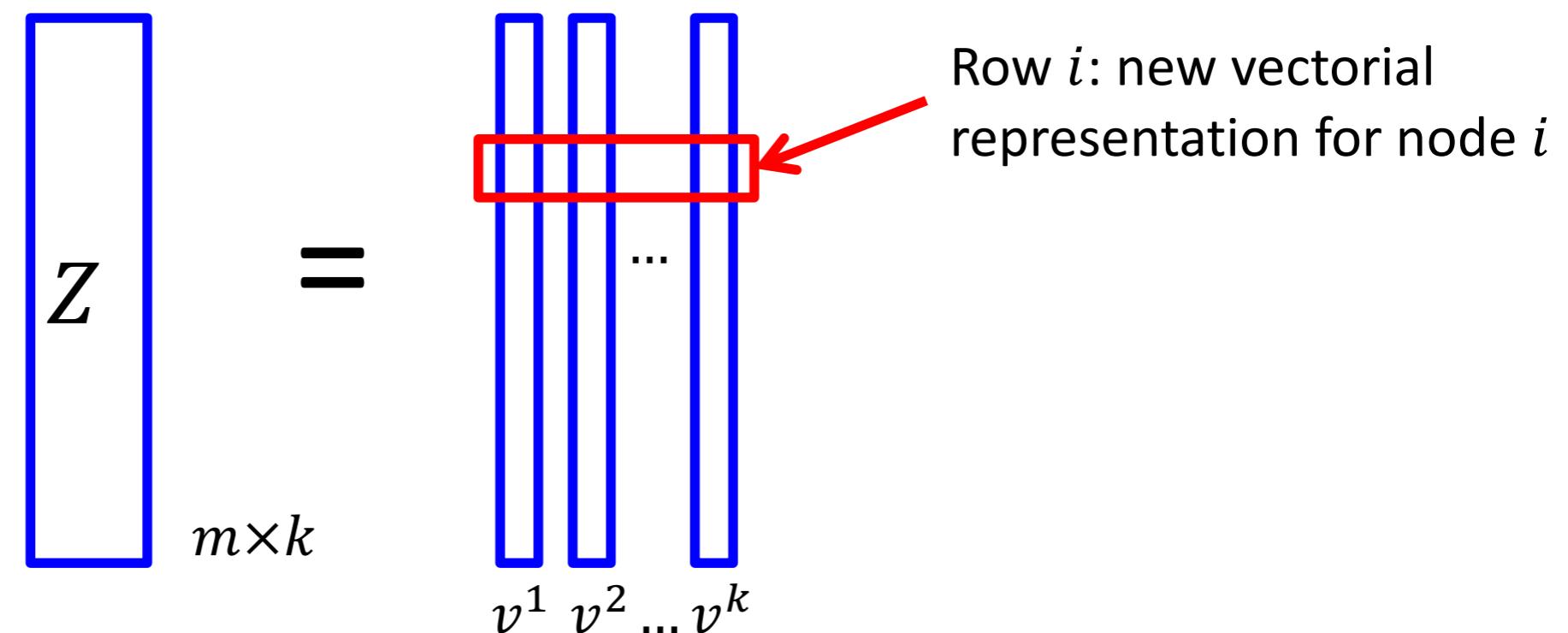
$$L = \begin{pmatrix} 3 & -1 & -1 & -1 \\ -1 & 1 & 0 & 0 \\ -1 & 0 & 1 & 0 \\ -1 & 0 & 0 & 1 \end{pmatrix}$$

# Spectral clustering algorithm (cont.)

- Step 3: compute  $k$  eigenvectors,  $v^1, v^2, \dots, v^k$ , of  $L$  corresponding to the  $k$  **smallest** eigenvalues ( $k \ll m$ )

$$Lv^1 = \begin{pmatrix} 1 & -1 & 0 & 0 \\ -1 & 1 & 0 & 0 \\ 0 & 0 & 1 & -1 \\ 0 & 0 & -1 & 1 \end{pmatrix} v^1 \stackrel{?}{=} \lambda_1 v^1$$

- Step 4: run kmeans algorithm on  $Z = (v^1, v^2, \dots, v^k)$  by treating each row as a new data point



# Why Laplacian matrix L?

## Step 1: Clustering as a Cut Problem on a Graph

We want to partition nodes into two groups  $S$  and  $\bar{S}$ . For a weighted graph with edge weights  $w_{ij}$ , define the cut:

$$\text{Cut}(S, \bar{S}) = \sum_{i \in S, j \in \bar{S}} w_{ij}.$$

Goal: make cross-group connections as weak as possible.

# Why Laplacian matrix L?

## Step 2: Encode the Partition with a Discrete Indicator

Define a discrete indicator vector  $x \in \{+1, -1\}^n$ :

$$x_i = \begin{cases} +1 & i \in S \\ -1 & i \in \bar{S} \end{cases}$$

Then for any pair  $(i, j)$ :

$$(x_i - x_j)^2 = \begin{cases} 0 & \text{same group} \\ 4 & \text{different groups} \end{cases}$$

So  $(x_i - x_j)^2$  acts like a “soft switch” that detects boundary edges.

# Why Laplacian matrix L?

## Step 3: A Cut-Counting Objective (No $L$ Yet)

Consider the objective

$$J(x) = \frac{1}{2} \sum_{i,j} w_{ij} (x_i - x_j)^2.$$

If  $x \in \{+1, -1\}^n$ , only cross-group pairs contribute:

$$J(x) = \frac{1}{2} \sum_{(i,j) \in \text{cut}} w_{ij} \cdot 4 = 2 \text{Cut}(S, \bar{S}).$$

Thus, for discrete indicators,

$$\min_{x \in \{+1, -1\}^n} J(x) \iff \min \text{Cut}(S, \bar{S}).$$

# Why Laplacian matrix L?

## Step 4: Why Relax to Continuous $x$ ?

The discrete optimization

$$\min_{x \in \{+1, -1\}^n} J(x)$$

is combinatorial and hard.

Spectral methods relax the constraint:

$$x \in \mathbb{R}^n,$$

while keeping the same objective form  $J(x)$ .

We also need constraints to avoid trivial solutions:

$$\|x\|^2 = 1, \quad x \perp \mathbf{1}.$$

# Why Laplacian matrix $L$ ?

Step 5: Now Introduce the Laplacian

Let  $W$  be the weighted adjacency matrix, and  $D$  be the degree matrix:

$$D_{ii} = \sum_j w_{ij}, \quad L = D - W.$$

A key identity shows the objective can be written compactly:

$$\frac{1}{2} \sum_{i,j} w_{ij} (x_i - x_j)^2 = x^\top L x.$$

So  $J(x) = x^\top L x$ . This is not a new objective: it is exactly the same cut-counting form.

# Why Laplacian matrix $L$ ?

## Step 6: Continuous Minimization Leads to Eigenvectors

We solve the relaxed problem:

$$\min_{x \in \mathbb{R}^n} x^\top L x \quad \text{s.t.} \quad \|x\|^2 = 1, \quad x \perp \mathbf{1}.$$

The solution is the eigenvector corresponding to the smallest eigenvalue that is not the trivial constant eigenvector. This is the second smallest eigenvector of  $L$ .

# Why smallest eigenvalue?

## Problem Statement

Let  $L = L^\top \in \mathbb{R}^{n \times n}$  be a symmetric matrix. We want to solve

$$\min_{\|x\|=1} x^\top L x.$$

Claim:

$$\min_{\|x\|=1} x^\top L x = \lambda_1,$$

where  $\lambda_1$  is the smallest eigenvalue of  $L$ , and the minimizer is an eigenvector associated with  $\lambda_1$ .

# Why smallest eigenvalue?

## Step A: Spectral Decomposition

Because  $L$  is symmetric, it is orthogonally diagonalizable:

$$L = U\Lambda U^\top, \quad U = [u_1, \dots, u_n], \quad \Lambda = \text{diag}(\lambda_1, \dots, \lambda_n),$$

with

$$U^\top U = I, \quad \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n.$$

# Why smallest eigenvalue?

Step B: Expand  $x$  in the Eigenbasis

Write any  $x \in \mathbb{R}^n$  as

$$x = \sum_{i=1}^n \alpha_i u_i.$$

Using orthonormality,

$$\|x\|^2 = x^\top x = \sum_{i=1}^n \alpha_i^2.$$

Thus the constraint  $\|x\| = 1$  is equivalent to

$$\sum_{i=1}^n \alpha_i^2 = 1.$$

# Why smallest eigenvalue?

## Step C: Quadratic Form Becomes a Weighted Sum

Compute:

$$x^\top L x = \left( \sum_i \alpha_i u_i \right)^\top L \left( \sum_j \alpha_j u_j \right) = \sum_{i,j} \alpha_i \alpha_j u_i^\top L u_j.$$

Since  $L u_j = \lambda_j u_j$  and  $u_i^\top u_j = \delta_{ij}$ ,

$$u_i^\top L u_j = u_i^\top (\lambda_j u_j) = \lambda_j \delta_{ij}.$$

Therefore,

$$x^\top L x = \sum_{i=1}^n \lambda_i \alpha_i^2.$$

So the original problem is equivalent to

$$\min_{\sum_i \alpha_i^2 = 1} \sum_{i=1}^n \lambda_i \alpha_i^2.$$

# Why smallest eigenvalue?

## Step D: Lower Bound via an Inequality (Method 1)

Because  $\lambda_i \geq \lambda_1$  for all  $i$  and  $\alpha_i^2 \geq 0$ ,

$$\sum_{i=1}^n \lambda_i \alpha_i^2 \geq \sum_{i=1}^n \lambda_1 \alpha_i^2 = \lambda_1 \sum_{i=1}^n \alpha_i^2 = \lambda_1.$$

Thus for every feasible  $x$  with  $\|x\| = 1$ ,

$$x^\top L x \geq \lambda_1.$$

So  $\lambda_1$  is a universal lower bound on the objective value.

# Why smallest eigenvalue?

Step D (continued):  $\lambda_1$  Is the Achievable Minimum

We now show the lower bound is achieved.

Take  $x = u_1$  (equivalently  $\alpha_1 = 1$  and  $\alpha_{i \neq 1} = 0$ ), which satisfies  $\|x\| = 1$ . Then

$$x^\top L x = u_1^\top L u_1 = u_1^\top (\lambda_1 u_1) = \lambda_1.$$

So the minimum value is at most  $\lambda_1$ , and together with the lower bound we conclude

$$\min_{\|x\|=1} x^\top L x = \lambda_1.$$

# Why smallest eigenvalue?

## Extra Argument by Contradiction

Assume  $\lambda_1$  is not the minimum value of the optimization. Then there exists a smaller value  $\lambda_1^* < \lambda_1$  achieved by some feasible  $x^*$  with  $\|x^*\| = 1$ :

$$(x^*)^\top L x^* = \lambda_1^* < \lambda_1.$$

But Step D already proved that for every feasible  $x$ ,

$$x^\top L x \geq \lambda_1.$$

This implies  $(x^*)^\top L x^* \geq \lambda_1$ , a contradiction. Therefore  $\lambda_1$  must be the achievable minimum.

# Why smallest eigenvalue?

## Minimizer Characterization

A minimizer can be chosen as an eigenvector corresponding to  $\lambda_1$ :

$$x^* = \pm u_1.$$

# Why smallest eigenvalue?

## One-Sentence Add-on: Multiplicity Case

If the smallest eigenvalue has multiplicity  $m$ ,

$$\lambda_1 = \lambda_2 = \cdots = \lambda_m < \lambda_{m+1},$$

then every unit vector in the subspace

$$\text{span}\{u_1, \dots, u_m\}$$

achieves the same minimum value  $\lambda_1$ .

# Graph Laplacian example



$$A = \begin{pmatrix} 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \end{pmatrix}$$

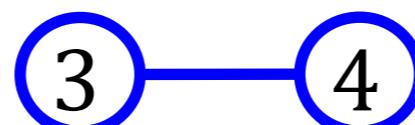
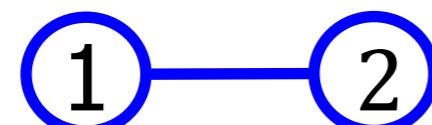
$$D = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}$$

$$L = \begin{pmatrix} 1 & -1 & 0 & 0 \\ -1 & 1 & 0 & 0 \\ 0 & 0 & 1 & -1 \\ 0 & 0 & -1 & 1 \end{pmatrix}$$

# Property I of Graph Laplacian

- $L = D - A$
- The multiplicity of the eigenvalue 0 corresponds to the number of connected components in the graph

• Example



$$L = \begin{pmatrix} 1 & -1 & 0 & 0 \\ -1 & 1 & 0 & 0 \\ 0 & 0 & 1 & -1 \\ 0 & 0 & -1 & 1 \end{pmatrix}$$

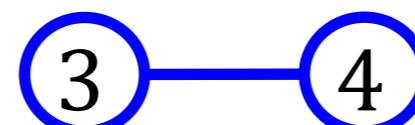
$$Lv_1 = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 & -1 & 0 & 0 \\ -1 & 1 & 0 & 0 \\ 0 & 0 & 1 & -1 \\ 0 & 0 & -1 & 1 \end{pmatrix} \begin{pmatrix} 1 \\ 1 \\ 0 \\ 0 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}$$

$$Lv_2 = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 & -1 & 0 & 0 \\ -1 & 1 & 0 & 0 \\ 0 & 0 & 1 & -1 \\ 0 & 0 & -1 & 1 \end{pmatrix} \begin{pmatrix} 0 \\ 0 \\ 1 \\ 1 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}$$

# Property II of Graph Laplacian

- $L = D - A$
- The eigenvectors with eigenvalue 0 contains cluster assignment information

- Example

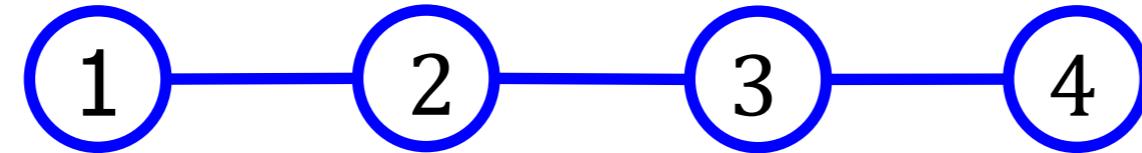


$$L = \begin{pmatrix} 1 & -1 & 0 & 0 \\ -1 & 1 & 0 & 0 \\ 0 & 0 & 1 & -1 \\ 0 & 0 & -1 & 1 \end{pmatrix}$$

$$Lv^1 = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 & -1 & 0 & 0 \\ -1 & 1 & 0 & 0 \\ 0 & 0 & 1 & -1 \\ 0 & 0 & -1 & 1 \end{pmatrix} \begin{pmatrix} 1 \\ 1 \\ 0 \\ 0 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}$$

$$Lv^2 = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 & -1 & 0 & 0 \\ -1 & 1 & 0 & 0 \\ 0 & 0 & 1 & -1 \\ 0 & 0 & -1 & 1 \end{pmatrix} \begin{pmatrix} 0 \\ 0 \\ 1 \\ 1 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}$$

# What if the graph has only 1 component



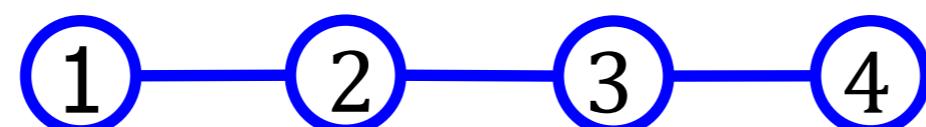
$$A = \begin{pmatrix} 0 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \\ 0 & 0 & 1 & 0 \end{pmatrix}$$

$$D = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 \\ 0 & 0 & 2 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}$$

$$L = \begin{pmatrix} 1 & -1 & 0 & 0 \\ -1 & 2 & -1 & 0 \\ 0 & -1 & 2 & -1 \\ 0 & 0 & -1 & 1 \end{pmatrix}$$

# Properties of Graph Laplacian

- $L = D - A$
- The smallest eigenvalue of  $L$  is 0, corresponding a constant eigenvector  $\frac{1}{\sqrt{m}} \mathbf{1}$
- Example



$$L = \begin{pmatrix} 1 & -1 & 0 & 0 \\ -1 & 2 & -1 & 0 \\ 0 & -1 & 2 & -1 \\ 0 & 0 & -1 & 1 \end{pmatrix}$$

$$\frac{1}{\sqrt{4}} L \mathbf{1} = \frac{1}{\sqrt{m}} \begin{pmatrix} 1 & -1 & 0 & 0 \\ -1 & 2 & -1 & 0 \\ 0 & -1 & 2 & -1 \\ 0 & 0 & -1 & 1 \end{pmatrix} \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}$$

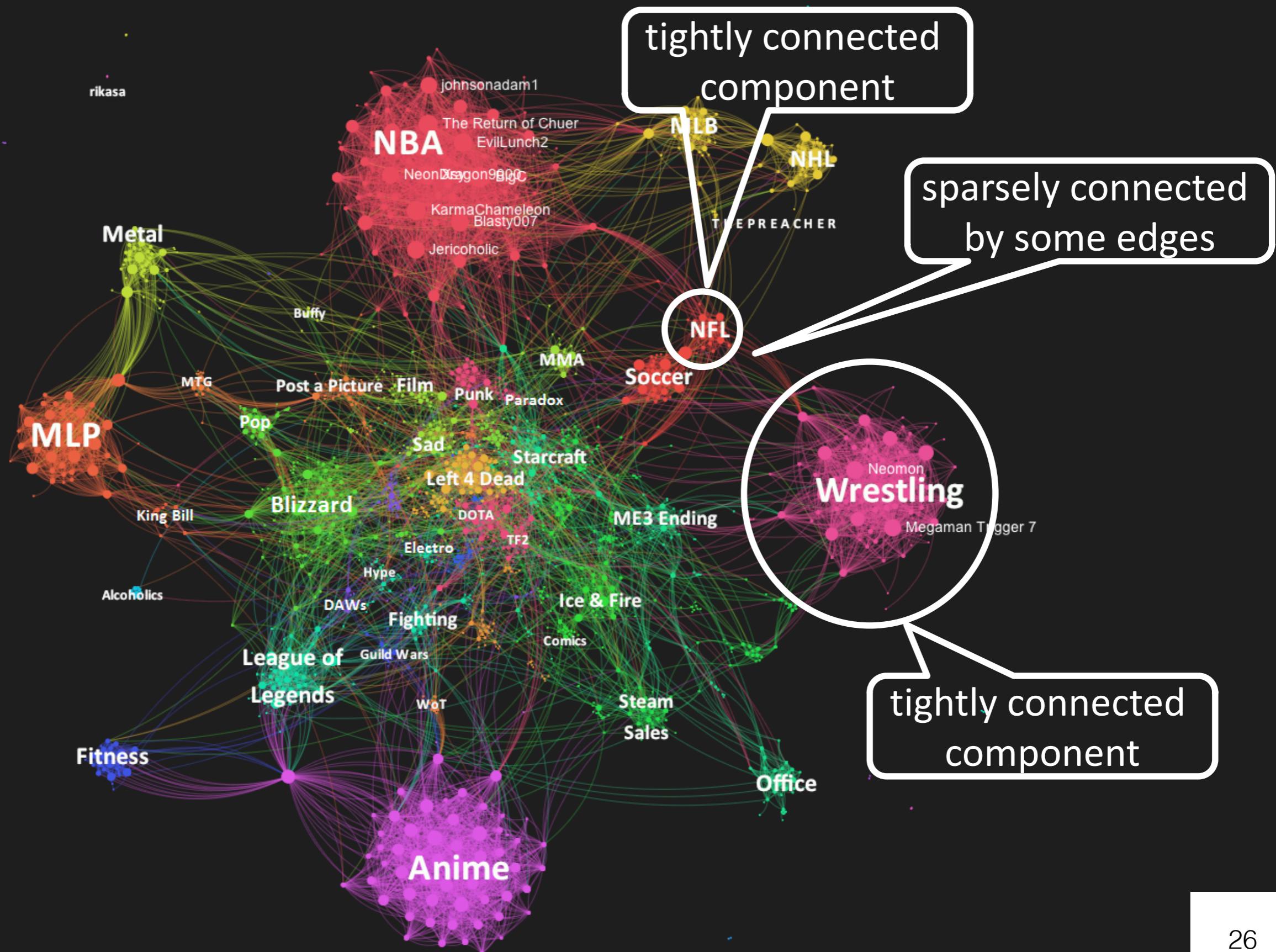
# What if the graph has $k$ components

- If a graph has  $k$  connected components (or  $k$  clusters)
- The graph Laplacian has  $k$  blocks

$$L = \begin{pmatrix} L_1 & 0 & 0 & 0 \\ 0 & L_2 & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & L_k \end{pmatrix}$$

- The graph Laplacian has  $k$  eigenvectors with zero eigenvalues
  - Eigenvector 1 is constant in block 1, but 0 in other blocks;  
eigenvector 2 is constant in block 2, but 0 in other blocks;
- ...

# Real world social networks



# In most real networks

- If a graph has  $k$  **tightly** connected components (or  $k$  clusters) with **sparsely** connected edges
- The graph Laplacian has **approximately**  $k$  blocks
- The graph Laplacian has  $k$  eigenvectors with **small** eigenvalues
- Eigenvector 1 is **approximately** constant in block 1, but 0 in other blocks; eigenvector 2 ...

# High level idea of spectral clustering

- Examine the properties of graph Laplacian for the perfect cases
  - The number of 0 eigenvalues corresponds to the number of connected components
  - Eigenvectors correspond to cluster assignment
- Then use the intuition from perfect cases to design algorithms for the imperfect case.
  - Eigenvectors no longer correspond exactly cluster indicator
  - Perform post processing to obtain cluster assignment

# Summary of spectral clustering

- Step 1: represent graph as adjacency matrix  $A \in R^{m \times m}$
- Step 2: form a special matrix  $L = D - A$ , the graph Laplacian
- Step 3: compute  $k$  eigenvectors,  $v^1, v^2, \dots, v^k$ , of  $L$  corresponding to the  $k$  **smallest** eigenvalues ( $k \ll m$ )
- Step 4: run kmeans algorithm on  $Z = (v^1, v^2, \dots, v^k)$  by treating each row as a new data point

$$Z_1 = (V_{11}, V_{21})$$

$$Z_2 = (V_{12}, V_{22})$$

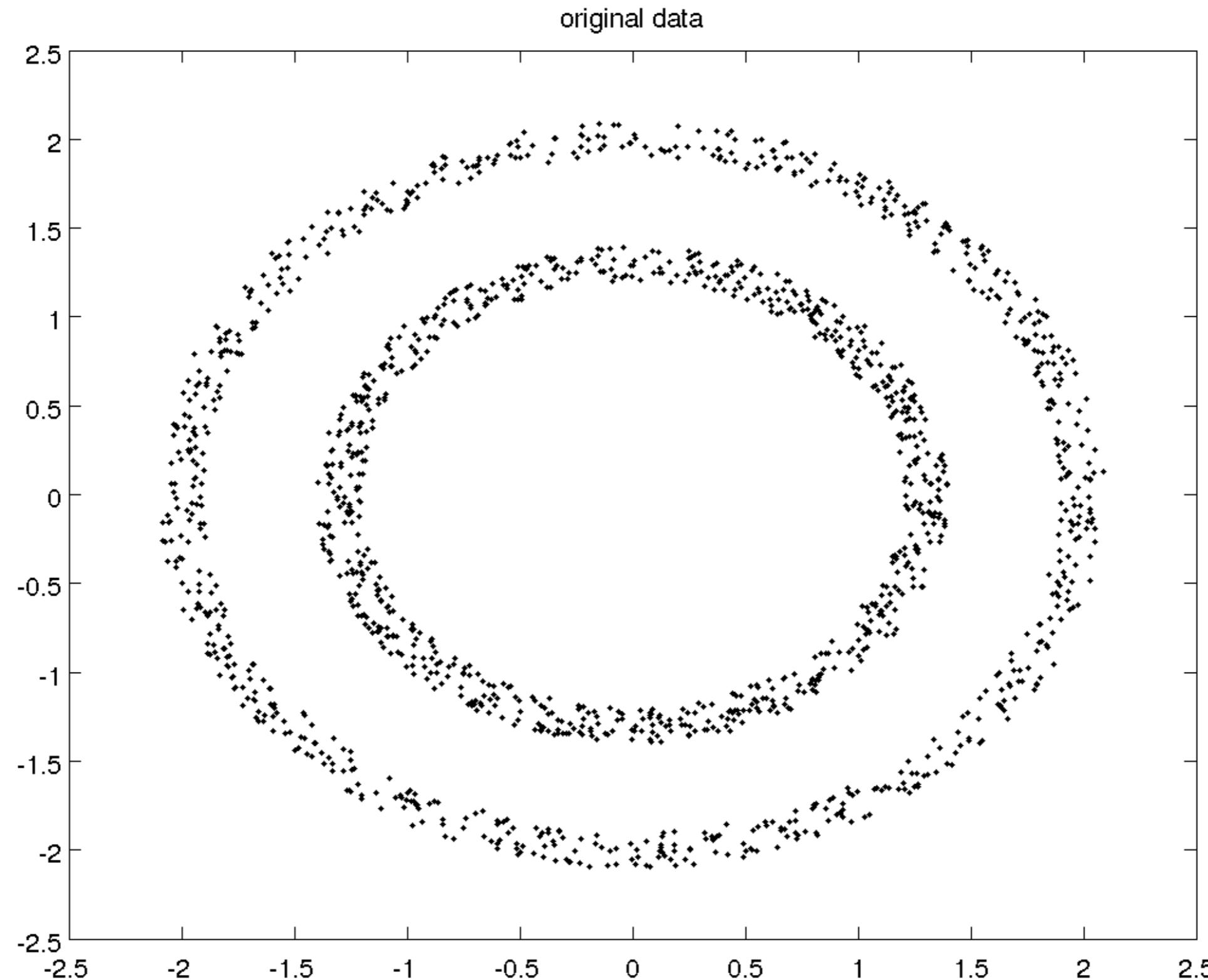
$$\vdots$$

$$Z_4 = (V_{14}, V_{24})$$

$$V_1 = \begin{pmatrix} 1.5 & 1 \\ 3 & 1.2 \\ 2.1 & 1.3 \\ 0.02 & 1.4 \end{pmatrix}$$

$$V_2 = \begin{pmatrix} 0.05 & 21 \\ 0.2 & 22 \\ 5 & 23 \\ 6.2 & 24 \end{pmatrix}$$

# How about this dataset?



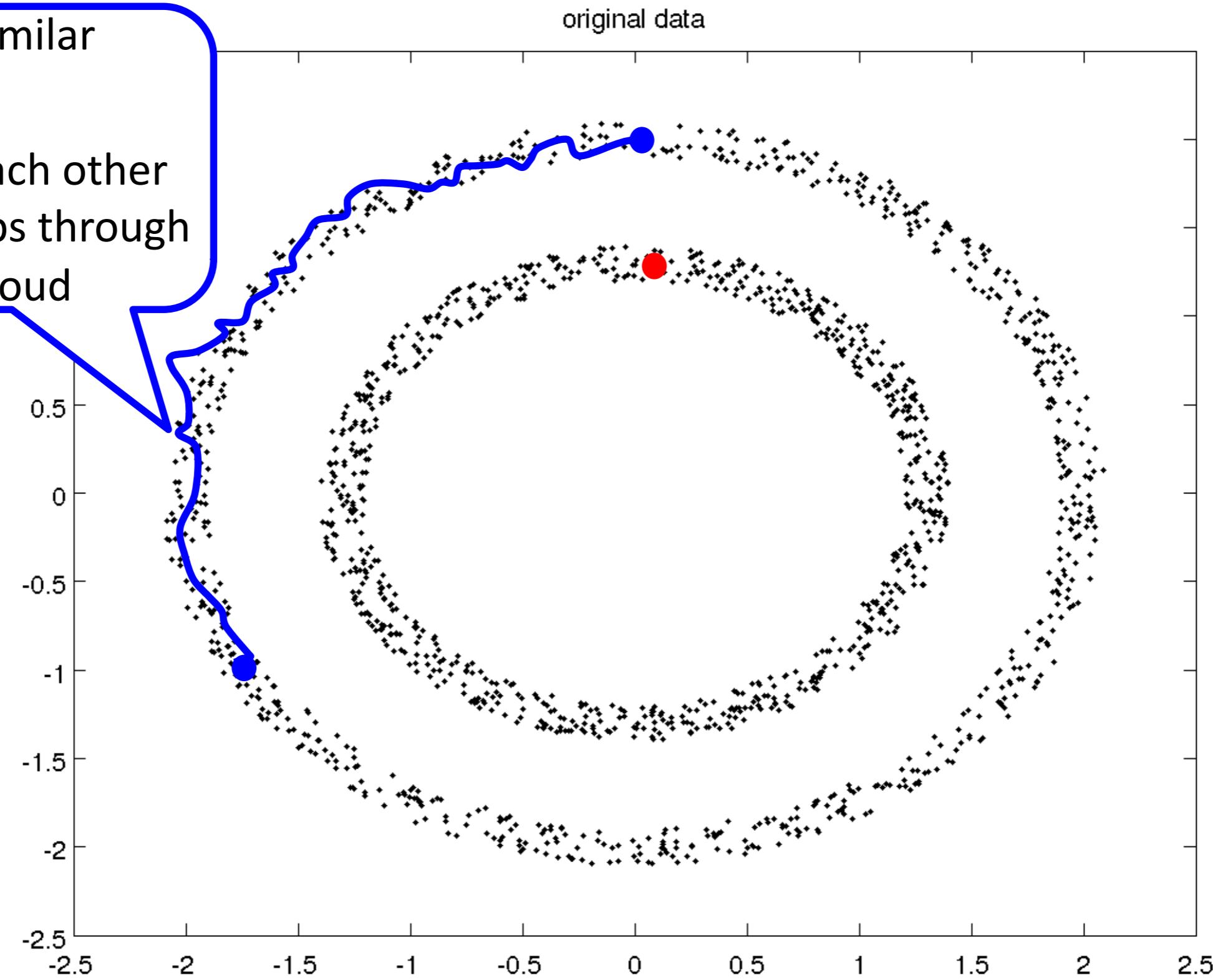
# What's a reasonable similarity measure?

points similar

=

can reach each other  
by small jumps through  
data cloud

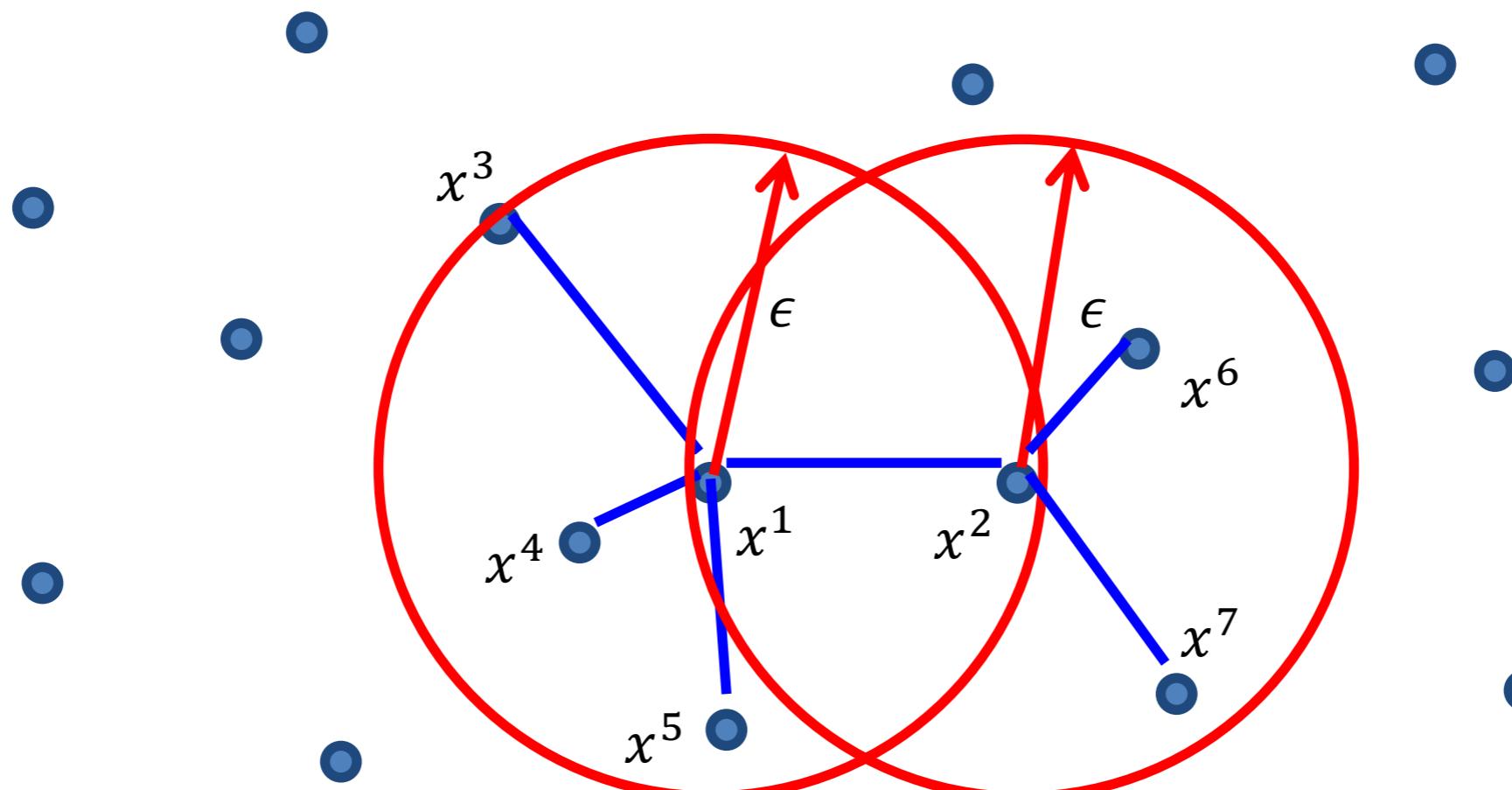
original data



# Neighbor graph

- Given  $m$  data points, threshold  $\epsilon$ , construct matrix  $A \in R^{m \times m}$

$$A^{ij} = \begin{cases} 1, & \text{if } \|x^i - x^j\| \leq \epsilon \\ 0, & \text{otherwise} \end{cases}$$



# Spectral clustering for vectorial data

- Given  $m$  nodes,  $\{x^1, x^2, \dots, x^m\} \in R^n$
- Step 1: build an adjacency matrix  $A$  using nearest neighbors
- Step 2: represent graph as adjacency matrix  $A \in R^{m \times m}$
- Step 3: form a special matrix  $L = D - A$ , the graph Laplacian
- Step 4: compute  $k$  eigenvectors,  $v^1, v^2, \dots, v^k$ , of  $L$  corresponding to the  $k$  **smallest** eigenvalues ( $k \ll m$ )
- Step 5: run kmeans algorithm on  $Z = (v^1, v^2, \dots, v^k)$  by treating each row as a new data point

# Variants of spectral clustering (Ng et al.)

- Given  $m$  data points (nodes),  $\{x^1, x^2, \dots, x^m\} \in R^n$

- Build an adjacency matrix  $A$  using **kernel functions**

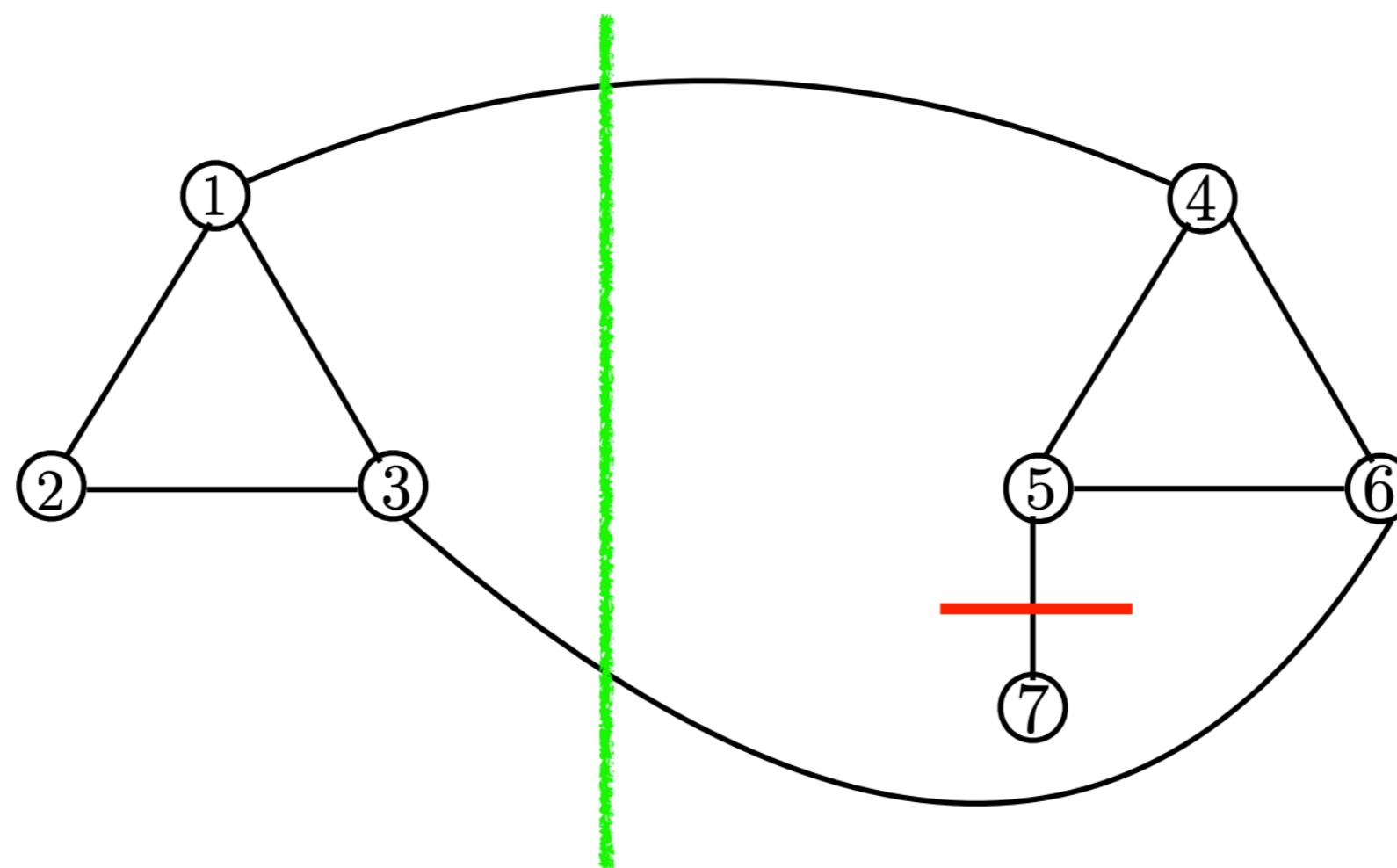
- Compute  $B = D^{-1/2}AD^{-1/2}$  (or  $B = I - D^{-1/2}AD^{-1/2}$ ),  
where  $D = \text{diag}(A)$

$$\begin{aligned} L &= D - A \\ D^{-1/2} L D^{-1/2} &= B \end{aligned}$$

- Compute  $k$  eigenvectors,  $v^1, v^2, \dots, v^k$ , of  $B$  corresponding to the  $k$  largest (or smallest) eigenvalues
- Run kmeans algorithm on  $Z = (v^1, v^2, \dots, v^k)$  by treating each row as a new data point

# Unnormalized vs normalized

- Unnormalized Spectral clustering aims to cluster based on minimizing cut
- cut: Number of edges that need to be deleted to have no links between the cluster and other nodes outside
- But is cut the right metric?



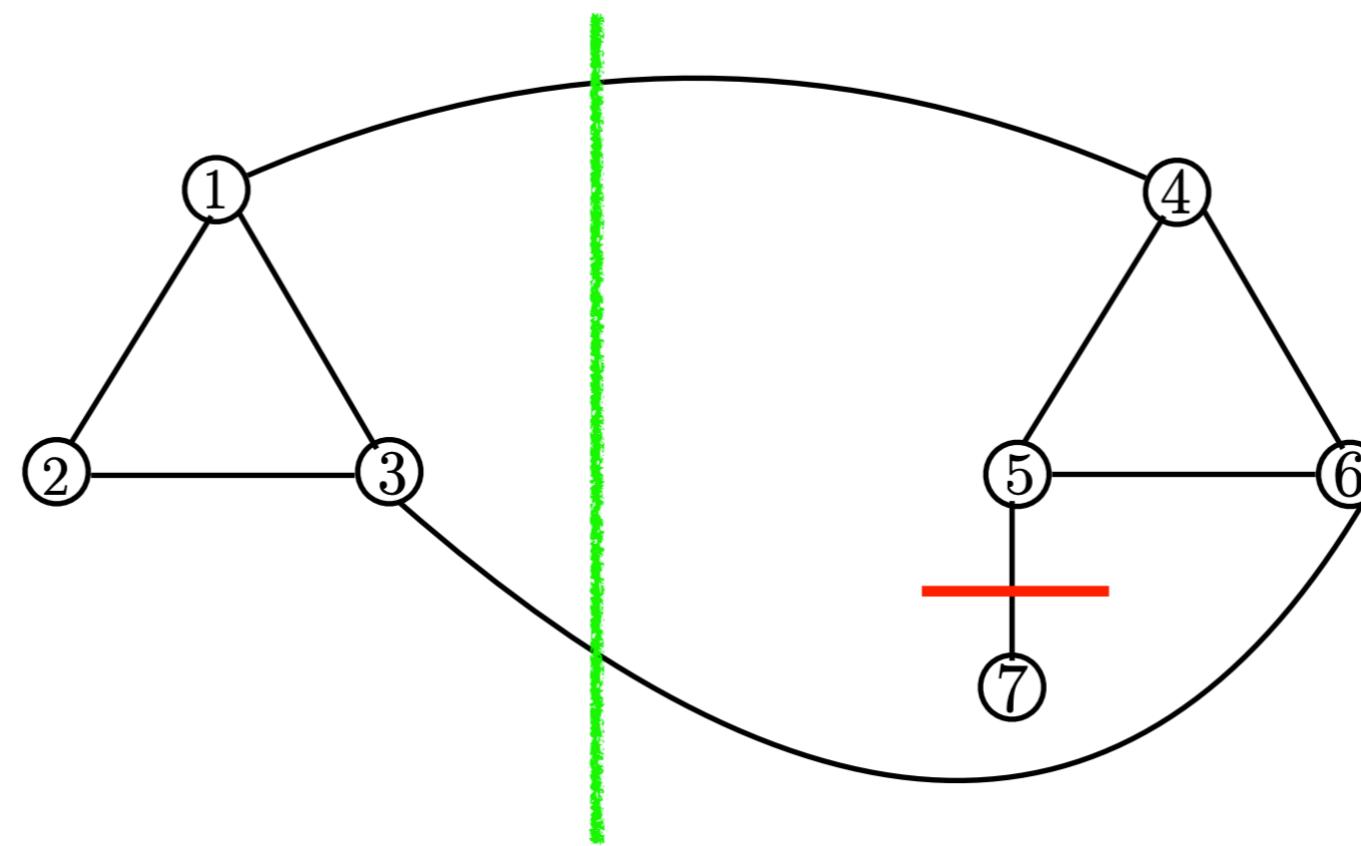
# Unnormalized vs normalized

- Normalized cut: Minimize sum of ratio of number of edges cut per cluster and number of edges within cluster

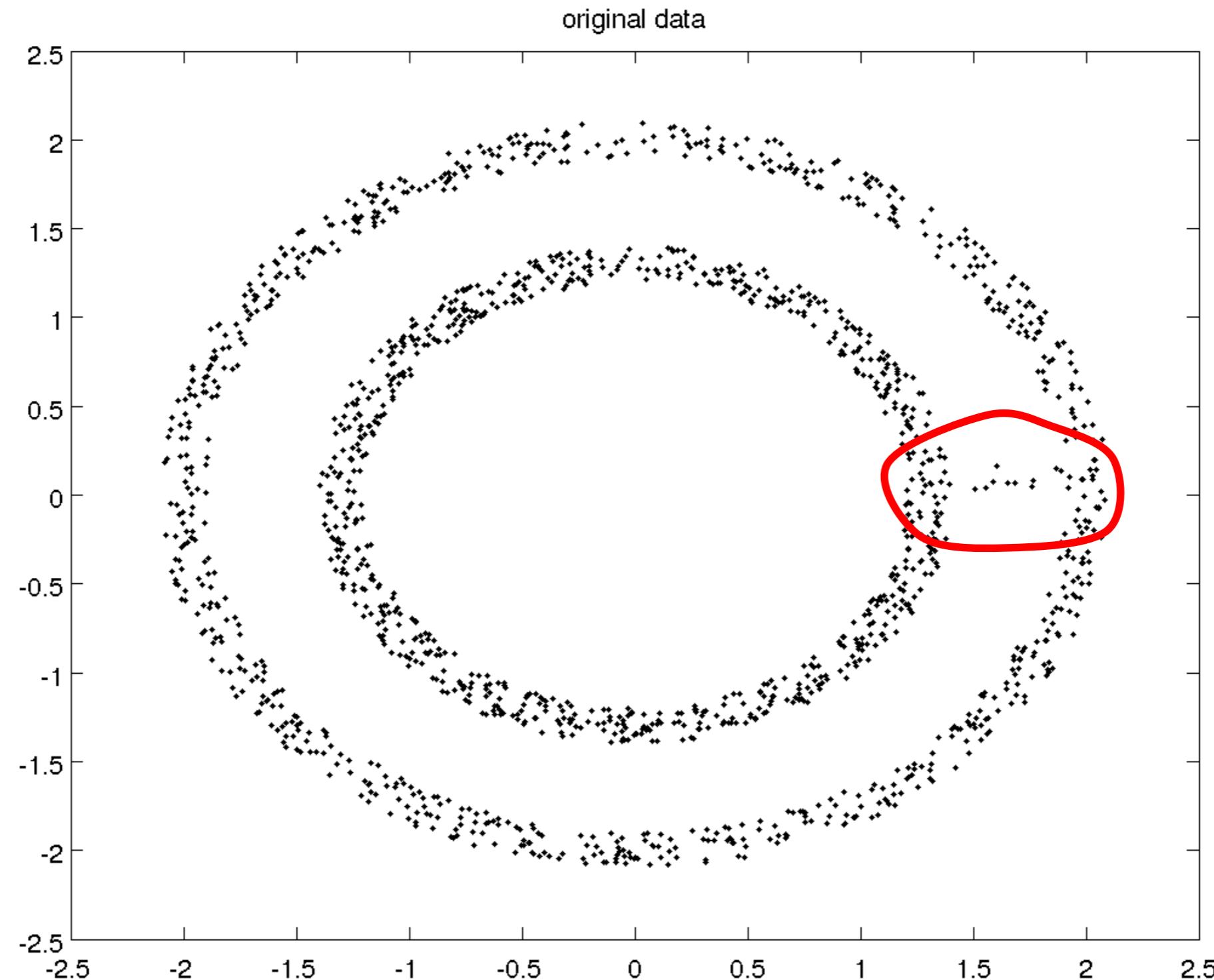
$$\text{NCUT} = \sum_j \frac{\text{CUT}(C_j)}{\text{Edges}(C_j)}$$

- Example K = 2

$$\text{CUT}(C_1, C_2) \left( \frac{1}{\text{Edges}(C_1)} + \frac{1}{\text{Edges}(C_2)} \right)$$



# What happens by adding more data points?



# What if my graph is large?

- Key challenge:
  - Eigen-decomposition of a large graph Laplacian is expensive
  - How to scale the algorithm up to millions of nodes?

- One solution:
  - Use randomized linear algebra to approximately find **top** eigenvectors of  $B = D^{-1/2}AD^{-1/2} \in R^{m \times m}$  (**big**)
    - Generate a Gaussian random matrix  $\Omega \in R^{d \times m}$  ( $d \ll m$ )
    - Find a set of orthonormal basis  $Q$  for column of

$$Y = B\Omega^\top$$

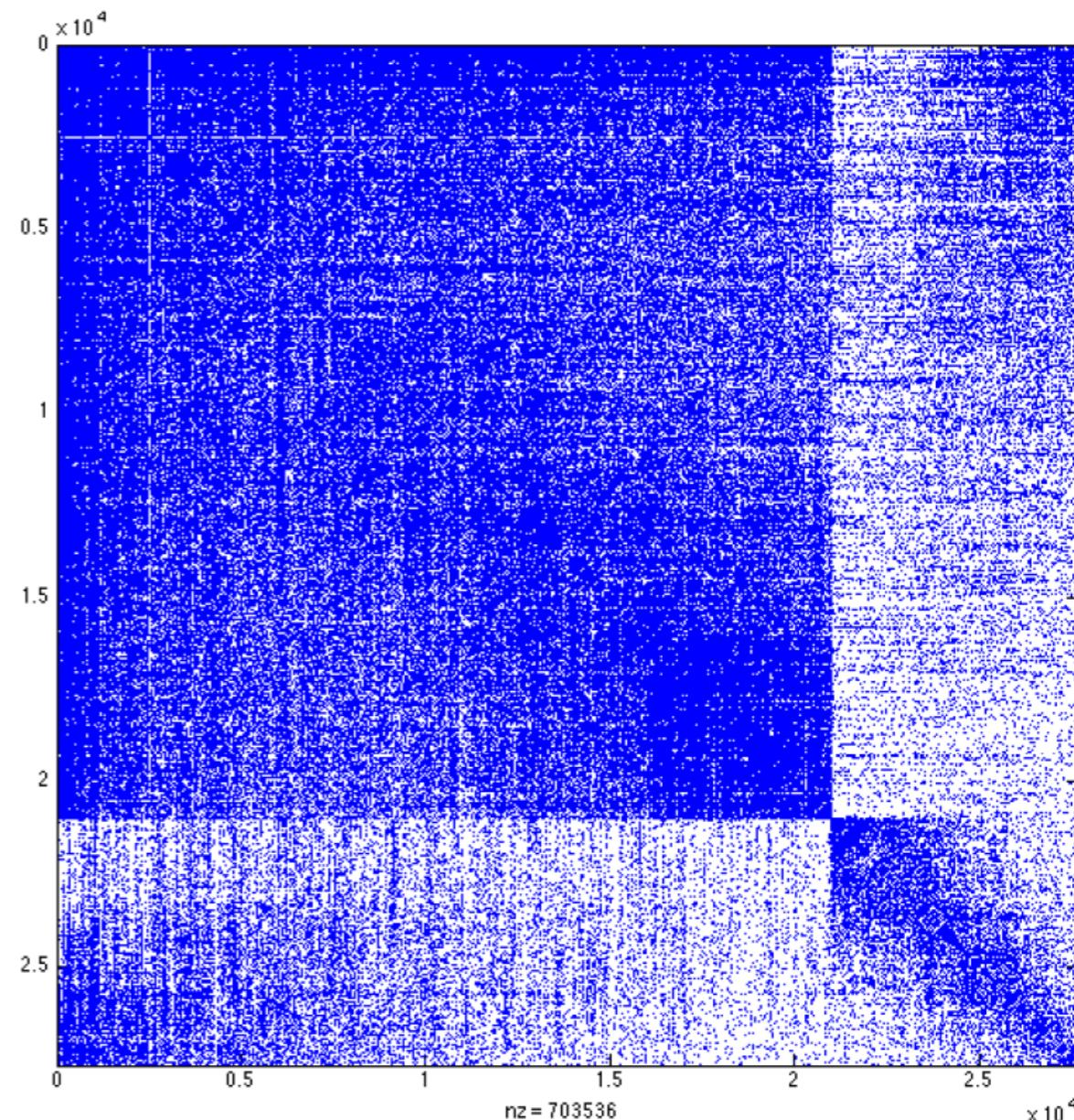
eg. using Gram-Schmidt orthogonalization

- Eigendecomposition of  $C = \Omega B \Omega^\top \in R^{d \times d}$  (**tiny**)
 
$$C = U \Lambda U^\top$$

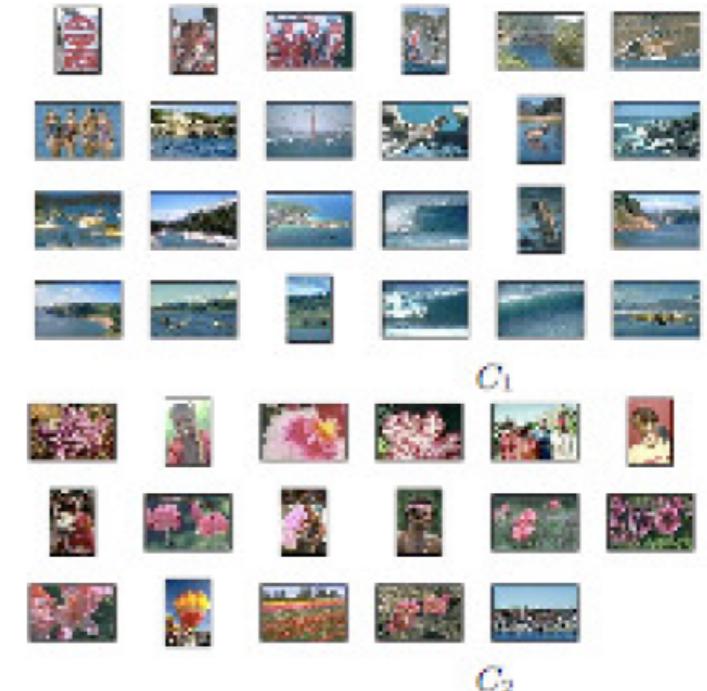
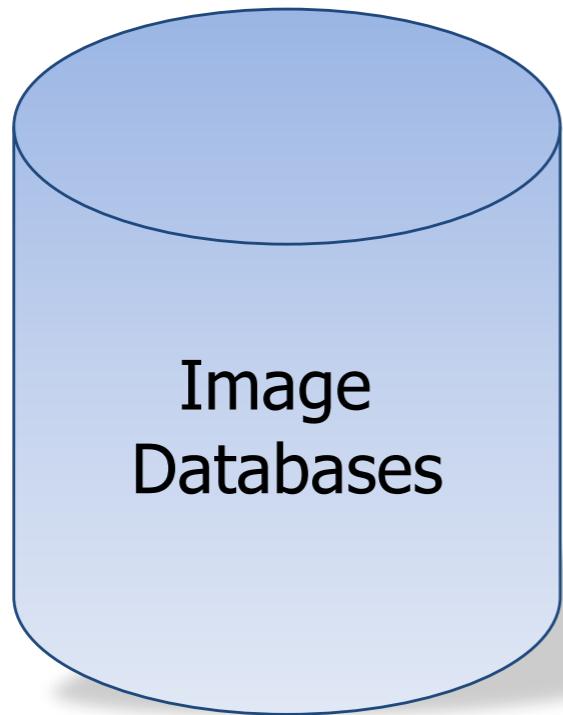
- $Z = QU$
- References: Halko, Martinsson and Tropp (2009).

# Run test\_citationgraph.m

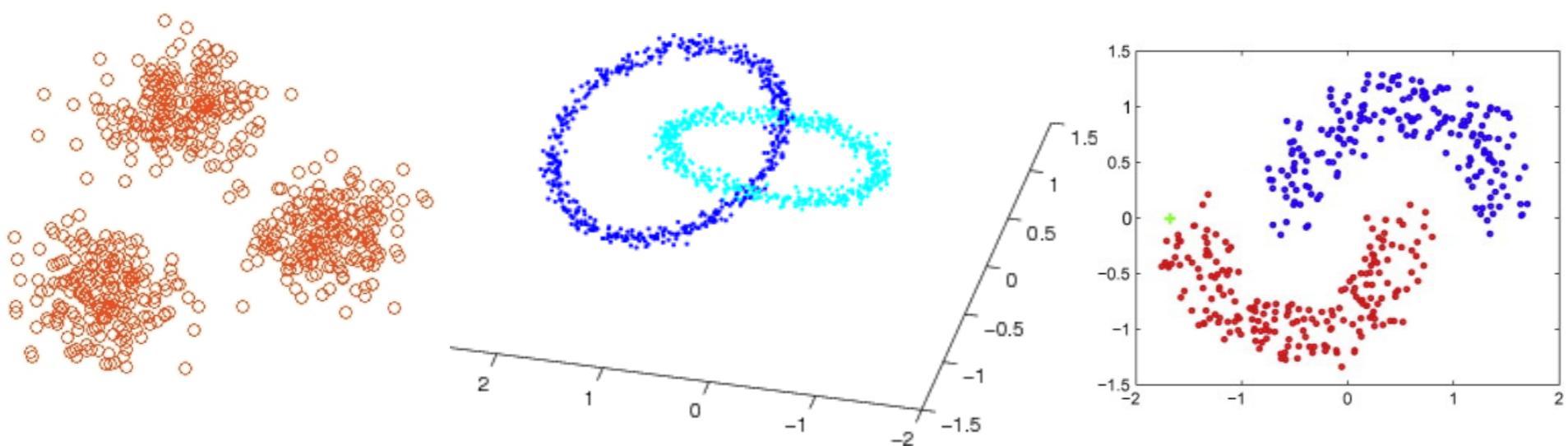
- Compare matlab eigs to the randomized eigendecomposition
- Compare matlab implementation of kmeans with our vectorized implementation
- Our codes run much faster.



# Image databases



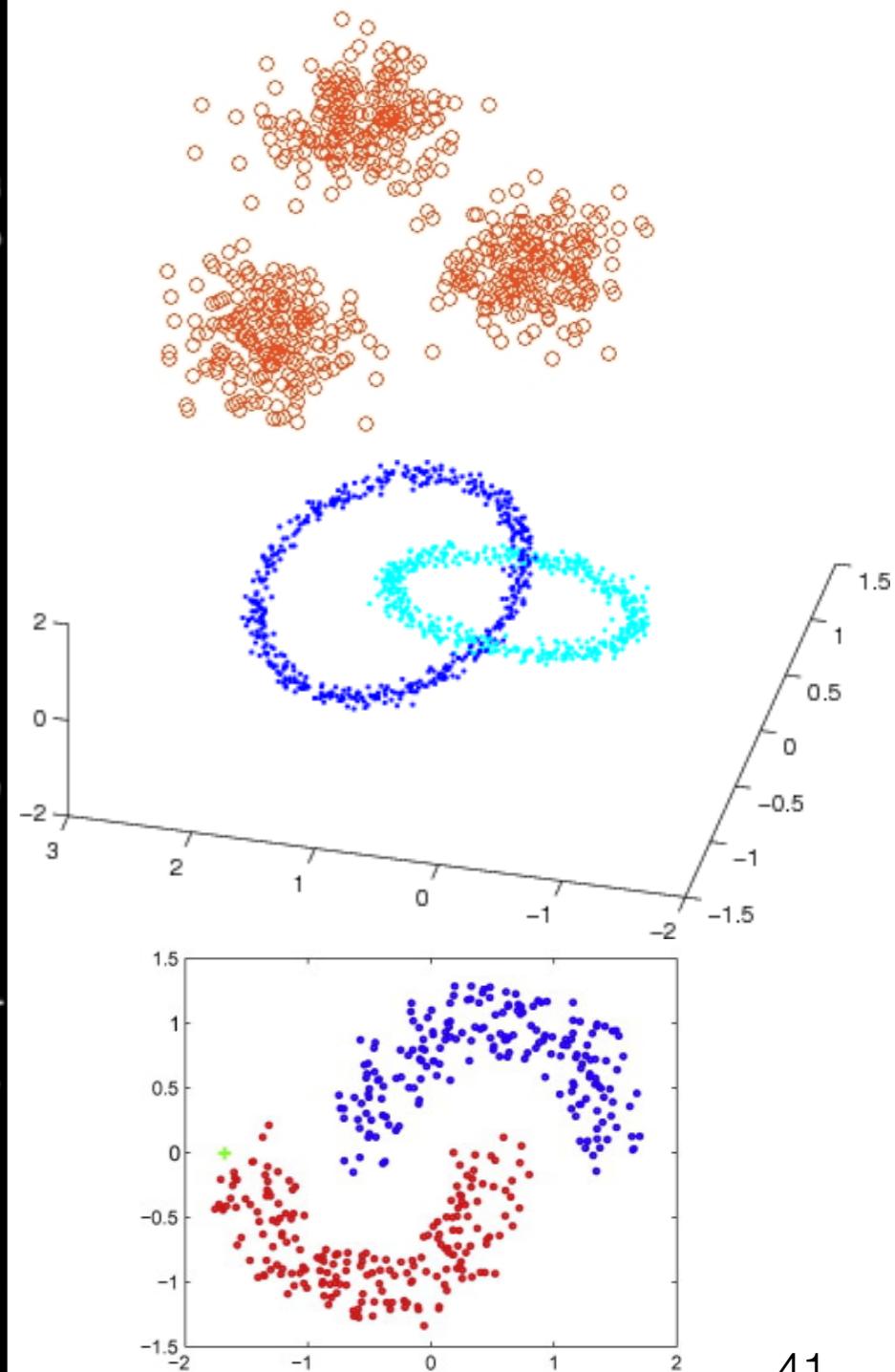
What are the relations  
between data points?



# Handwritten digits

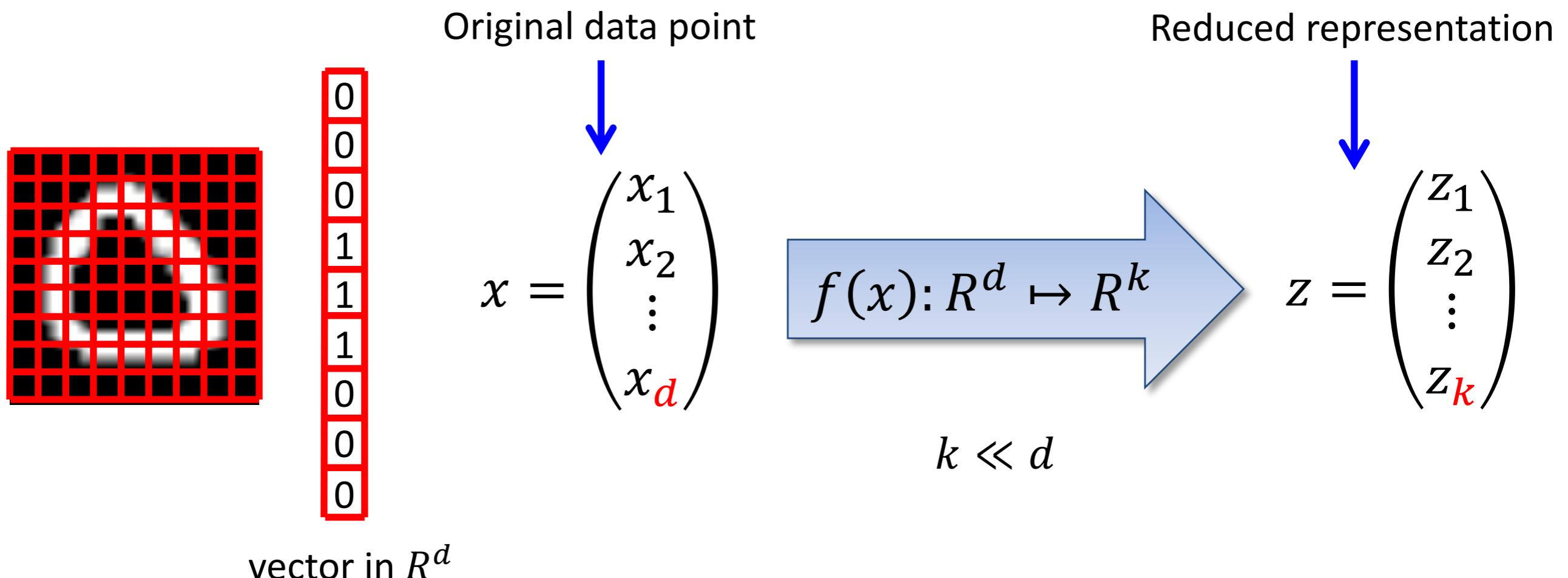
72104149590690159784966540740131
34727121174235124463556041957853
74643070291732977627847361369314
17696054992194873974449254767905
85665781016467317182029955156034
46546545144723271818185084250111
090316423611139529459390365573227
12841733887922415987230442419577
28268577918180301994182129759264
15829204002847124027433003196535
12930420711215339786561381051315
56185179462250656372088541140337
61621928619525442838245031773797
19219292049148184599837600302664
93332391268056663882758961841259
19754089914523789406395213136571
22632654897130383193446421825488
40023277087447969098046063548339
33378037170654380963809968685786
02402231975108462479309822927359
18020511376712580371409186774349
19317397691378336728585114431077
07944855408215845040613326726931
46259206217341054311749948402451
16471942415538314568941538032512
83440883317358632613607217142821
79611248177480231310770355276692
83522560829288887493066321322930
0578144602914747398847121223232323
91740355865267663279112564951334
78911691445406223151203812671623
90122089

What are the relations between data points?



# What is dimensionality reduction?

- The process of reducing the number of random variables under consideration
  - One can combine, transform or select variables
  - One can use linear or nonlinear operations



# Why dimensionality reduction and how to think

- The dimension-reduced data can be used for
  - Visualizing, exploring and understanding the data
  - Aggregating weak signals in the data
  - Cleaning the data
  - Speeding up subsequent learning task
  - Building simpler model later
- Key questions of a dimensionality reduction algorithm
  - What is the criterion for carrying out the reduction process?
  - What are the algorithm steps?

# Use what criterion for reduction?

- There are many criteria (geometric based, information theory based, etc.)
- One criterion: want to capture **variation** in data
  - variations are “signals” or information in the data
  - need to normalize each variables first
- In the process, also discover variables or dimensions highly **correlated**
  - represent highly related phenomena
  - combine them to form a stronger signal
  - lead to simpler presentation

# How to formulate the problem

- Given  $m$  data points,  $\{x^1, x^2, \dots, x^m\} \in R^d$ , with their mean  $\mu = \frac{1}{m} \sum_{i=1}^m x^i$
- Find a direction  $w \in R^d$  where  $\|w\| = 1$
- Such that the variance (or variation) of the data along direction  $w$  is maximized

$$\max_{w: \|w\|=1} \frac{1}{m} \sum_{i=1}^m (w^\top x^i - w^\top \mu)^2$$


variance

# Is it an easy optimization problem?

- Manipulate the objective with linear algebra

$$\begin{aligned} & \frac{1}{m} \sum_{i=1}^m (w^\top x^i - w^\top \mu)^2 \\ &= \frac{1}{m} \sum_{i=1}^m (w^\top (x^i - \mu))^2 \\ &= \frac{1}{m} \sum_{i=1}^m w^\top (x^i - \mu)(x^i - \mu)^\top w \\ &= w^\top \left( \frac{1}{m} \sum_{i=1}^m (x^i - \mu)(x^i - \mu)^\top \right) w \end{aligned}$$

  
covariance matrix  $C$

# Landscape of the optimization problem

- Suppose the data has two dimension ( $d = 2$ )
- $C$  is a diagonal matrix

$$C = \begin{pmatrix} 1 & 0 \\ 0 & 2 \end{pmatrix}$$

- The optimization problem becomes

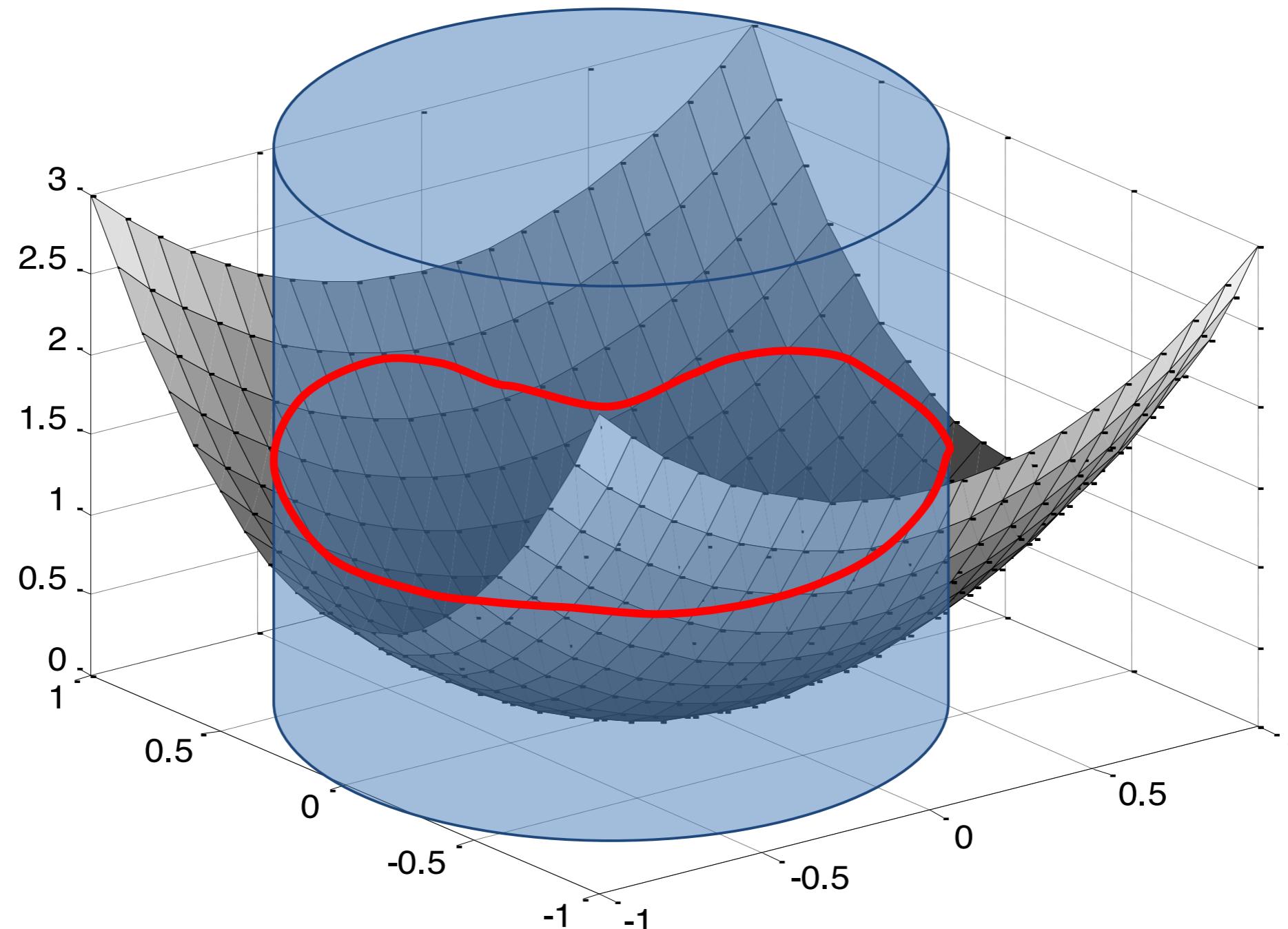
$$\max_{w: \|w\|=1} w^\top C w$$

$$= \max_{w: \|w\|=1} (w_1, w_2) \begin{pmatrix} 1 & 0 \\ 0 & 2 \end{pmatrix} \begin{pmatrix} w_1 \\ w_2 \end{pmatrix}$$

$$= \max_{w: \|w\|=1} w_1^2 + 2w_2^2$$

# Landscape of the optimization problem

- $f(w_1, w_2) = w_1^2 + 2w_2^2$



# Eigen-value problem

- Eigen-value problem
  - Given a symmetric matrix  $C \in R^{d \times d}$
  - Find a vector  $w \in R^d$  and  $\|w\| = 1$
  - Such that
$$Cw = \lambda w$$
- There will be multiple solution of  $w^1, w^2, \dots$  with different  $\lambda_1, \lambda_2, \dots$ 
  - They are ortho-normal:  $w^{i^\top} w^i = 1, w^{i^\top} w^j = 0$

# Equivalent to eigen-value problem

- Claim:

$$\max_{w: \|w\|=1} w^T C w \Leftrightarrow Cw = \lambda w$$

- Form lagrangian function of the optimization problem

$$L(w, \lambda) = w^T C w + \lambda(1 - \|w\|^2)$$

Necessary  
condition

- If  $w$  is a maximum of the original optimization problem, then there exists a  $\lambda$ , where  $(w, \lambda)$  is a **stationary point** of  $L(w, \lambda)$
- This implies that

$$\frac{\partial L}{\partial w} = 0 = 2Cw - 2\lambda w$$

# Variance of in the principal direction

- Principal direction  $w$  satisfies

$$Cw = \lambda w$$

- Variance in principal direction is

# Variance of in the principal direction

- Principal direction  $w$  satisfies

$$Cw = \lambda w$$

- Variance in principal direction is

$$w^T C w$$

$$= \lambda w^T w$$

$$= \lambda$$

eigen-value

# Multiple principal directions

- Directions  $w^1, w^2, \dots$  which has
  - the largest variances
  - but are **orthogonal** to each other
- Take the eigenvectors  $w^1, w^2, \dots$  of  $C$  corresponding to
  - the largest eigenvalue  $\lambda_1$ ,
  - the second largest eigenvalue  $\lambda_2$
  - ...

# Principal component analysis

- Given  $m$  data points,  $\{x^1, x^2, \dots, x^m\} \in R^d$ , with mean
- Step 1: Estimate the mean and covariance matrix from data

$$\mu = \frac{1}{m} \sum_{i=1}^m x^i \quad \text{and} \quad C = \frac{1}{m} \sum_{i=1}^m (x^i - \mu)(x^i - \mu)^T$$

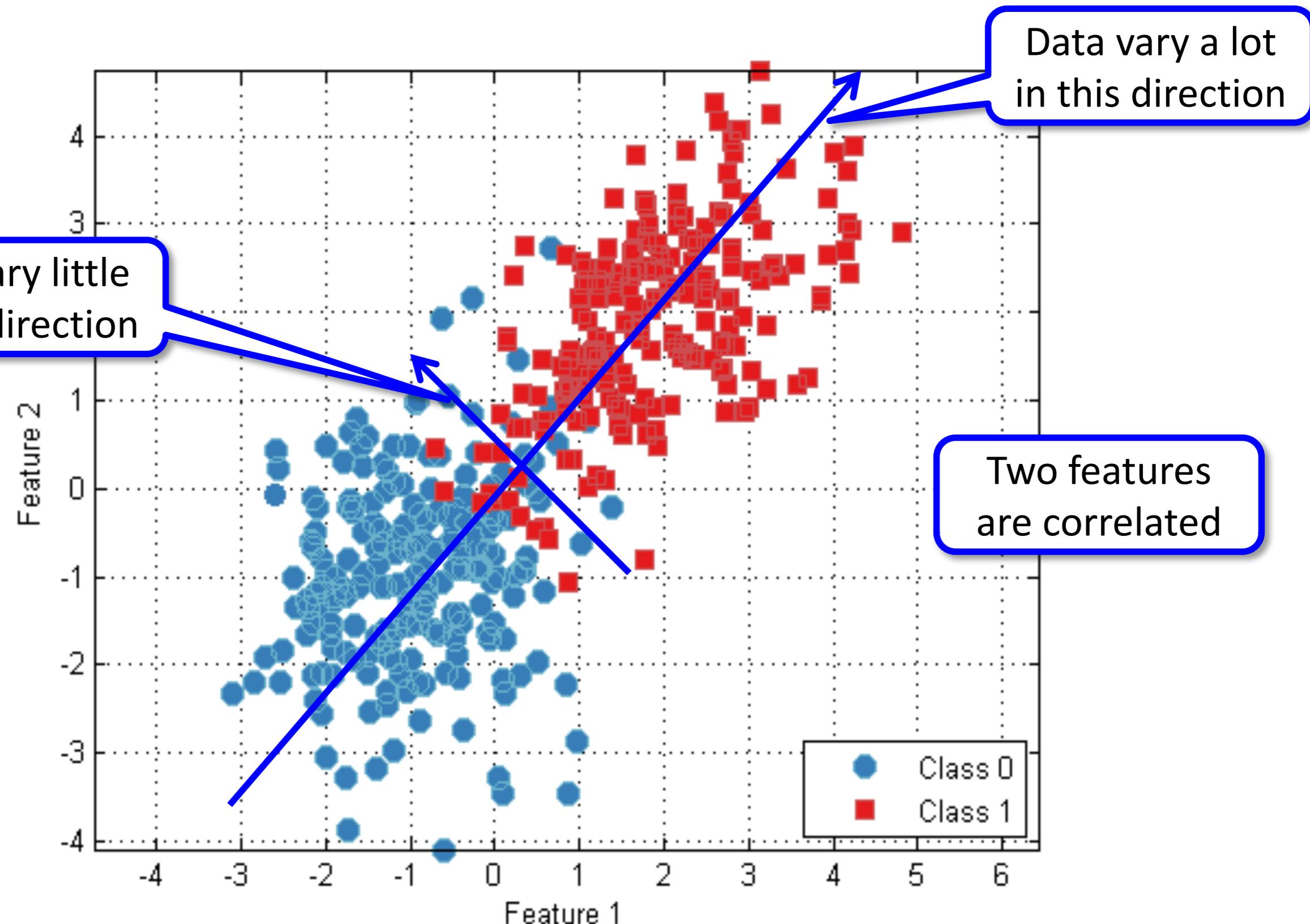
Principal directions

- Step 2: Take the eigenvectors  $w^1, w^2, \dots$  of  $C$  corresponding to the largest eigenvalue  $\lambda_1$ , the second largest eigenvalue  $\lambda_2 \dots$
- Step 3: Compute reduced representation

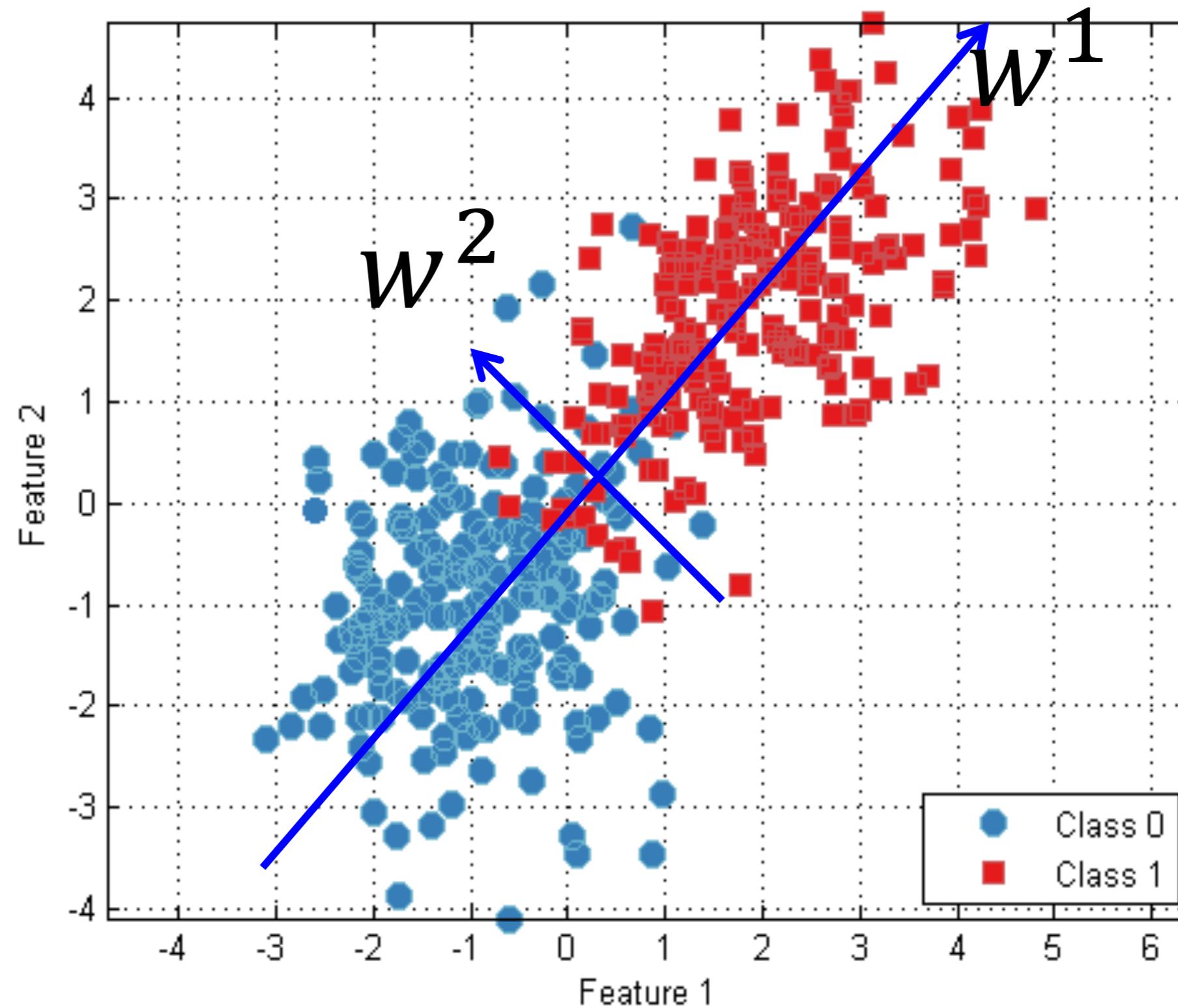
$$z^i = \begin{pmatrix} w^1^T(x^i - \mu) / \sqrt{\lambda_1} \\ w^2^T(x^i - \mu) / \sqrt{\lambda_2} \\ \vdots \end{pmatrix}$$

Normalize by  
standard deviation

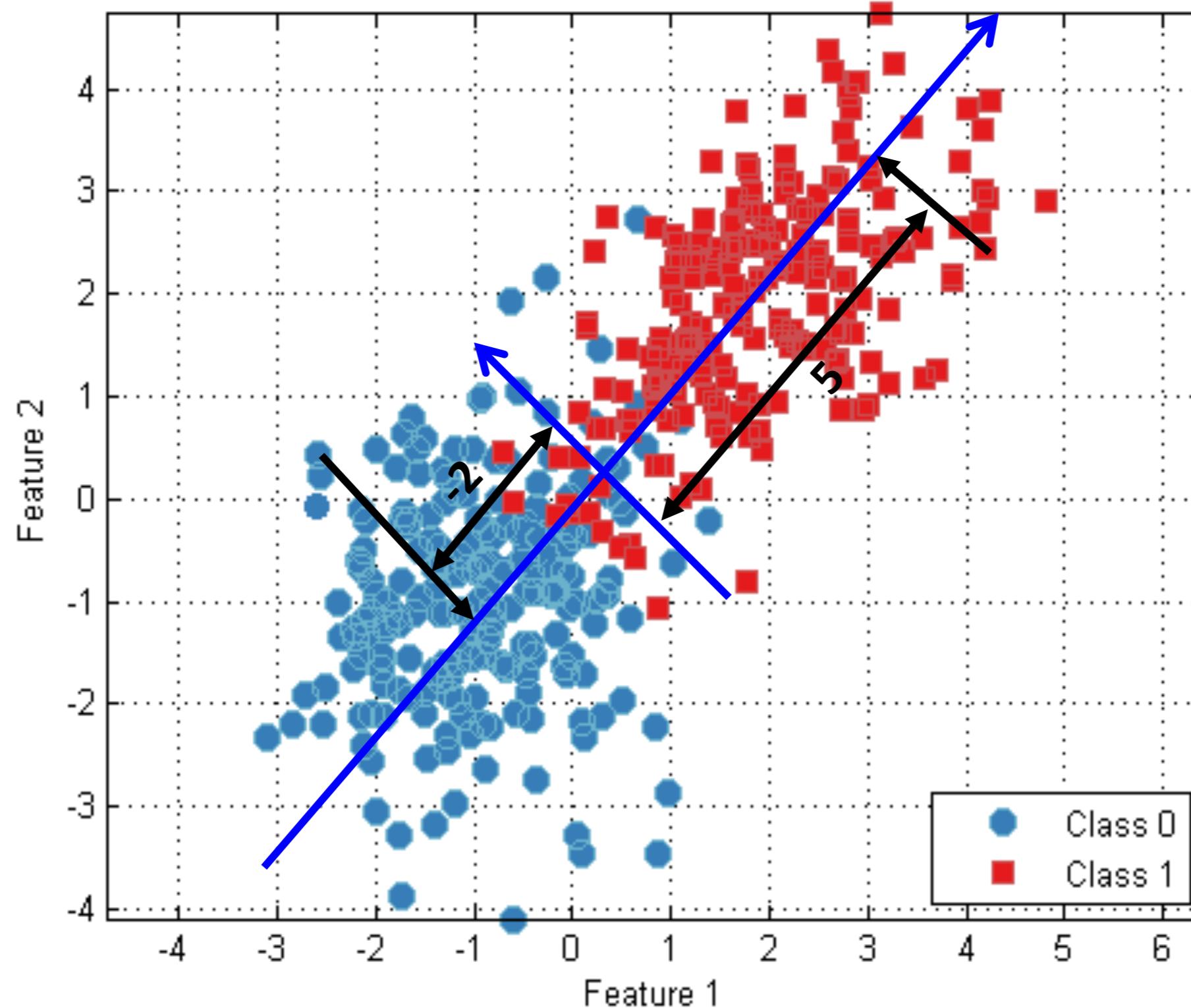
# An example



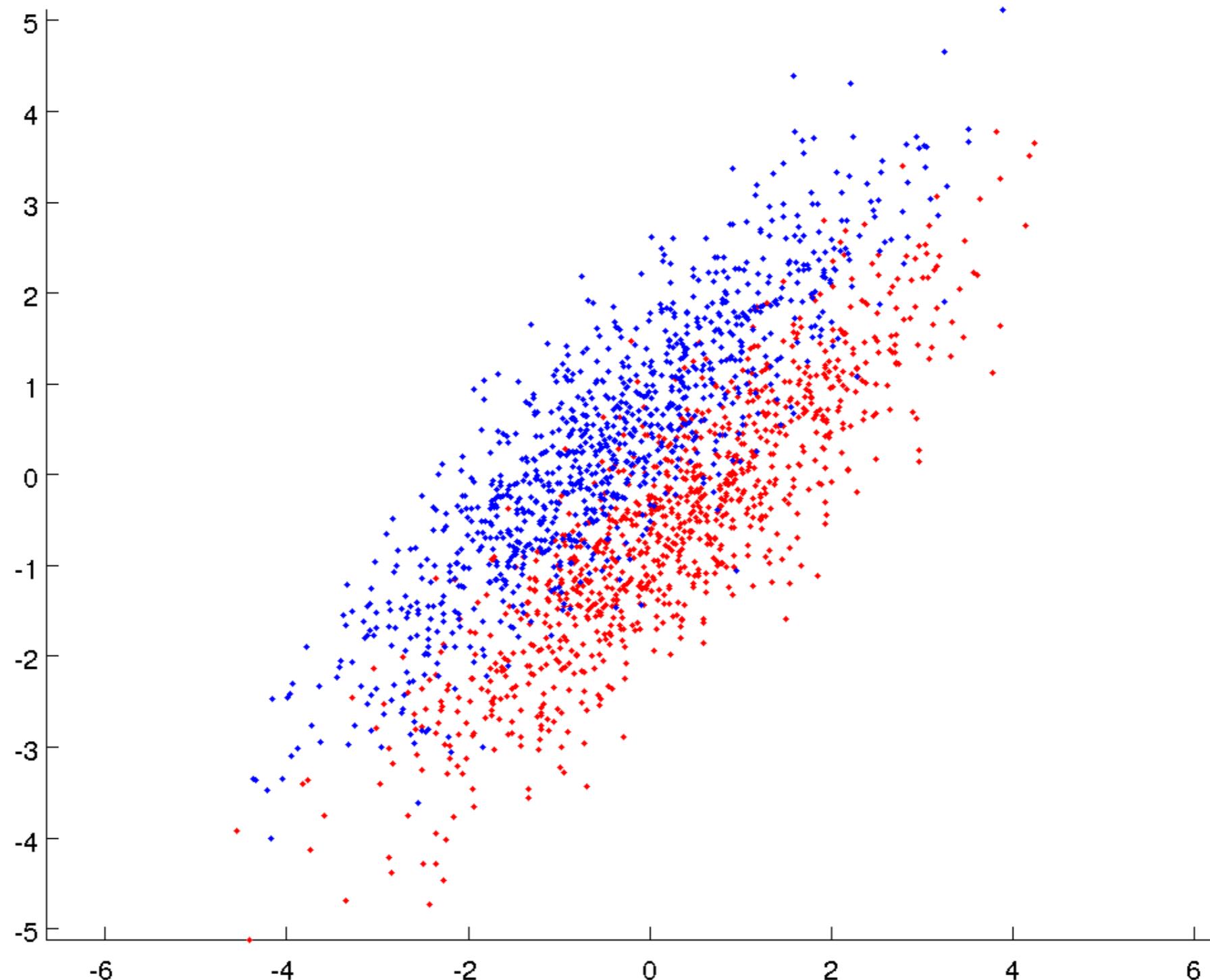
# Principal direction of the data



# Reduce to 1 dimension

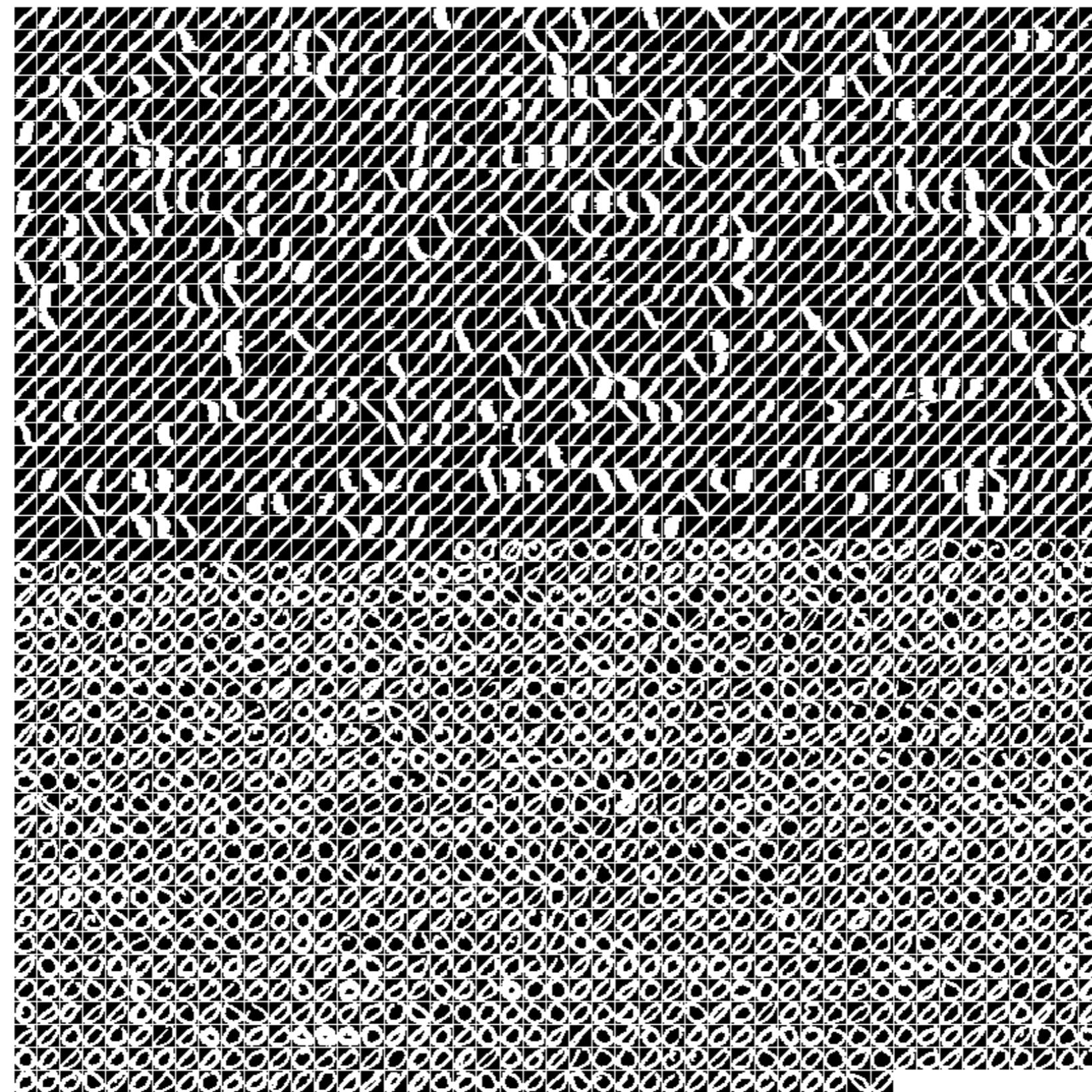


# Are principal components good for classification?



# Run demo PCA\_digits.m

digit 1 and 0



# Run demo PCA\_leaf.m



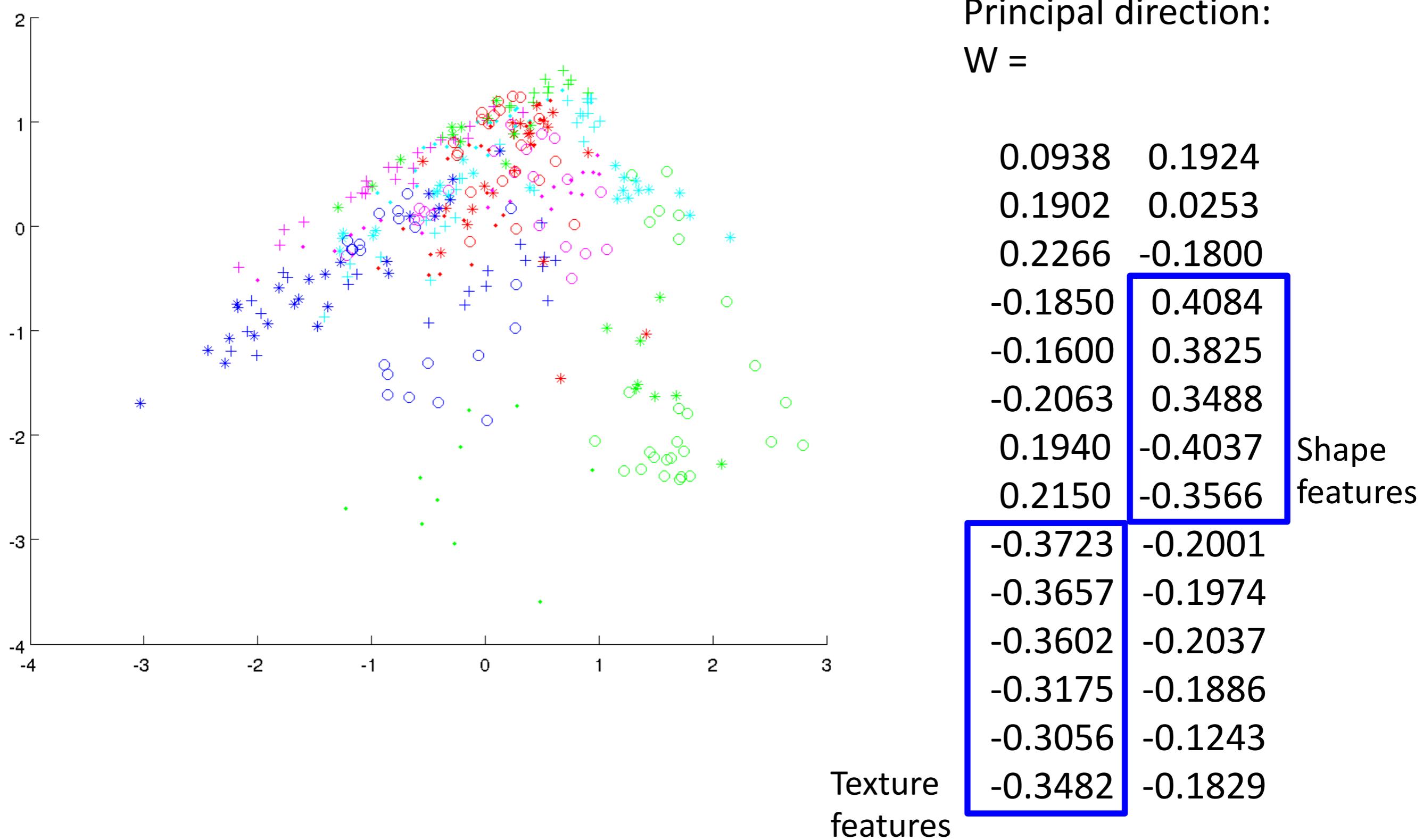
# Input features (representation)

Shape feature	Description
<i>Eccentricity</i>	Eccentricity of the ellipse with identical second moments to $I$ . This value ranges from 0 to 1.
<i>Aspect Ratio</i>	Consider any $X, Y \in \partial I$ . Choose $X$ and $Y$ such that $d(X, Y) = D(I)$ . Find $Z, W \in \partial I$ maximizing $D^\perp = d(Z, W)$ on the set of all pairs of $\partial I$ that define a segment orthogonal to $[XY]$ . The aspect ratio is defined as the quotient $D(I)/D^\perp$ . Values close to 0 indicate an elongated shape.
<i>Elongation</i>	Compute the maximum escape distance $d_{\max} = \max_{X \in I} d(X, \partial I)$ . Elongation is obtained as $1 - 2d_{\max}/D(I)$ and ranges from 0 to 1. The minimum is achieved for a circular region. Note that the ratio $2d_{\max}/D(I)$ is the quotient between the diameter of the largest inscribed circle and the diameter of the smallest circumscribed circle.
<i>Solidity</i>	The ratio $A(I)/A(H(I))$ is computed, which can be understood as a certain measure of convexity. It measures how well $I$ fits a convex shape.
<i>Stochastic Convexity</i>	This variable extends the usual notion of convexity in topological sense, using sampling to perform the calculation. The aim is to estimate the probability of a random segment $[XY]$ , $X, Y \in I$ , to be fully contained in $I$ .
<i>Isoperimetric Factor</i>	The ratio $4\pi A(I)/L(\partial I)^2$ is calculated. The maximum value of 1 is reached for a circular region. Curvy intertwined contours yield low values.
<i>Maximal Indentation Depth</i>	Let $C_{H(I)}$ and $L(H(I))$ denote the centroid and arclength of $H(I)$ . The distances $d(X, C_{H(I)})$ and $d(Y, C_{H(I)})$ are computed $\forall X \in H(I)$ and $\forall Y \in \partial I$ . The indentation function can then be defined as $[d(X, C_{H(I)}) - d(Y, C_{H(I)})]/L(H(I))$ , which is sampled at one degree intervals. The maximal indentation depth $\mathfrak{D}$ is the maximum of this function.
<i>Lobedness</i>	The Fourier Transform of the indentation function above is computed after mean removal. The resulting spectrum is normalized by the total energy. Calculate lobedness as $F \times \mathfrak{D}^2$ , where $F$ stands for the smallest frequency at which the cumulated energy exceeds 80%. This feature characterizes how lobed a leaf is.

Texture feature	Description
<i>Average Intensity</i>	Average intensity is defined as the mean of the intensity in $I$ .
<i>Average Contrast</i>	Average contrast is the standard deviation of the intensity in $I$ , $\sigma = \sqrt{\mu_2(z)}$ .
<i>Smoothness</i>	Smoothness is defined as $R = 1 - 1/(1 + \sigma^2)$ and measures the relative smoothness of the intensities in a given region. For regions of constant intensity, $R$ takes the value 0 and $R$ approaches 1 for regions exhibiting larger disparities in intensity values. $\sigma^2$ is normalized by $(L - 1)^2$ to ensure that $R \in [0, 1]$ .
<i>Third moment</i>	$\mu_3$ is a measure of the intensity histogram's skewness. This is generally normalized by $(L - 1)^2$ like smoothness.
<i>Uniformity</i>	Defined as $U = \sum_{i=0}^{L-1} p^2(z_i)$ , uniformity's maximum is reached when all intensity levels are equal.
<i>Entropy</i>	A measure of intensity randomness.

8 shape features and  
6 texture features

# Reduce representation



# Singular Value Decomposition

- Singular value decomposition, known as SVD, is a factorization of a real matrix with applications in calculating pseudo-inverse, rank, solving linear equations, and many others.
- For a matrix  $M \in \mathbb{R}^{m \times n}$  assume  $n \leq m$ 
  - $M = U\Sigma V^T$  where  $U \in \mathbb{R}^{m \times m}$ ,  $V^T \in \mathbb{R}^{n \times n}$ ,  $\Sigma \in \mathbb{R}^{m \times n}$
  - The  $m$  columns of  $U$ , and the  $n$  columns of  $V$  are called the left and right singular vectors of  $M$ . The diagonal elements of  $\Sigma$ ,  $\Sigma_{ii}$  are known as the singular values of  $M$ .
  - Let  $v$  be the  $i^{\text{th}}$  column of  $V$ , and  $u$  be the  $i^{\text{th}}$  column of  $U$ , and  $\sigma$  be the  $i^{\text{th}}$  diagonal element of  $\Sigma$

$$Mv = \sigma u \quad \text{and} \quad M^T u = \sigma v$$

# Singular Value Decomposition - II

- $$M = [u_1 \ u_2 \ \dots \ u_m] \begin{bmatrix} \Sigma_{11} & \cdots & \Sigma_{1n} \\ \vdots & \ddots & \vdots \\ \Sigma_{m1} & \cdots & \Sigma_{mn} \end{bmatrix} [v_1 \ v_2 \ \dots \ v_n]^T$$

principal directions

Scaling factor

Projection in principal directions
- Singular value decomposition is related to eigenvalue decomposition
  - Suppose  $M = [x_1 - \mu \ x_2 - \mu \ \dots \ x_m - \mu] \in \mathbb{R}^{m \times n}$
  - Then covariance matrix is  $C = \frac{1}{m} MM^T$
  - Starting from singular vector pair
    - $M^T u = \sigma v$
    - $\Rightarrow MM^T u = \sigma M v$
    - $\Rightarrow MM^T u = \sigma^2 u$
    - $\Rightarrow Cu = \lambda u$