# Maximum Variance Subspace

- Let $v_1, v_2, \ldots, v_p$ denote the $p$ principal components
$$v_i \cdot v_j = 0 \text{ for } i \neq j \text{ and } v_i \cdot v_j = 1 \text{ for } i = j$$
- Assume that data $X$ is centered (mean = 0)
  - if not, we can preprocess it: $\tilde{X} = X - \frac{1}{n}\mathbf{1}_n\mathbf{1}_n^T X$
- Find vector that maximizes sample variance of projected data

$$\frac{1}{n-1}\sum_{i=1}^{n}(v^T x_i)^2 = v^T X^T X v$$

$$\max_{v} v^T X^T X v \quad s.t. \quad v^T v = 1$$

- Lagrangian: $L(v, \lambda) = \frac{1}{n-1}v^T X^T X v - \lambda(v^T v - 1)$

$$\frac{\partial L}{\partial v} = 0 \qquad \boxed{\left(\frac{1}{n-1}X^T X - \lambda I\right) v = 0}$$

Eigen Problem

*(handwritten annotations in red):*

$$\checkmark$$

$$= \frac{1}{n-1} v^T X^T X v$$

$$C = \frac{1}{n-1} X^T X$$

# Eigen Decomposition

$$V = [v_1, \ldots, v_p] \in \mathbb{R}^{d \times p}, Y = XV$$

$$y_i = V^T x_i = \sum_{j=1}^{p} v_j^T x_i$$

- Eigen-decomposition

$$\max_{v_1, \cdots, v_p} trace(V^T C V) . s.t. V^T V = I \rightarrow max_V \ trace\left(V^T C V - \Lambda(V^T V - I)\right)$$

The sum of variance after projection

○  $\frac{1}{n-1}\sum_{i=1}^{n} y_i^T y_i = \frac{1}{n-1}\sum_{i=1}^{n}\sum_{j=1}^{p}(v_j^T x_i)^2 = \frac{1}{n-1} trace(V^T X^T X V), \ y_i = (v_1^T x_i, \ldots, v_p^T x_i)$

○  $cov(Y, Y) = \frac{1}{n-1}Y^T Y = \frac{1}{n-1}V^T X^T X V = V^T V \Lambda V^T V = \Lambda$    Statistically uncorrelated

# Eigen Decomposition

- In $Cv = \lambda v$, $v$ (the first PC) is the eigenvector of sample covariance matrix $C = \frac{1}{n-1}X^T X$
  - Sample variance of projection becomes $v^T C v = \lambda v^T v = \lambda$
- For the eigenvectors and values for the sample variance matrix $C$,
  - Eigenvalues $\lambda_1 \geq \lambda_2 \geq \lambda_3 \geq \cdots$
  - The 1$^{st}$ PC $v_1$ is the eigenvector of the sample covariance matrix $C$ associated with the largest eigenvalue
  - The 2$^{nd}$ PC $v_2$ is the eigenvector of the sample covariance matrix $C$ associated with the second largest eigenvalue
  - $\sum_j var(y_j) = trace\big(cov(Y,Y)\big) = \lambda_1 + \lambda_2 + \cdots$
    - $cov(Y,Y) = \frac{1}{n-1}Y^T Y = \frac{1}{n-1}V^T X^T X V = V^T V \Lambda V^T V = \Lambda$

$p$

$trace\big(cov(Y,Y)\big)$

$= \lambda_1 + \cdots + \lambda_p$

# Singular Value Decomposition

- For an n $\times$ d matrix $X$ of rank r there exists a factorization (Singular Value Decomposition = SVD) as follows:

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

  - The columns of $U$ are orthogonal eigenvectors of $XX^T$
  - The columns of $V$ are orthogonal eigenvectors of $X^T X$
  - Eigenvalues $\mu_1, ..., \mu_r$ of $XX^T$ are the eigenvalues of $X^T X$ $(rank(X) = r)$
  - Singular values $\sigma_i = \sqrt{\mu_i}$ where $\Sigma = diag(\sigma_1, ..., \sigma_r)$
  - For eigenvalue $\lambda_i$ of $C = \frac{1}{n-1} X^T X$, $\lambda_i = \frac{\sigma_i^2}{n-1}$

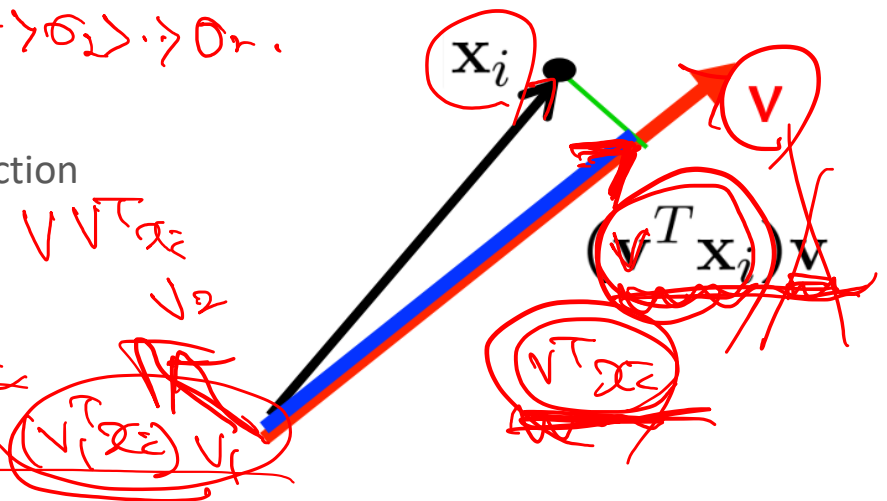- Minimum Reconstruction Error
  - Find the subspace that yields minimum MSE reconstruction
    - MSE: $\frac{1}{n}\sum_{i=1}^{n} \|x_i - VV^T x_i\|^2$ s.t. $V^T V = I$
  - Low rank approximation
    - $\min_{A:rank(A)=p} \|X - A\|_F = \sigma_{p+1}$

blue² + green² = black²

black² is fixed (it's just the data)

So, maximizing blue² is equivalent to minimizing green²

# Dimensionality Reduction using PCA

- Original representation
  - Data point: $x_i = (x_i^1, \dots, x_i^d)$
- Transformed representation
  - Principal components (PCs): a normalized linear combination of the original features $(v_1, v_2, \dots, v_p)$
  - PC loadings: how much each original variable contributes to that principal component
    - Loading for the i-th PC are given by the column of V ($v_i$)
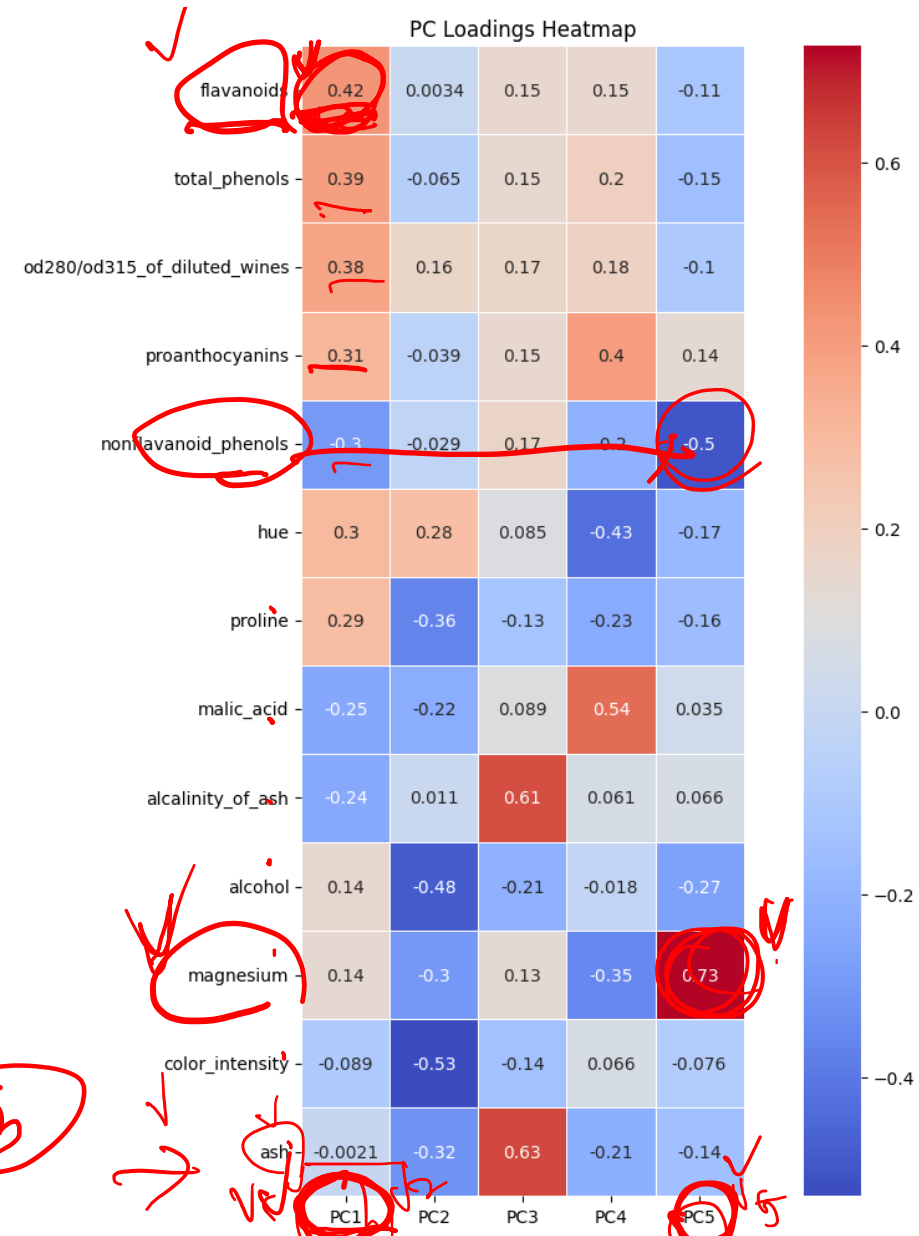    - The j-th variable on the i-th PC is $v_{ij}$
  - PC score (Projection): $V^T x_i = (v_1^T x_i, \dots, v_p^T x_i)$
- Ratio of Explained variance: $\sum_{i=1}^{p} \lambda_i / trace(C)$



PC Loadings Heatmap

| | PC1 | PC2 | PC3 | PC4 | PC5 |
|---|---|---|---|---|---|
| flavanoids | 0.42 | 0.0034 | 0.15 | 0.15 | -0.11 |
| total_phenols | 0.39 | -0.065 | 0.15 | 0.2 | -0.15 |
| od280/od315_of_diluted_wines | 0.38 | 0.16 | 0.17 | 0.18 | -0.1 |
| proanthocyanins | 0.31 | -0.039 | 0.15 | 0.4 | 0.14 |
| nonflavanoid_phenols | -0.3 | -0.029 | 0.17 | 0.2 | -0.5 |
| hue | 0.3 | 0.28 | 0.085 | -0.43 | -0.17 |
| proline | 0.29 | -0.36 | -0.13 | -0.23 | -0.16 |
| malic_acid | -0.25 | -0.22 | 0.089 | 0.54 | 0.035 |
| alcalinity_of_ash | -0.24 | 0.011 | 0.61 | 0.061 | 0.066 |
| alcohol | 0.14 | -0.48 | -0.21 | -0.018 | -0.27 |
| magnesium | 0.14 | -0.3 | 0.13 | -0.35 | 0.73 |
| color_intensity | -0.089 | -0.53 | -0.14 | 0.066 | -0.076 |
| ash | -0.0021 | -0.32 | 0.63 | -0.21 | -0.14 |

# Dimensionality Reduction using PCA

- In high-dimensional problems, data sometimes lies near a linear subspace, as noise introduces small variability
- Only keep data projections onto principal components with large eigenvalues
  - Can ignore the components of smaller significance
  - Might lose some info, but if eigenvalues are small, do not lose much



$$\frac{\lambda_1}{\text{tr}(C)} \qquad \frac{\lambda_2}{\text{tr}(C)}$$

Eigenfaces from 7562 images:

top left image is linear combination of rest.

Sirovich & Kirby (1987)
Turk & Pentland (1991)

# PCA wrap-up

- Strengths
  - Eigenvector method
  - No tuning of the parameters
  - No local optima

- Weaknesses
  - Limited to second order statistics
  - Limited to linear projections
    - If the structure of data is non-linear, the transformed representation can not capture the original data structure and can lose the important information.

# PCA vs LDA

| | PCA | LDA |
|---|---|---|
| **Objective** | Maximize variance in the dataset | Maximize separability between known classes |
| **Data Requirements** | Feature matrix only | Feature matrix and corresponding class labels |
| **Number of Components** | Up to the number of original features (d)` | At most C−1 where C is the number of classes (min(C-1, d)) |
| **Usage** | Dimensionality reduction, visualization, noise reduction, feature extraction | Classification and dimensionality reduction for classification |
| **Optimization Criteria** | Maximizes total variance | Maximizes ratio of between-class variance to within-class variance |

# Kernel PCA

- Non-linear feature mapping
  - PCA aims to find the subspace that consists of eigenvectors with largest eigenvalues of covariance matrix
- Covariance matrix for non-linear feature mapping
  - Consider a non-linear feature mapping $x \to \Phi(x)$

$$C_\Phi = \frac{1}{n-1} \sum_{i=1}^{n} \Phi(x_i)\Phi(x_i)^T$$

  - Eigenvectors of $C_\Phi$ are linear combinations of the feature vectors $\Phi(x_1), \dots, \Phi(x_n)$

$$\lambda v = C_\Phi v = \frac{1}{n-1} \sum_{i=1}^{n} \Phi(x_i)\Phi(x_i)^T v = \frac{1}{n-1} \sum_{i=1}^{n} \langle \Phi(x_i), v \rangle \Phi(x_i) = \lambda \sum_{i=1}^{n} \alpha_i \Phi(x_i)$$

  - $\alpha_i = \frac{1}{\lambda(n-1)} \langle \Phi(x_i), v \rangle$

# Kernel PCA

$$\langle \Phi(x_k), \Phi(x_i) \rangle = K(x_k, x_i)$$

$$\alpha_i = \frac{1}{\lambda(n-1)}\langle \Phi(x_i), v \rangle, \quad v = \sum_i \alpha_i \Phi(x_i)$$

$$\langle \Phi(x_i), v \rangle = \alpha_i \, \lambda(n-1)$$

$$\widetilde{1_n} = \begin{bmatrix} \ddots & \\ & \ddots \end{bmatrix}$$

$$\langle a, b+c \rangle = \langle a, b \rangle + \langle a, c \rangle$$

- Kernel PCA

  ○ Suppose that $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_n)$ and kernel matrix $\boldsymbol{K}$ where $K_{ij} = \langle \Phi(\boldsymbol{x}_i), \Phi(\boldsymbol{x}_j) \rangle$

$$\sum_{i=1}^{n} \langle \Phi(\boldsymbol{x}_k), \Phi(\boldsymbol{x}_i) \rangle \alpha_i = \left\langle \Phi(\boldsymbol{x}_k), \sum_{i=1}^{n} \alpha_i \Phi(\boldsymbol{x}_i) \right\rangle = \langle \Phi(\boldsymbol{x}_k), \boldsymbol{v} \rangle = (n-1)\lambda \alpha_k$$

$$\begin{bmatrix} K(x_1,x_1) & K(x_1,x_2) \\ & \ddots \end{bmatrix} \begin{bmatrix} \alpha_1 \\ \vdots \\ \alpha_n \end{bmatrix}$$

k-th row

  ○ Thus, we need to solve the kernel eigenvalue problem

$$K\boldsymbol{\alpha} = (n-1)\lambda \boldsymbol{\alpha}$$

$$(n-1)\lambda \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \vdots \end{bmatrix}$$

  ○ Normalizing Eigenvectors $\langle \boldsymbol{v}, \boldsymbol{v} \rangle = 1$ leads to the condition: $\sum_{i,j} \alpha_i \alpha_j \Phi(x_i)^T \Phi(x_j) =$

$$\sum_{i,j} \alpha_i \alpha_j k(x_i, x_j) = \boldsymbol{\alpha} K \boldsymbol{\alpha} = (n-1)\lambda \langle \boldsymbol{\alpha}, \boldsymbol{\alpha} \rangle = 1$$

$$\langle v, v \rangle$$

$$k(x) = \frac{1}{(n-1)\lambda}$$

$$\boldsymbol{\beta} = \begin{bmatrix} \text{---} \, \boldsymbol{\alpha}_1^T \, \text{---} \\ \vdots \\ \text{---} \, \boldsymbol{\alpha}_p^T \, \text{---} \end{bmatrix} \begin{bmatrix} k(\boldsymbol{x}, \boldsymbol{x}^{(1)}) \\ k(\boldsymbol{x}, \boldsymbol{x}^{(2)}) \\ \vdots \\ k(\boldsymbol{x}, \boldsymbol{x}^{(N)}) \end{bmatrix}$$

- For new test data $\boldsymbol{x}_{new}$

PC score

$$\langle v, \Phi(x_{new}) \rangle = \left\langle \sum_{i=1}^{n} \alpha_i \Phi(\boldsymbol{x}_i), \Phi(x_{new}) \right\rangle = \sum_{i=1}^{n} \alpha_i \langle \Phi(\boldsymbol{x}_i), \Phi(x_{new}) \rangle = \sum_{i=1}^{n} \alpha_i K(\boldsymbol{x}_i, \boldsymbol{x}_{new})$$

- How to obtain the centered feature map?

$$K \rightarrow \widetilde{K} \qquad K = XX^T \qquad \Phi \Phi^T$$

$$\widetilde{X} = X - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^T X = X - \widetilde{\mathbf{1}}_n X$$

$$\widetilde{K} = \widetilde{X}\widetilde{X}^T = (X - \widetilde{\mathbf{1}}_n X)(X - \widetilde{\mathbf{1}}_n X)^T$$

$$\tilde{K} = K - \tilde{1}_n K - K\tilde{1}_n + \tilde{1}_n K \tilde{1}_n$$
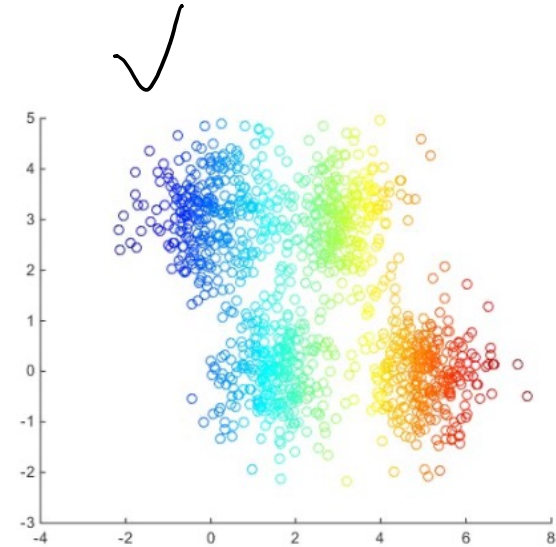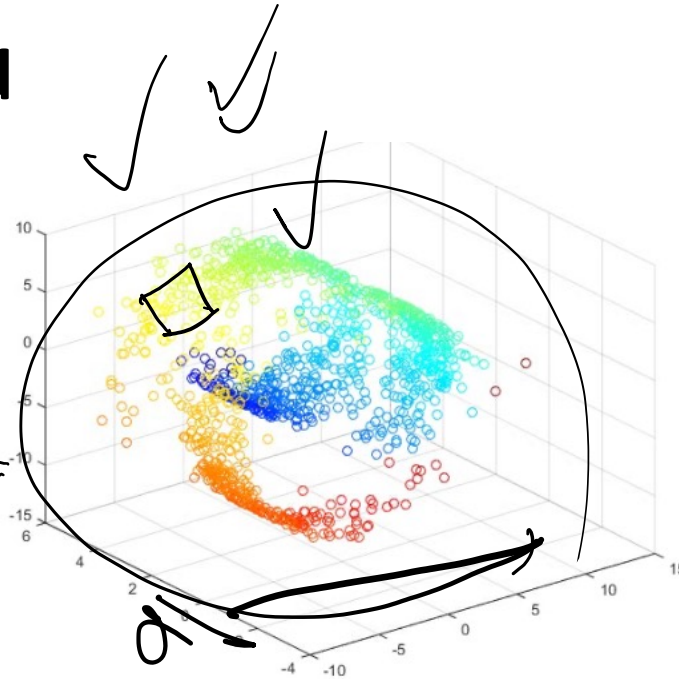
$$K = XX^T$$

# **Manifold Learning**

$$\Phi = \begin{bmatrix} - \phi(x_1) - \\ \vdots \\ - \phi(x_n) - \end{bmatrix} \in \mathbb{R}^{n \times m}$$

$$K = \Phi\Phi^T \in \mathbb{R}^{n \times n}$$

- Isomap
- Multi-dimensional Scaling
- Laplacian Eigenmap

RBF

**Data Mining**
**Prof. Saerom Park**

# Manifold



- **Key assumption**: High dimensional data actually resides in a *lower dimensional space* that is *locally Euclidean*

- Informally, the manifold is a subset of points in the high-dimensional space that locally looks like a low-dimensional space

# Similarity graph models data geometry
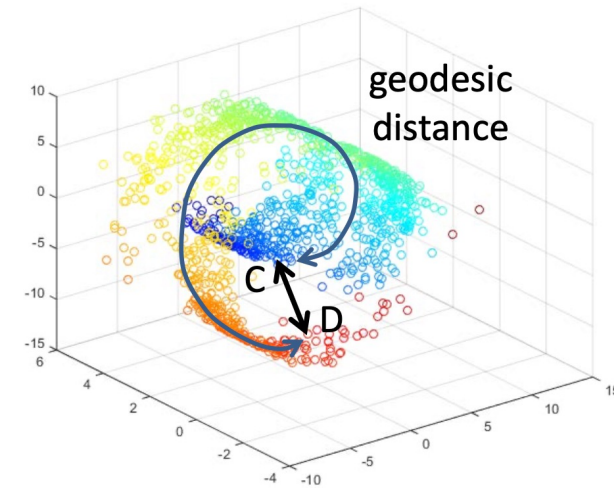


geodesic distance

- Manifold Structure
  - Global distance ignores the manifold structure of high-dimensional data
  - Geodesic distance: we can consider the distance along the manifold
- Manifold Learning
  - The idea is that (hopefully) once the manifold is "unfolded", the analysis, such as clustering becomes easy
  - In essence, "unfolding" a manifold is achieved via dimensionality reduction by considering the relationship between local data points
- Graph based on local distances    *datapoints*
  - Graph $G = \{V, E\}$, where $V$ is a set of vertices, and $E \subset V \times V$ is a set of edges
  - Graph models pairwise relations between objects (similarity or distance)
  - Weighted graph: each vertex $v_{ij}$ has an associated weight $w_{ij}$ (the strength of the relation between objects)

# Similarity graph models data geometry

- We will consider weighted undirected graphs with non-negative weights $w_{ij} \geq 0$. Moreover, we will assume that $\mathrm{w_{ij}} = 0$, if and only if vertices $i$ and $j$ are not connected
- The degree of a vertex $v_i \in \boldsymbol{V}$ is defined as

$$\deg(i) = \sum_{j=1}^{n} w_{ij}$$

- The weighted undirected graph can be represented with a weighted adjacency matrix $\boldsymbol{W}$ that contains weights $w_{ij}$ as its elements
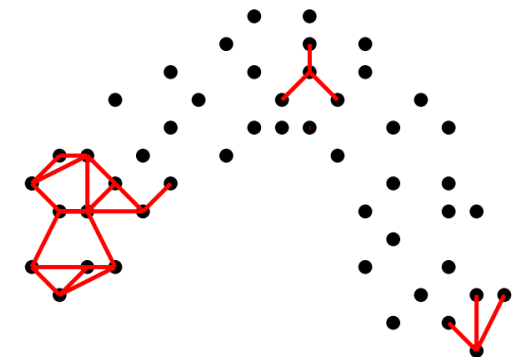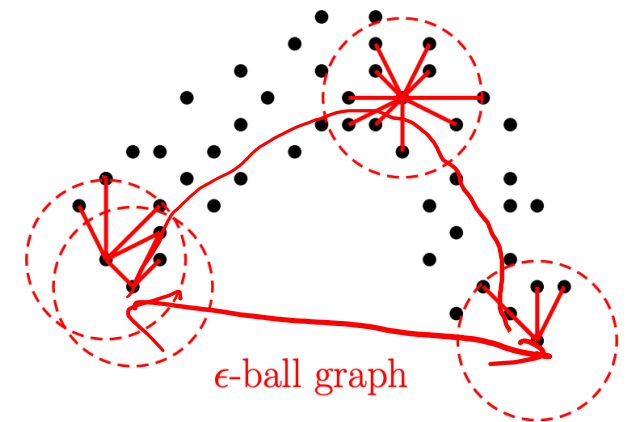
# Similarity graph models data geometry

- Geodesic distances can be approximated using a graph in which vertices represent data points
- Let $d(\boldsymbol{x}_i, \boldsymbol{x}_j)$ be the Euclidian distance between the points in the original data space

- **Option 1**: define some local radius $\epsilon$

  Connect vertices $i$ and $j$ with an edge if $d(\boldsymbol{x}_i, \boldsymbol{x}_j) \leq \epsilon$

- **Option 2**: define nearest neighbor threshold k

  Connect vertices $i$ and $j$ if $i$ is among the $k$ nearest neighbors of $j$ OR $j$ is among the $k$ nearest neighbors of $i$

$\epsilon$-ball graph

$r$NN graph ($r = 3$)

# Isomap

- Given $n$ data points $x_1, \ldots, x_n \in \mathbb{R}^d$
    1. Compute the $k$-nearest neighbor graph $G = (V, E)$
    2. Compute graph distances (geodesic distance) $d_G(x_i, x_j)$ between all points using Dijkstra's algorithm
    3. Embed the points into low dimensions using metric MDS