

# Data Mining and Machine Learning: Fundamental Concepts and Algorithms

dataminingbook.info

Mohammed J. Zaki<sup>1</sup>    Wagner Meira Jr.<sup>2</sup>

<sup>1</sup>Department of Computer Science  
Rensselaer Polytechnic Institute, Troy, NY, USA

<sup>2</sup>Department of Computer Science  
Universidade Federal de Minas Gerais, Belo Horizonte, Brazil

## Chapter 7: Dimensionality Reduction

# The curse of (attribute) dimensionality

- Data may have large number of attributes
  - Data matrix with many columns d

$$\mathbf{D} = \left( \begin{array}{c|cccc} & X_1 & X_2 & \cdots & X_d \\ \mathbf{x}_1 & x_{11} & x_{12} & \cdots & x_{1d} \\ \mathbf{x}_2 & x_{21} & x_{22} & \cdots & x_{2d} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \mathbf{x}_n & x_{n1} & x_{n2} & \cdots & x_{nd} \end{array} \right)$$

- Challenges
  - Models have large number of parameters
    - $d(d-1)/2$  variances / covariances in a Gaussian model
    - Cumbersome, more difficult to estimate well
- Dimension reduction
  - Identify small number of attributes (or combinations of attributes) that account for most of the data variation
  - Other dimensions of the data are viewed as noise,

# Dimensionality Reduction

The goal of dimensionality reduction is to find a lower dimensional representation of the data matrix  $\mathbf{D}$  to avoid the curse of dimensionality.

Given  $n \times d$  data matrix, each point  $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{id})^T$  is a vector in the ambient  $d$ -dimensional vector space spanned by the  $d$  standard basis vectors  $\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_d$ .

Given any other set of  $d$  orthonormal vectors  $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_d$  we can re-express each point  $\mathbf{x}$  as

$$\mathbf{x} = a_1 \mathbf{u}_1 + a_2 \mathbf{u}_2 + \cdots + a_d \mathbf{u}_d$$

where  $\mathbf{a} = (a_1, a_2, \dots, a_d)^T$  represents the coordinates of  $\mathbf{x}$  in the new basis. More compactly:

$$\mathbf{x} = \mathbf{U}\mathbf{a}$$

where  $\mathbf{U}$  is the  $d \times d$  orthogonal matrix, whose  $i$ th column comprises the  $i$ th basis vector  $\mathbf{u}_i$ . Thus  $\mathbf{U}^{-1} = \mathbf{U}^T$ , and we have

$$\mathbf{a} = \mathbf{U}^T \mathbf{x}$$

# Optimal Basis: Projection in Lower Dimensional Space

There are potentially infinite choices for the orthonormal basis vectors. Our goal is to choose an *optimal* basis that preserves essential information about  $\mathbf{D}$ .

We are interested in finding the optimal  $r$ -dimensional representation of  $\mathbf{D}$ , with  $r \ll d$ . Projection of  $\mathbf{x}$  onto the first  $r$  basis vectors is given as

$$\mathbf{x}' = a_1 \mathbf{u}_1 + a_2 \mathbf{u}_2 + \cdots + a_r \mathbf{u}_r = \sum_{i=1}^r a_i \mathbf{u}_i = \mathbf{U}_r \mathbf{a}_r$$

where  $\mathbf{U}_r$  and  $\mathbf{a}_r$  comprises the  $r$  basis vectors and coordinates, respv. Also, restricting  $\mathbf{a} = \mathbf{U}^T \mathbf{x}$  to  $r$  terms, we have

$$\mathbf{a}_r = \mathbf{U}_r^T \mathbf{x}$$

The  $r$ -dimensional projection of  $\mathbf{x}$  is thus given as:

$$\mathbf{x}' = \mathbf{U}_r \mathbf{U}_r^T \mathbf{x} = \mathbf{P}_r \mathbf{x}$$

where  $\mathbf{P}_r = \mathbf{U}_r \mathbf{U}_r^T = \sum_{i=1}^r \mathbf{u}_i \mathbf{u}_i^T$  is the *orthogonal projection matrix* for the subspace spanned by the first  $r$  basis vectors.

# Optimal Basis: Error Vector

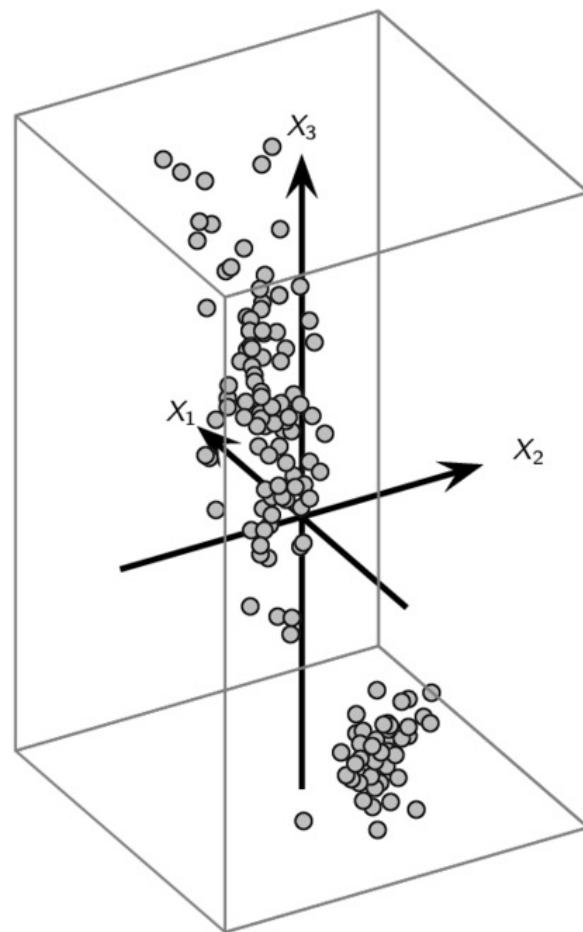
Given the projected vector  $\mathbf{x}' = \mathbf{P}_r \mathbf{x}$ , the corresponding *error vector*, is the projection onto the remaining  $d - r$  basis vectors

$$\epsilon = \sum_{i=r+1}^d a_i \mathbf{u}_i = \mathbf{x} - \mathbf{x}'$$

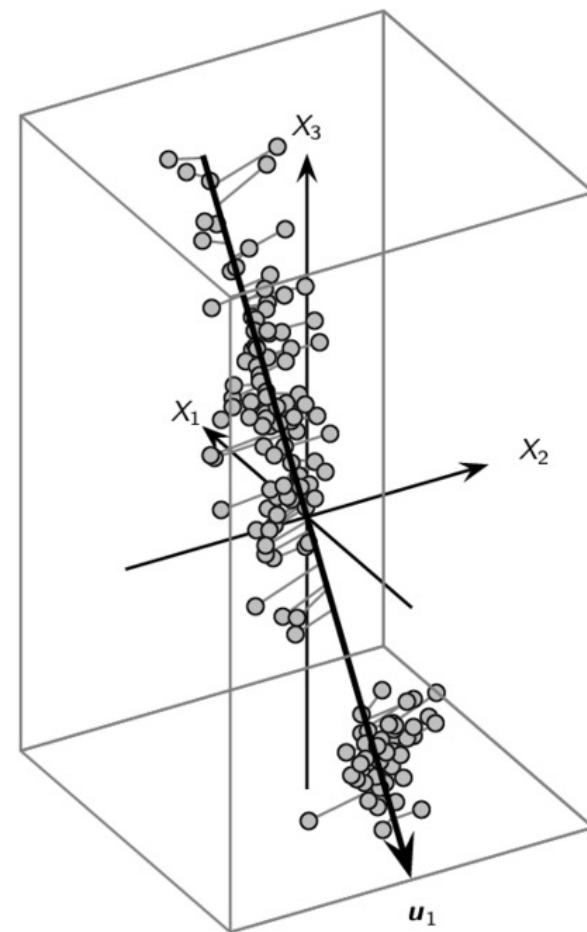
The error vector  $\epsilon$  is orthogonal to  $\mathbf{x}'$ .

The goal of dimensionality reduction is to seek an  $r$ -dimensional basis that gives the best possible approximation  $\mathbf{x}'$ ; over all the points  $\mathbf{x}_i \in \mathcal{D}$ . Alternatively, we seek to minimize the error  $\epsilon_i = \mathbf{x}_i - \mathbf{x}'_i$  over all the points.

# Iris Data: Optimal One-dimensional Basis

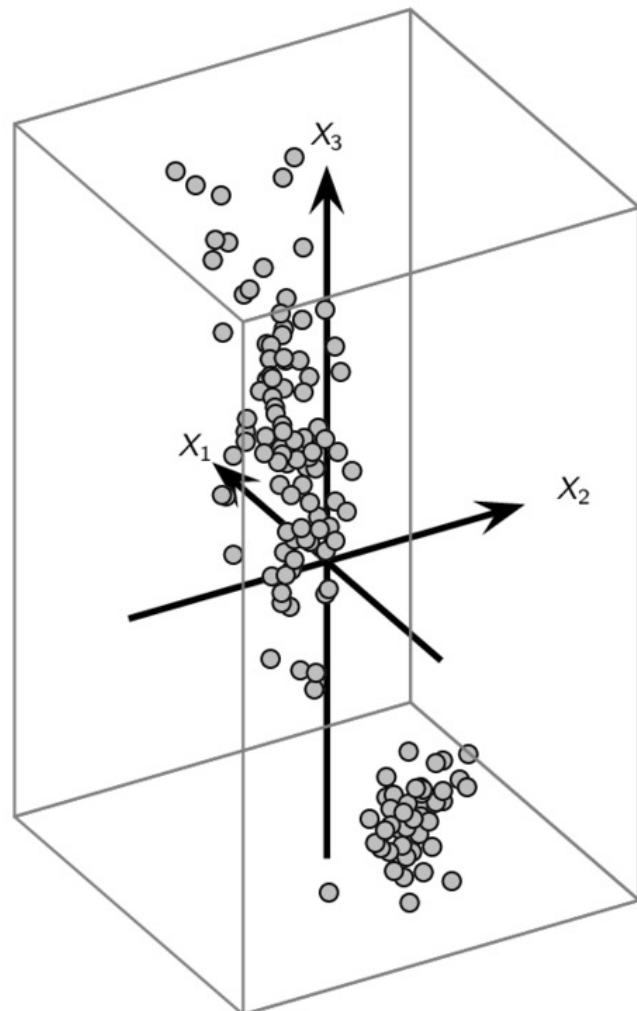


Iris Data: 3D

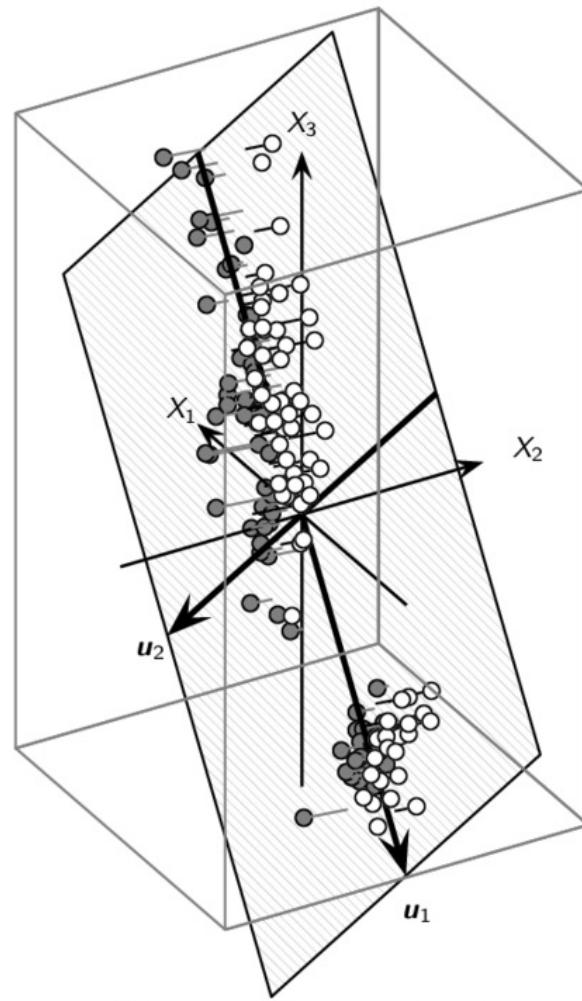


Optimal 1D Basis

# Iris Data: Optimal 2D Basis



Iris Data (3D)



Optimal 2D Basis

# Principal Component Analysis

Principal Component Analysis (PCA) is a technique that seeks a  $r$ -dimensional basis that best captures the variance in the data.

The direction with the largest projected variance is called the first principal component.

The orthogonal direction that captures the second largest projected variance is called the second principal component, and so on.

The direction that maximizes the variance is also the one that minimizes the mean squared error.

# Principal Component: Direction of Most Variance

We seek to find the unit vector  $\mathbf{u}$  that maximizes the projected variance of the points. Let  $\mathbf{D}$  be centered, and let  $\Sigma$  be its covariance matrix.

The projection of  $\mathbf{x}_i$  on  $\mathbf{u}$  is given as

$$\mathbf{x}'_i = \left( \frac{\mathbf{u}^T \mathbf{x}_i}{\mathbf{u}^T \mathbf{u}} \right) \mathbf{u} = (\mathbf{u}^T \mathbf{x}_i) \mathbf{u} = a_i \mathbf{u}$$

$$a_i = \mathbf{u}^T \mathbf{x}_i$$

Across all the points, the projected variance along  $\mathbf{u}$  is

projected mean is 0  
because  $\mathbf{D}$  is centered

$$\sigma_{\mathbf{u}}^2 = \frac{1}{n} \sum_{i=1}^n (a_i - \mu_{\mathbf{u}})^2 = \frac{1}{n} \sum_{i=1}^n \mathbf{u}^T (\mathbf{x}_i \mathbf{x}_i^T) \mathbf{u} = \mathbf{u}^T \left( \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T \right) \mathbf{u} = \mathbf{u}^T \Sigma \mathbf{u}$$

We have to find the optimal basis vector  $\mathbf{u}$  that maximizes the projected variance  $\sigma_{\mathbf{u}}^2 = \mathbf{u}^T \Sigma \mathbf{u}$ , subject to the constraint that  $\mathbf{u}^T \mathbf{u} = 1$ . The maximization objective is given as

$\alpha$  is Lagrange Multiplier used  
in constrained optimization.

$$\max_{\mathbf{u}} J(\mathbf{u}) = \mathbf{u}^T \Sigma \mathbf{u} - \alpha(\mathbf{u}^T \mathbf{u} - 1)$$

See [https://en.wikipedia.org/wiki/Lagrange\\_multipliers](https://en.wikipedia.org/wiki/Lagrange_multipliers)

# Principal Component: Direction of Most Variance

Given the objective  $\max_{\mathbf{u}} J(\mathbf{u}) = \mathbf{u}^T \Sigma \mathbf{u} - \alpha(\mathbf{u}^T \mathbf{u} - 1)$ , we solve it by setting the derivative of  $J(\mathbf{u})$  with respect to  $\mathbf{u}$  to the zero vector, to obtain

$$\frac{\partial}{\partial \mathbf{u}} (\mathbf{u}^T \Sigma \mathbf{u} - \alpha(\mathbf{u}^T \mathbf{u} - 1)) = 0$$

that is,  $2\Sigma \mathbf{u} - 2\alpha \mathbf{u} = 0$

which implies

$$\Sigma \mathbf{u} = \alpha \mathbf{u}$$

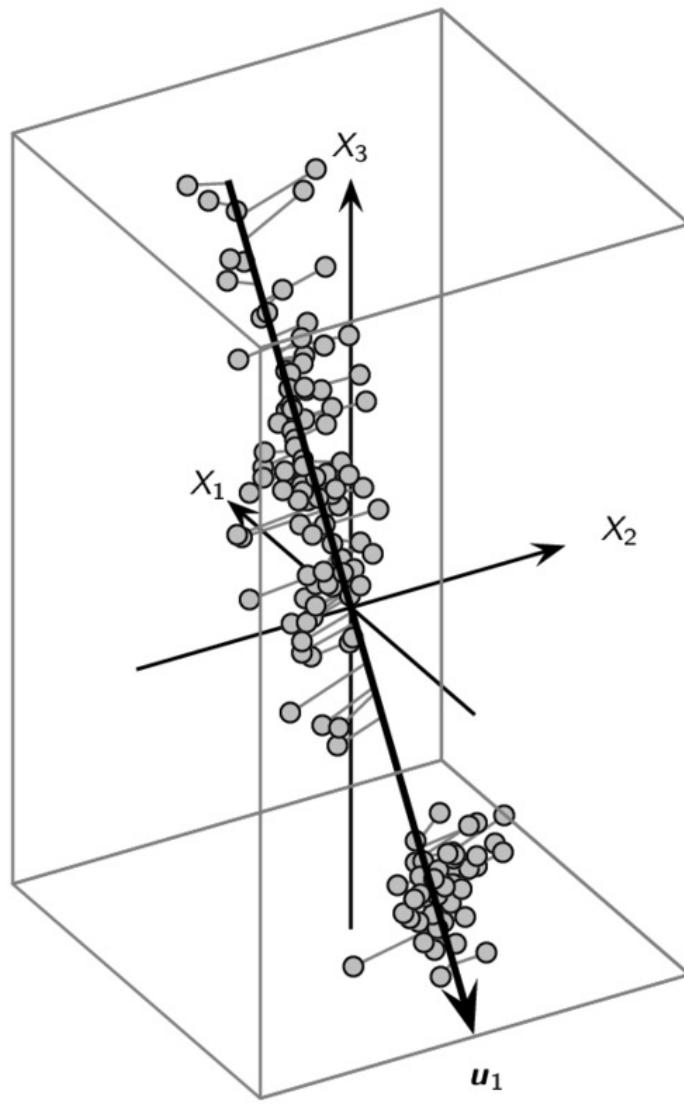
Thus  $\alpha$  is an eigenvalue of the covariance matrix  $\Sigma$ , with the associated eigenvector  $\mathbf{u}$ .

Taking the dot product with  $\mathbf{u}$  on both sides, we have

$$\sigma_{\mathbf{u}}^2 = \mathbf{u}^T \Sigma \mathbf{u} \mathbf{u}^T \alpha \mathbf{u} = \alpha \mathbf{u}^T \mathbf{u} = \alpha$$

To maximize the projected variance  $\sigma_{\mathbf{u}}^2$ , we thus choose the largest eigenvalue  $\lambda_1$  of  $\Sigma$ , and the dominant eigenvector  $\mathbf{u}_1$  specifies the direction of most variance, also called the *first principal component*.

# Iris Data: First Principal Component



# Minimum Squared Error Approach

The direction that maximizes the projected variance is also the one that minimizes the average squared error. The mean squared error (*MSE*) optimization condition is

$$MSE(\mathbf{u}) = \frac{1}{n} \sum_{i=1}^n \|\epsilon_i\|^2 = \frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i - \mathbf{x}'_i\|^2 = \sum_{i=1}^n \frac{\|\mathbf{x}_i\|^2}{n} - \mathbf{u}^T \Sigma \mathbf{u}$$

Exercise!

Since the first term is fixed for a dataset  $\mathcal{D}$ , we see that the direction  $\mathbf{u}_1$  that maximizes the variance is also the one that minimizes the MSE. Further,

$$\sum_{i=1}^n \frac{\|\mathbf{x}_i\|^2}{n} - \mathbf{u}^T \Sigma \mathbf{u} = var(\mathcal{D}) = tr(\Sigma) = \sum_{i=1}^d \sigma_i^2$$

Thus, for the direction  $\mathbf{u}_1$  that minimizes MSE, we have

$$MSE(\mathbf{u}_1) = var(\mathcal{D}) - \mathbf{u}_1^T \Sigma \mathbf{u}_1 = var(\mathcal{D}) - \lambda_1$$

# Best 2-dimensional Approximation

The best 2D subspace that captures the most variance in  $\mathbf{D}$  comprises the eigenvectors  $\mathbf{u}_1$  and  $\mathbf{u}_2$  corresponding to the largest and second largest eigenvalues  $\lambda_1$  and  $\lambda_2$ , respv.

Let  $\mathbf{U}_2 = (\mathbf{u}_1 \quad \mathbf{u}_2)$  be the matrix whose columns correspond to the two principal components. Given the point  $\mathbf{x}_i \in \mathbf{D}$  its projected coordinates are computed as follows:

$$\mathbf{a}_i = \mathbf{U}_2^T \mathbf{x}_i$$

Let  $\mathbf{A}$  denote the projected 2D dataset. The total projected variance for  $\mathbf{A}$  is given as

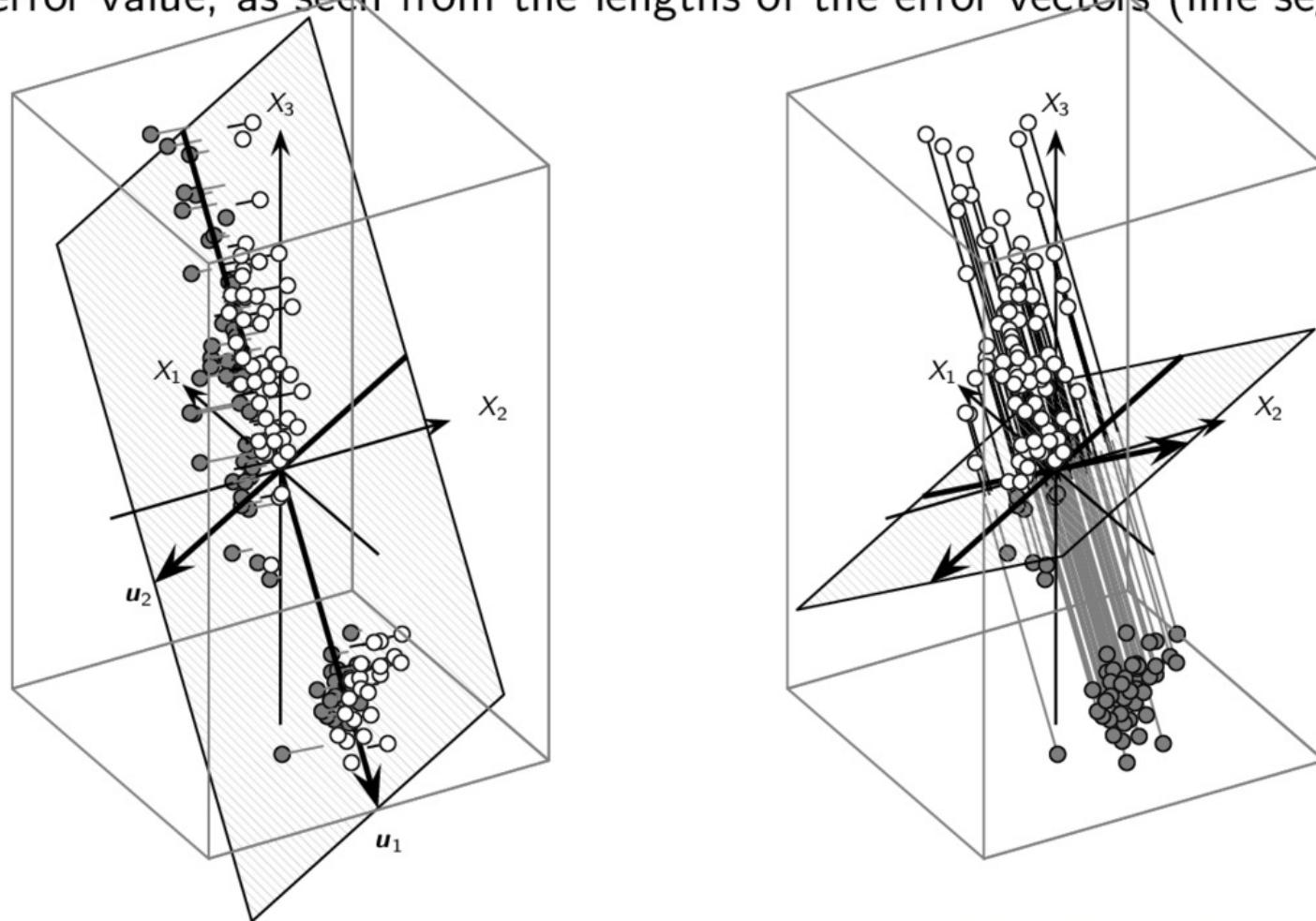
$$var(\mathbf{A}) = \mathbf{u}_1^T \Sigma \mathbf{u}_1 + \mathbf{u}_2^T \Sigma \mathbf{u}_2 = \mathbf{u}_1^T \lambda_1 \mathbf{u}_1 + \mathbf{u}_2^T \lambda_2 \mathbf{u}_2 = \lambda_1 + \lambda_2$$

The first two principal components also minimize the mean square error objective, since

$$MSE = \frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i - \mathbf{x}'_i\|^2 = var(\mathbf{D}) - \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i^T \mathbf{P}_2 \mathbf{x}_i) = var(\mathbf{D}) - var(\mathbf{A})$$

# Optimal and Non-optimal 2D Approximations

The optimal subspace maximizes the variance, and minimizes the squared error, whereas the nonoptimal subspace captures less variance, and has a high mean squared error value, as seen from the lengths of the error vectors (line segments).



# Best $r$ -dimensional Approximation

To find the best  $r$ -dimensional approximation to  $\mathbf{D}$ , we compute the eigenvalues of  $\Sigma$ . Because  $\Sigma$  is positive semidefinite, its eigenvalues are non-negative

$$\lambda_1 \geq \lambda_2 \geq \cdots \lambda_r \geq \lambda_{r+1} \cdots \geq \lambda_d \geq 0$$

We select the  $r$  largest eigenvalues, and their corresponding eigenvectors to form the best  $r$ -dimensional approximation.

**Total Projected Variance:** Let  $\mathbf{U}_r = (\mathbf{u}_1 \ \cdots \ \mathbf{u}_r)$  be the  $r$ -dimensional basis vector matrix, with the projection matrix given as  $\mathbf{P}_r = \mathbf{U}_r \mathbf{U}_r^T = \sum_{i=1}^r \mathbf{u}_i \mathbf{u}_i^T$ .

Let  $\mathbf{A}$  denote the dataset formed by the coordinates of the projected points in the  $r$ -dimensional subspace. The projected variance is given as

$$var(\mathbf{A}) = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i^T \mathbf{P}_r \mathbf{x}_i = \sum_{i=1}^r \mathbf{u}_i^T \Sigma \mathbf{u}_i = \sum_{i=1}^r \lambda_i$$

**Mean Squared Error:** The mean squared error in  $r$  dimensions is

$$MSE = \frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i - \mathbf{x}'_i\|^2 = var(\mathbf{D}) - \sum_{i=1}^r \lambda_i = \sum_{i=1}^d \lambda_i - \sum_{i=1}^r \lambda_i$$

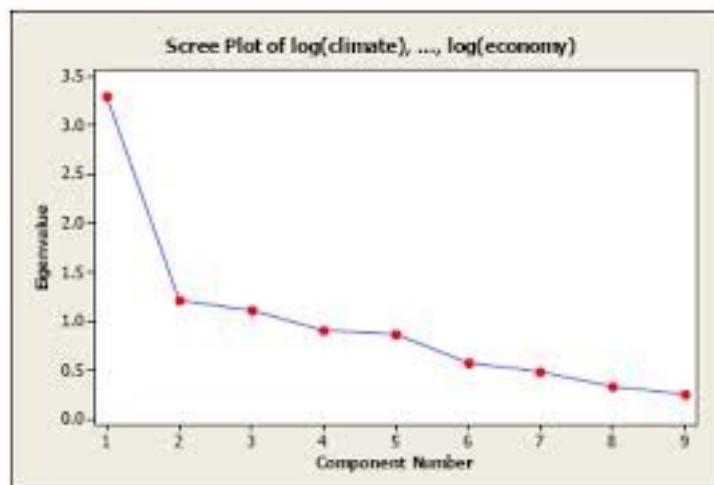
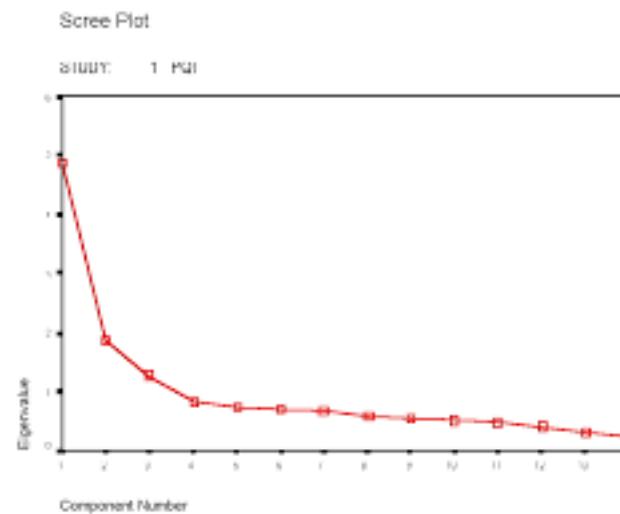
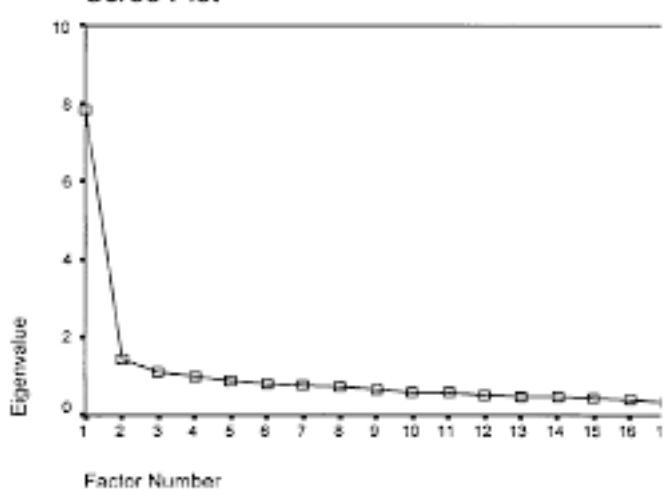
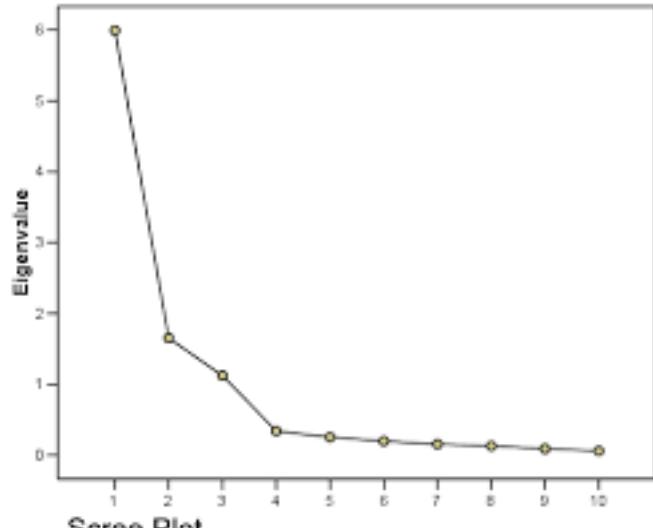
# Choosing the Dimensionality

One criteria for choosing  $r$  is to compute the fraction of the total variance captured by the first  $r$  principal components, computed as

$$f(r) = \frac{\lambda_1 + \lambda_2 + \cdots + \lambda_r}{\lambda_1 + \lambda_2 + \cdots + \lambda_d} = \frac{\sum_{i=1}^r \lambda_i}{\sum_{i=1}^d \lambda_i} = \frac{\sum_{i=1}^r \lambda_i}{var(\mathbf{D})}$$

Given a certain desired variance threshold, say  $\alpha$ , starting from the first principal component, we keep on adding additional components, and stop at the smallest value  $r$ , for which  $f(r) \geq \alpha$ . In other words, we select the fewest number of dimensions such that the subspace spanned by those  $r$  dimensions captures at least  $\alpha$  fraction (say 0.9) of the total variance.

# Scree Plots: eigenvalues in decreasing order



# Principal Component Analysis: Algorithm

**PCA ( $D, \alpha$ ):**

- 1  $\mu = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i$  // compute mean
- 2  $\mathbf{Z} = \mathbf{D} - \mathbf{1} \cdot \mu^T$  // center the data
- 3  $\Sigma = \frac{1}{n} (\mathbf{Z}^T \mathbf{Z})$  // compute covariance matrix
- 4  $(\lambda_1, \lambda_2, \dots, \lambda_d) = \text{eigenvalues}(\Sigma)$  // compute eigenvalues
- 5  $\mathbf{U} = (\mathbf{u}_1 \quad \mathbf{u}_2 \quad \cdots \quad \mathbf{u}_d) = \text{eigenvectors}(\Sigma)$  // compute eigenvectors
- 6  $f(r) = \frac{\sum_{i=1}^r \lambda_i}{\sum_{i=1}^d \lambda_i}$ , for all  $r = 1, 2, \dots, d$  // fraction of total variance
- 7 Choose smallest  $r$  so that  $f(r) \geq \alpha$  // choose dimensionality
- 8  $\mathbf{U}_r = (\mathbf{u}_1 \quad \mathbf{u}_2 \quad \cdots \quad \mathbf{u}_r)$  // reduced basis
- 9  $\mathbf{A} = \{\mathbf{a}_i \mid \mathbf{a}_i = \mathbf{U}_r^T \mathbf{x}_i, \text{for } i = 1, \dots, n\}$  // reduced dimensionality data

# Iris Principal Components

Covariance matrix:

$$\Sigma = \begin{pmatrix} 0.681 & -0.039 & 1.265 \\ -0.039 & 0.187 & -0.320 \\ 1.265 & -0.32 & 3.092 \end{pmatrix}$$

The eigenvalues and eigenvectors of  $\Sigma$

$$\lambda_1 = 3.662$$

$$\lambda_2 = 0.239$$

$$\lambda_3 = 0.059$$

$$\mathbf{u}_1 = \begin{pmatrix} -0.390 \\ 0.089 \\ -0.916 \end{pmatrix}$$

$$\mathbf{u}_2 = \begin{pmatrix} -0.639 \\ -0.742 \\ 0.200 \end{pmatrix}$$

$$\mathbf{u}_3 = \begin{pmatrix} -0.663 \\ 0.664 \\ 0.346 \end{pmatrix}$$

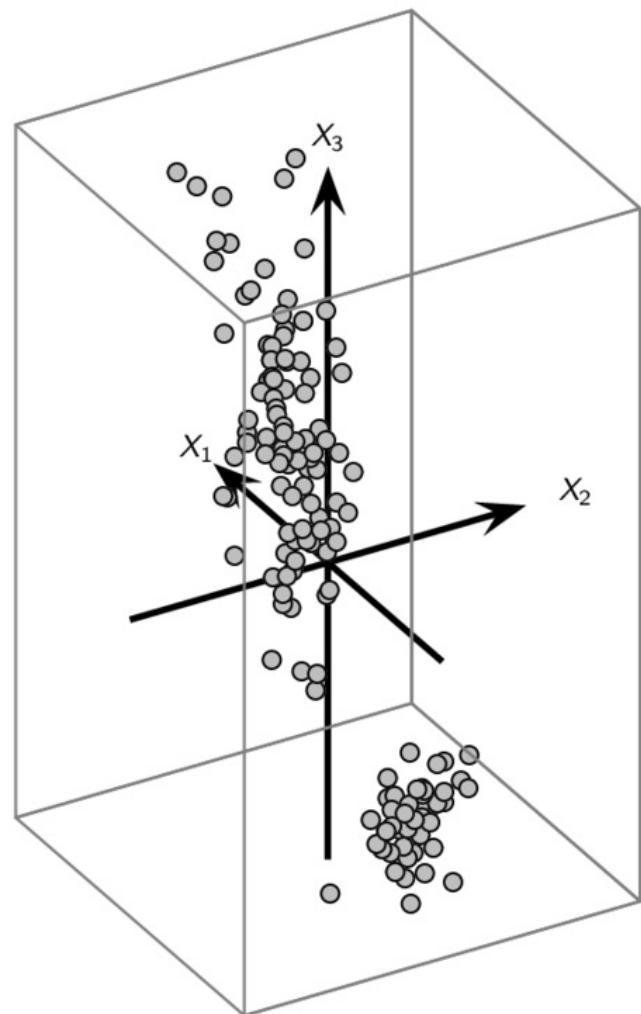
The total variance is therefore  $\lambda_1 + \lambda_2 + \lambda_3 = 3.662 + 0.239 + 0.059 = 3.96$ .

The fraction of total variance for different values of  $r$  is given as

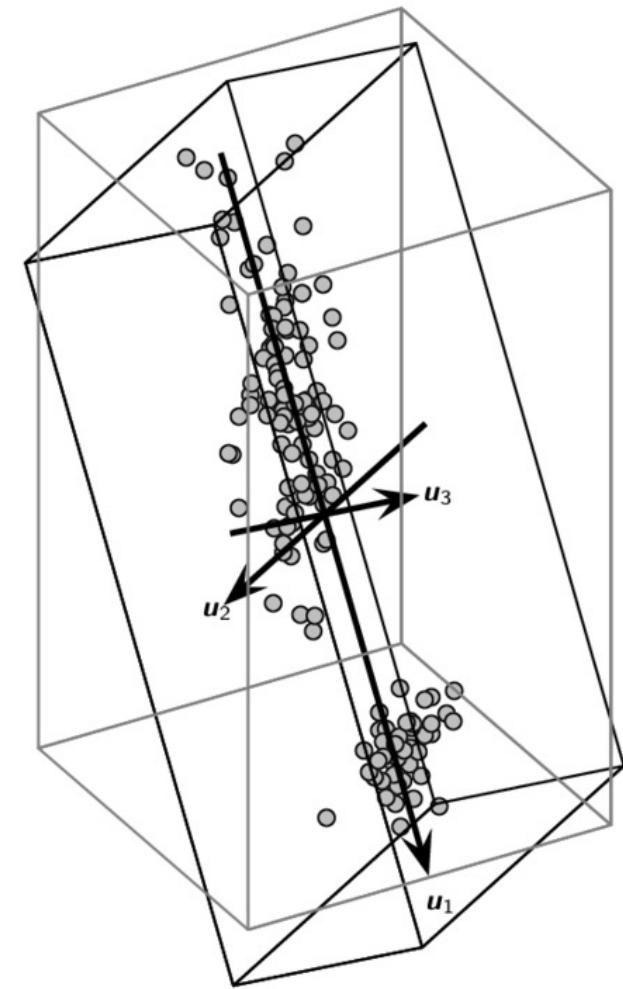
$r$	1	2	3
$f(r)$	0.925	0.985	1.0

This  $r = 2$  PCs are need to capture  $\alpha = 0.95$  fraction of variance.

# Iris Data: Optimal 3D PC Basis

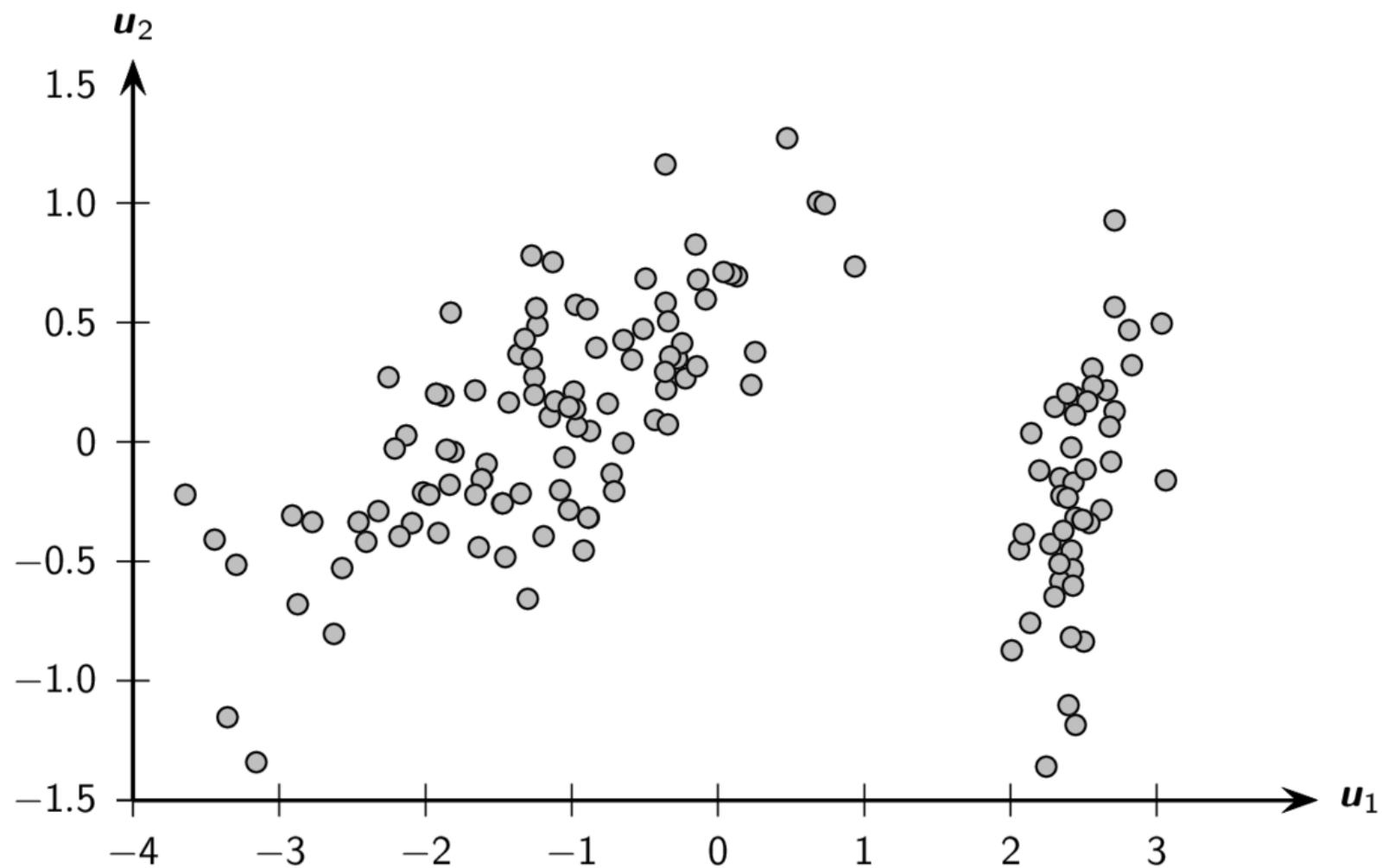


Iris Data (3D)



Optimal 3D Basis

# Iris Principal Components: Projected Data (2D)



# Geometry of PCA

Geometrically, when  $r = d$ , PCA corresponds to a orthogonal change of basis, so that the total variance is captured by the sum of the variances along each of the principal directions  $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_d$ , and further, all covariances are zero.

Let  $\mathbf{U}$  be the  $d \times d$  orthogonal matrix  $\mathbf{U} = (\mathbf{u}_1 \quad \mathbf{u}_2 \quad \cdots \quad \mathbf{u}_d)$ , with  $\mathbf{U}^{-1} = \mathbf{U}^T$ . Let  $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_d)$  be the diagonal matrix of eigenvalues. Each principal component  $\mathbf{u}_i$  corresponds to an eigenvector of the covariance matrix  $\Sigma$

$$\Sigma \mathbf{u}_i = \lambda_i \mathbf{u}_i \text{ for all } 1 \leq i \leq d$$

which can be written compactly in matrix notation:

$$\Sigma \mathbf{U} = \mathbf{U} \Lambda \text{ which implies } \Sigma = \mathbf{U} \Lambda \mathbf{U}^T$$

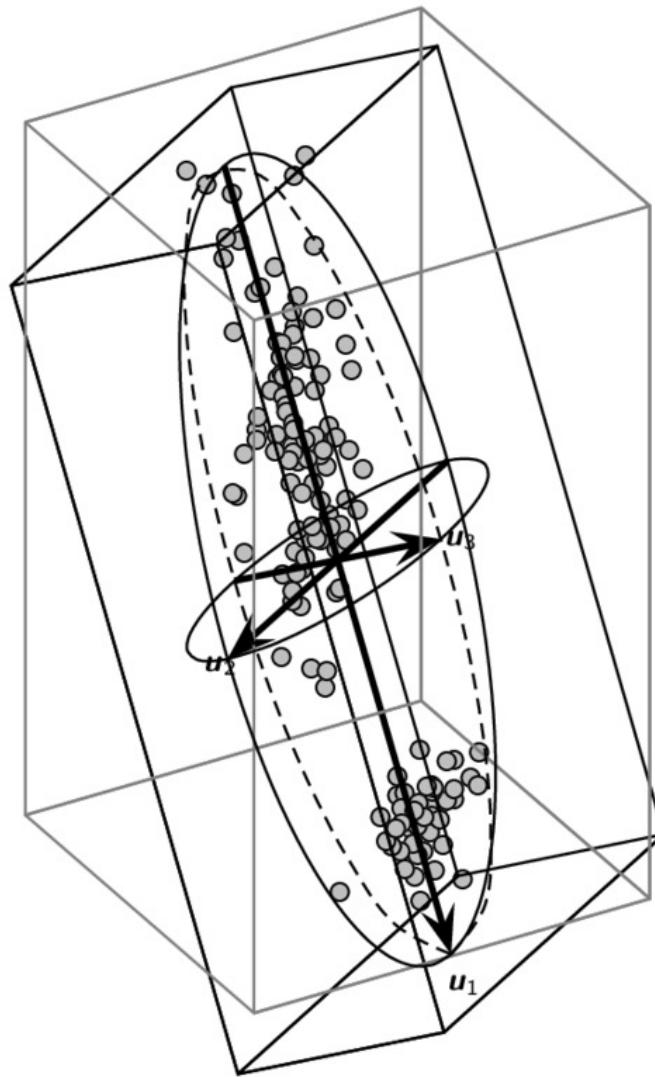
Thus,  $\Lambda$  represents the covariance matrix in the new PC basis.

In the new PC basis, the equation

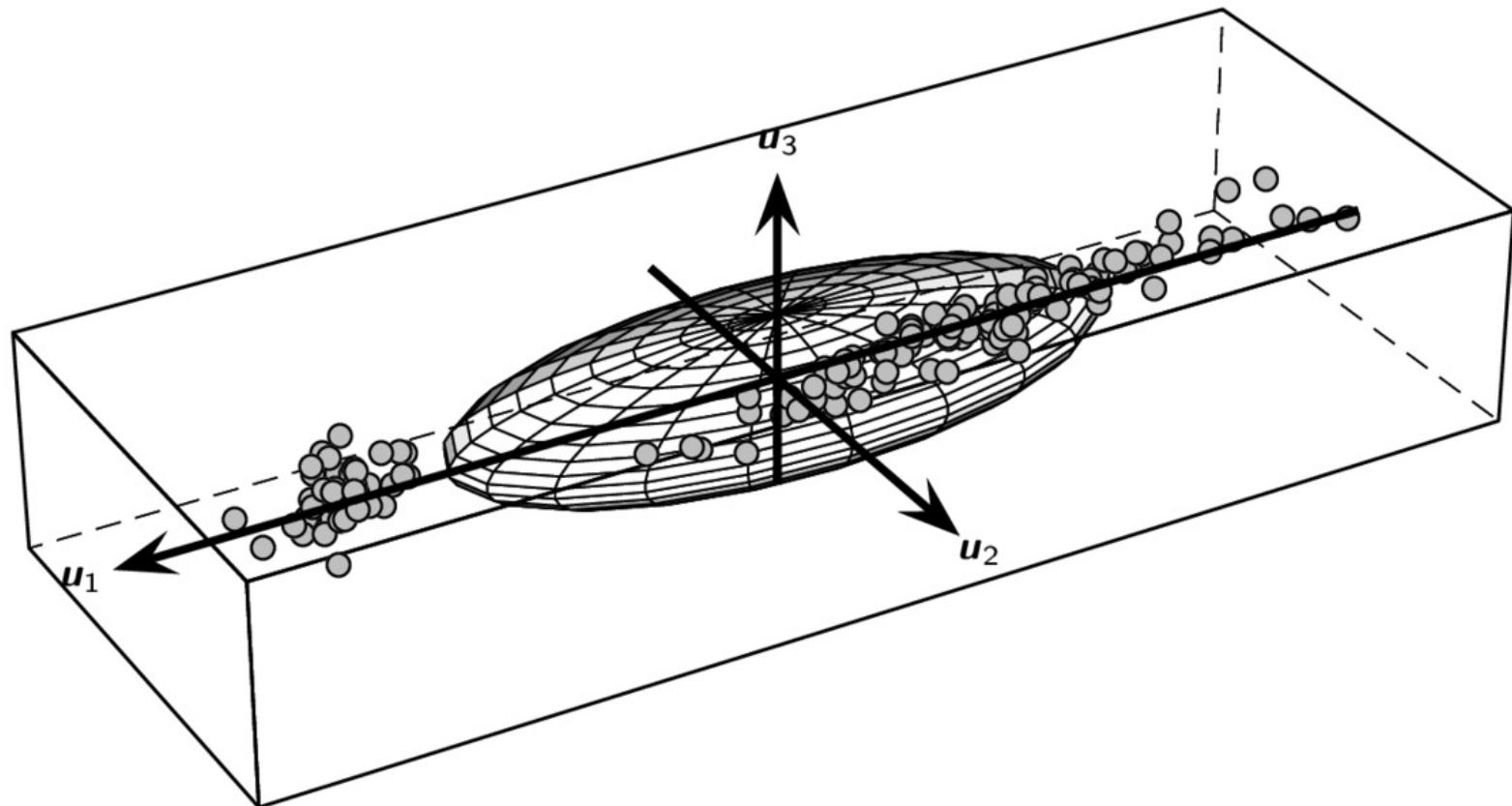
$$\mathbf{x}^T \Sigma^{-1} \mathbf{x} = 1$$

defines a  $d$ -dimensional ellipsoid (or hyper-ellipse). The eigenvectors  $\mathbf{u}_i$  of  $\Sigma$ , that is, the principal components, are the directions for the principal axes of the ellipsoid. The square roots of the eigenvalues, that is,  $\sqrt{\lambda_i}$ , give the lengths of the semi-axes.

# Iris: Elliptic Contours in Standard Basis



# Iris: Axis-Parallel Ellipsoid in PC Basis



Stop here!

# Kernel Principal Component Analysis

Principal component analysis can be extended to find nonlinear “directions” in the data using kernel methods. Kernel PCA finds the directions of most variance in the feature space instead of the input space. Using the *kernel trick*, all PCA operations can be carried out in terms of the kernel function in input space, without having to transform the data into feature space.

Let  $\phi$  be a function that maps a point  $\mathbf{x}$  in input space to its image  $\phi(\mathbf{x}_i)$  in feature space. Let the points in feature space be centered and let  $\Sigma_\phi$  be the covariance matrix. The first PC in feature space correspond to the dominant eigenvector

$$\Sigma_\phi \mathbf{u}_1 = \lambda_1 \mathbf{u}_1$$

where

$$\Sigma_\phi = \frac{1}{n} \sum_{i=1}^n \phi(\mathbf{x}_i) \phi(\mathbf{x}_i)^T$$

# Kernel Principal Component Analysis

It can be shown that  $\mathbf{u}_1 = \sum_{i=1}^n c_i \phi(\mathbf{x}_i)$ . That is, the PC direction in feature space is a linear combination of the transformed points.

The coefficients are captured in the weight vector

$$\mathbf{c} = (c_1, c_2, \dots, c_n)^T$$

Substituting into the eigen-decomposition of  $\Sigma_\phi$  and simplifying, we get:

$$\mathbf{K}\mathbf{c} = n\lambda_1 \mathbf{c} = \eta_1 \mathbf{c}$$

Thus, the weight vector  $\mathbf{c}$  is the eigenvector corresponding to the largest eigenvalue  $\eta_1$  of the kernel matrix  $\mathbf{K}$ .

# Kernel Principal Component Analysis

The weight vector  $\mathbf{c}$  can be used to then find  $\mathbf{u}_1$  via  $\mathbf{u}_1 = \sum_{i=1}^n c_i \phi(\mathbf{x}_i)$ .

The only constraint we impose is that  $\mathbf{u}_1$  should be normalized to be a unit vector, which implies  $\|\mathbf{c}\|^2 = \frac{1}{\eta_1}$ .

We cannot compute directly the principal direction, but we can project any point  $\phi(\mathbf{x})$  onto the principal direction  $\mathbf{u}_1$ , as follows:

$$\mathbf{u}_1^T \phi(\mathbf{x}) = \sum_{i=1}^n c_i \phi(\mathbf{x}_i)^T \phi(\mathbf{x}) = \sum_{i=1}^n c_i K(\mathbf{x}_i, \mathbf{x})$$

which requires only kernel operations.

We can obtain the additional principal components by solving for the other eigenvalues and eigenvectors of

$$\mathbf{K}\mathbf{c}_j = n\lambda_j \mathbf{c}_j = \eta_j \mathbf{c}_j$$

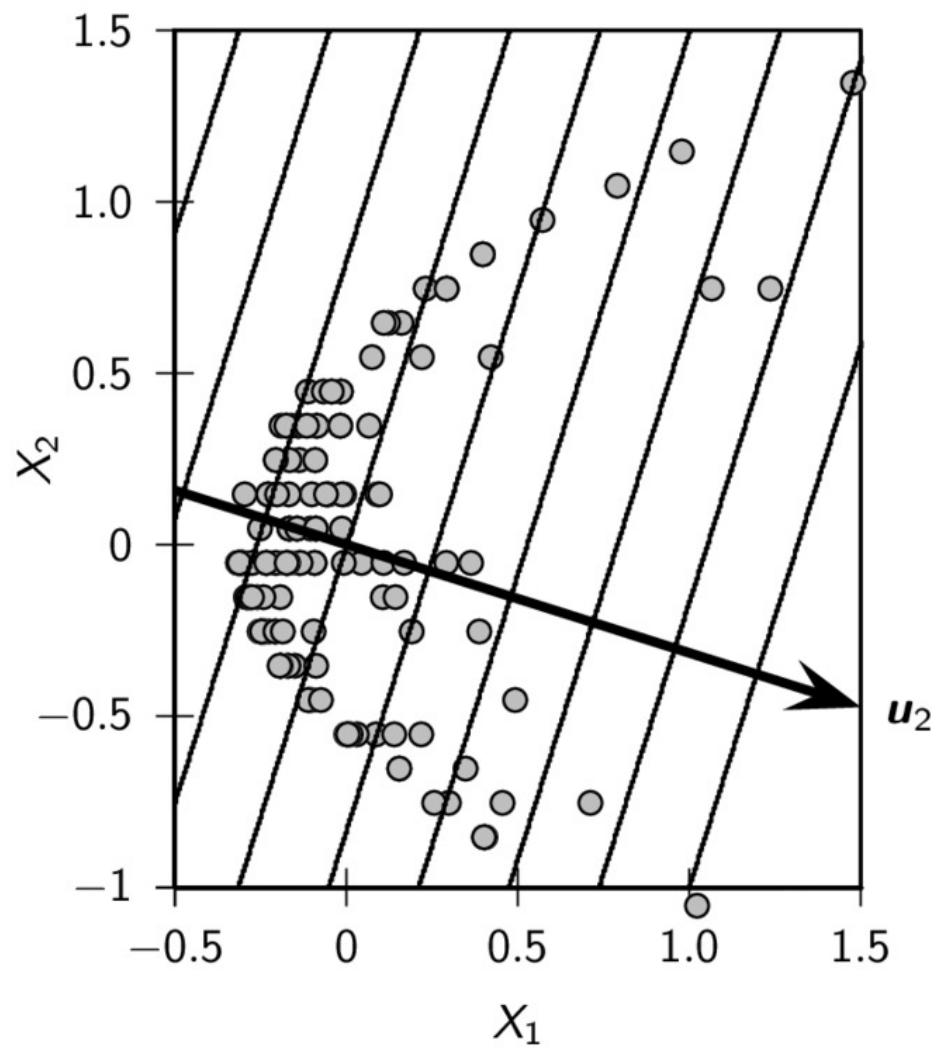
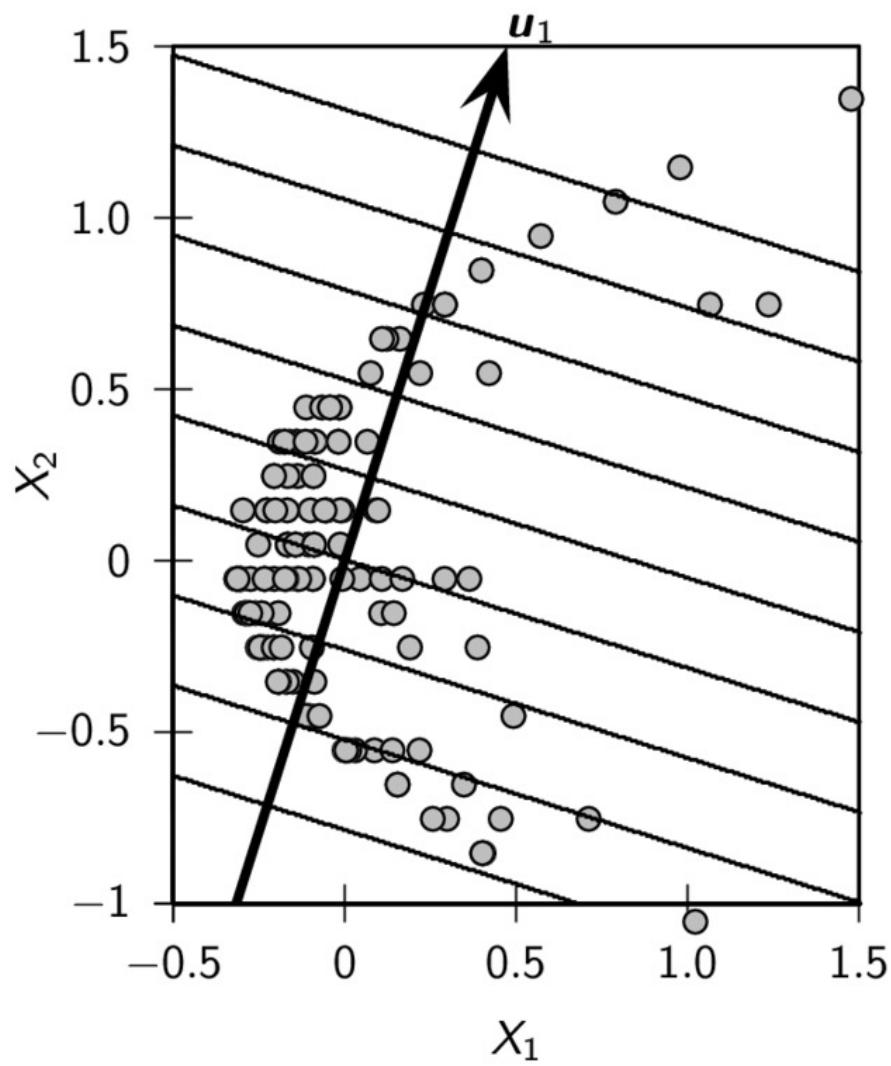
If we sort the eigenvalues of  $\mathbf{K}$  in decreasing order  $\eta_1 \geq \eta_2 \geq \dots \geq \eta_n \geq 0$ , we can obtain the  $j$ th principal component as the corresponding eigenvector  $\mathbf{c}_j$ . The variance along the  $j$ th principal component is given as  $\lambda_j = \frac{\eta_j}{n}$ .

# Kernel PCA Algorithm

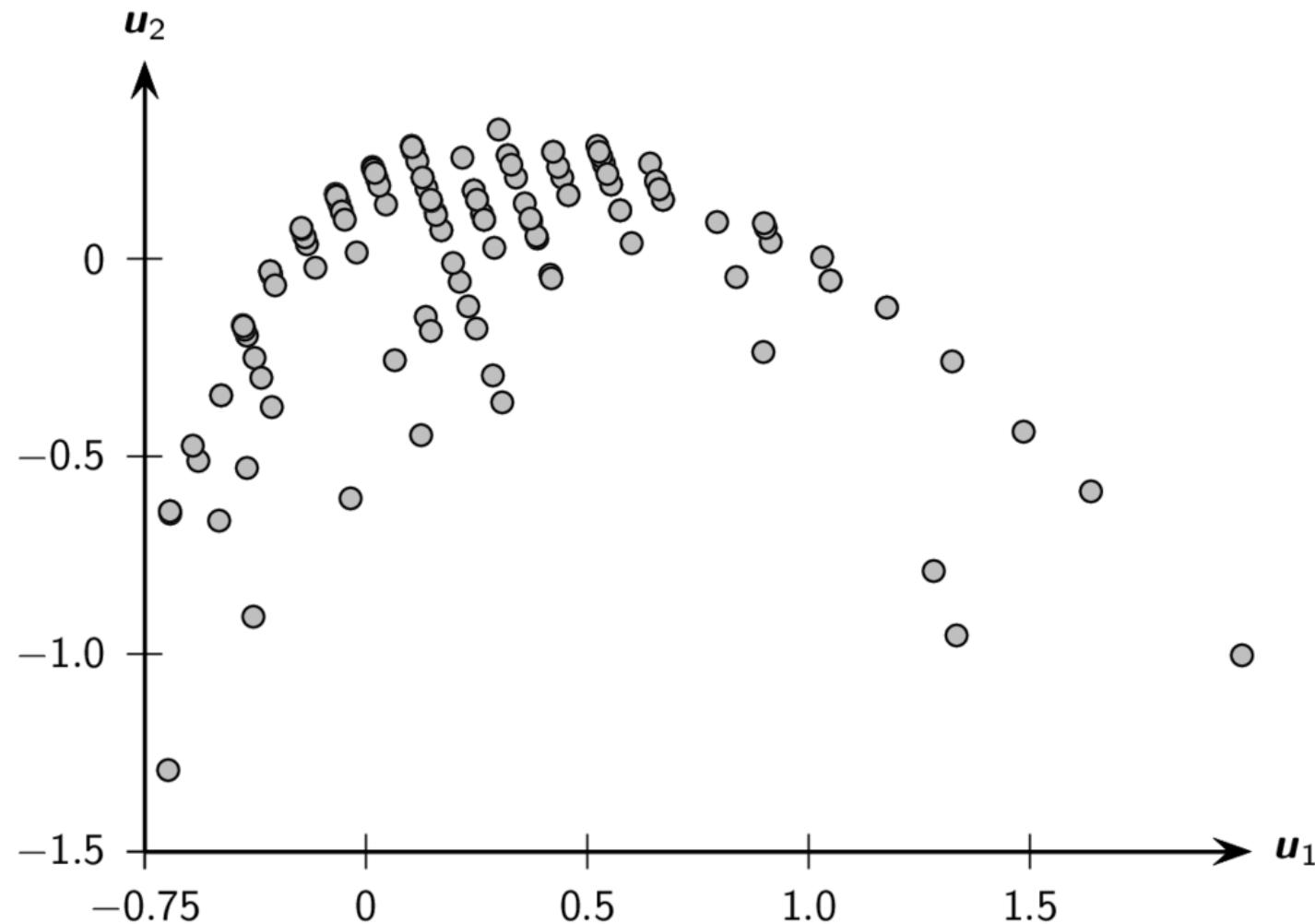
**KernelPCA ( $D, K, \alpha$ ):**

- 1  $\mathbf{K} = \{K(\mathbf{x}_i, \mathbf{x}_j)\}_{i,j=1,\dots,n}$  // compute  $n \times n$  kernel matrix
- 2  $\mathbf{K} = (\mathbf{I} - \frac{1}{n}\mathbf{1}_{n \times n})\mathbf{K}(\mathbf{I} - \frac{1}{n}\mathbf{1}_{n \times n})$  // center the kernel matrix
- 3  $(\eta_1, \eta_2, \dots, \eta_d) = \text{eigenvalues}(\mathbf{K})$  // compute eigenvalues
- 4  $(\mathbf{c}_1 \quad \mathbf{c}_2 \quad \dots \quad \mathbf{c}_n) = \text{eigenvectors}(\mathbf{K})$  // compute eigenvectors
- 5  $\lambda_i = \frac{\eta_i}{n}$  for all  $i = 1, \dots, n$  // compute variance for each component
- 6  $\mathbf{c}_i = \sqrt{\frac{1}{\eta_i}} \cdot \mathbf{c}_i$  for all  $i = 1, \dots, n$  // ensure that  $\mathbf{u}_i^T \mathbf{u}_i = 1$
- 7  $f(r) = \frac{\sum_{i=1}^r \lambda_i}{\sum_{i=1}^d \lambda_i}$ , for all  $r = 1, 2, \dots, d$  // fraction of total variance
- 8 Choose smallest  $r$  so that  $f(r) \geq \alpha$  // choose dimensionality
- 9  $\mathbf{C}_r = (\mathbf{c}_1 \quad \mathbf{c}_2 \quad \dots \quad \mathbf{c}_r)$  // reduced basis
- 10  $\mathbf{A} = \{\mathbf{a}_i \mid \mathbf{a}_i = \mathbf{C}_r^T \mathbf{K}_i, \text{for } i = 1, \dots, n\}$  // reduced dimensionality data

# Nonlinear Iris Data: PCA in Input Space

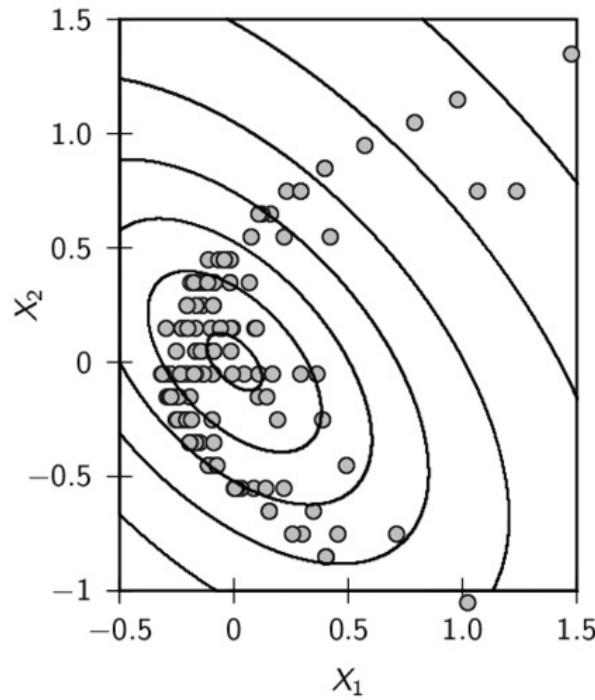


# Nonlinear Iris Data: Projection onto PCs

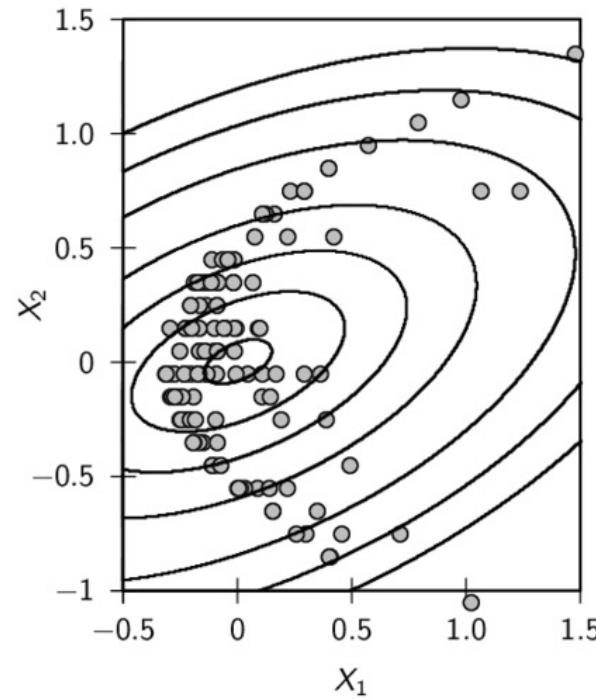


# Kernel PCA: 3 PCs (Contours of Constant Projection)

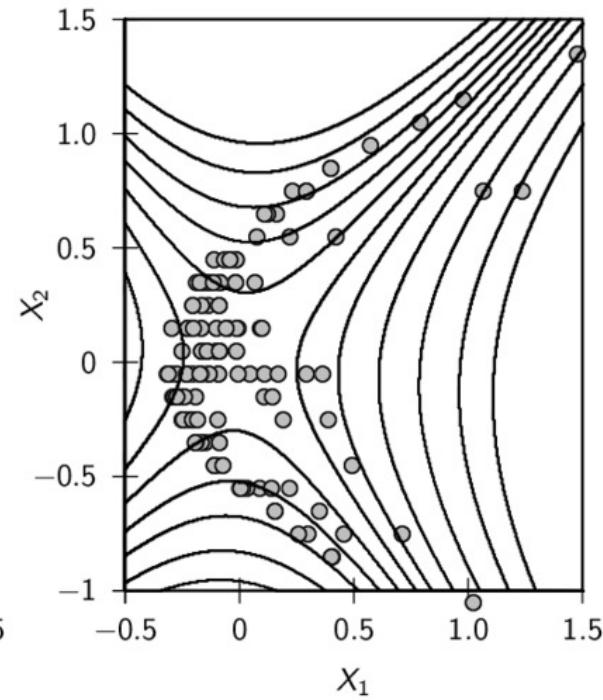
Homogeneous Quadratic Kernel:  $K(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i^T \mathbf{x}_j)^2$



(a)  $\lambda_1 = 0.2067$



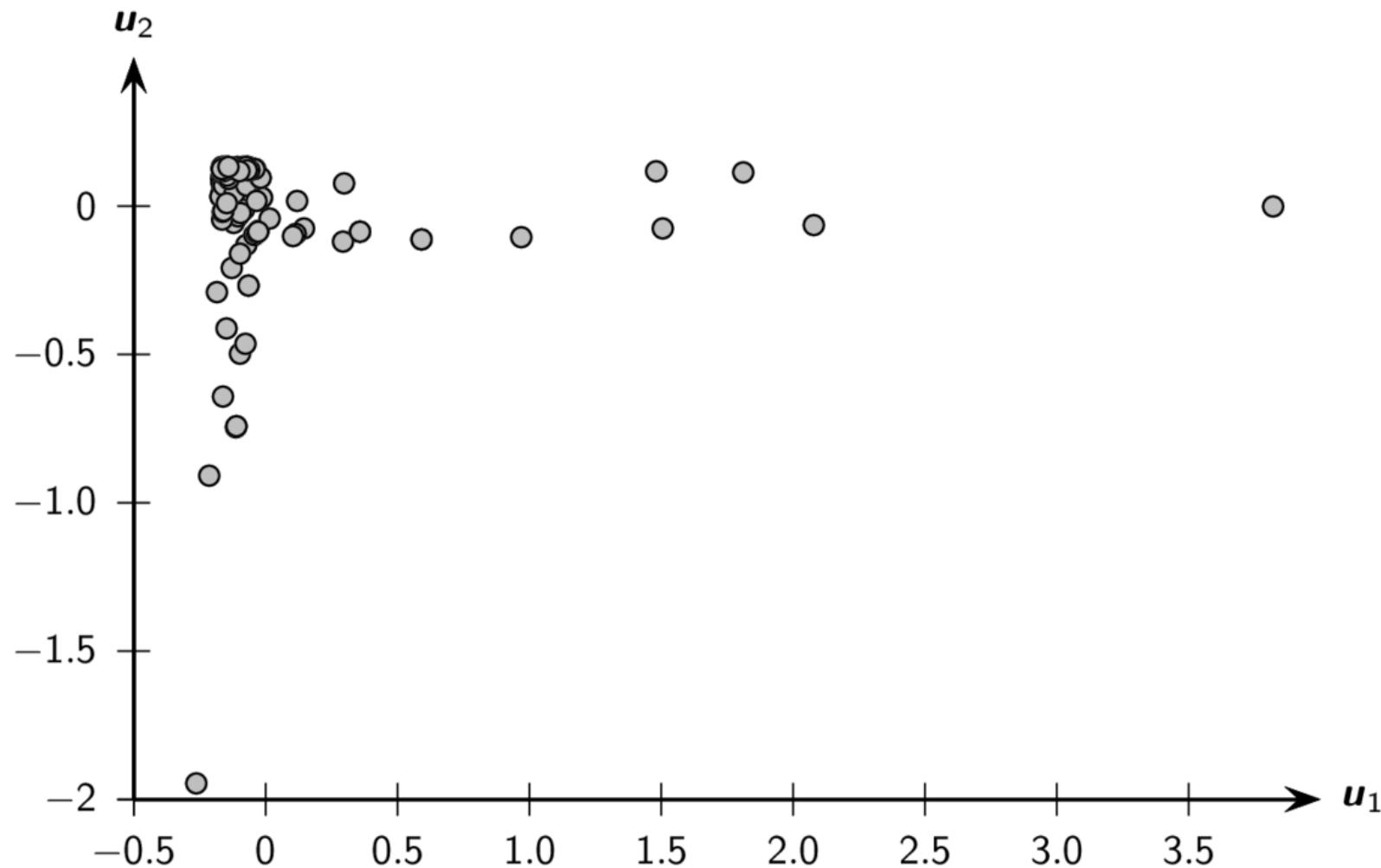
(b)  $\lambda_2 = 0.0596$



(c)  $\lambda_3 = 0.0184$

# Kernel PCA: Projected Points onto 2 PCs

Homogeneous Quadratic Kernel:  $K(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i^T \mathbf{x}_j)^2$



# Singular Value Decomposition

Principal components analysis is a special case of a more general matrix decomposition method called *Singular Value Decomposition (SVD)*. PCA yields the following decomposition of the covariance matrix:

$$\Sigma = \mathbf{U} \Lambda \mathbf{U}^T$$

where the covariance matrix has been factorized into the orthogonal matrix  $\mathbf{U}$  containing its eigenvectors, and a diagonal matrix  $\Lambda$  containing its eigenvalues (sorted in decreasing order).

SVD generalizes the above factorization for any matrix. In particular for an  $n \times d$  data matrix  $\mathbf{D}$  with  $n$  points and  $d$  columns, SVD factorizes  $\mathbf{D}$  as follows:

$$\mathbf{D} = \mathbf{L} \Delta \mathbf{R}^T$$

where  $\mathbf{L}$  is a orthogonal  $n \times n$  matrix,  $\mathbf{R}$  is an orthogonal  $d \times d$  matrix, and  $\Delta$  is an  $n \times d$  “diagonal” matrix, defined as  $\Delta(i, i) = \delta_i$ , and 0 otherwise. The columns of  $\mathbf{L}$  are called the *left singular and the columns of  $\mathbf{R}$  (or rows of  $\mathbf{R}^T$ ) are called the right singular vectors*. The entries  $\delta_i$  are called the *singular values* of  $\mathbf{D}$ , and they are all non-negative.

# Reduced SVD

If the rank of  $\mathbf{D}$  is  $r \leq \min(n, d)$ , then there are only  $r$  nonzero singular values, ordered as follows:  $\delta_1 \geq \delta_2 \geq \dots \geq \delta_r > 0$ .

We discard the left and right singular vectors that correspond to zero singular values, to obtain the *reduced SVD* as

$$\mathbf{D} = \mathbf{L}_r \Delta_r \mathbf{R}_r^T$$

where  $\mathbf{L}_r$  is the  $n \times r$  matrix of the left singular vectors,  $\mathbf{R}_r$  is the  $d \times r$  matrix of the right singular vectors, and  $\Delta_r$  is the  $r \times r$  diagonal matrix containing the positive singular vectors.

The reduced SVD leads directly to the *spectral decomposition* of  $\mathbf{D}$  given as

$$\mathbf{D} = \sum_{i=1}^r \delta_i \mathbf{l}_i \mathbf{r}_i^T$$

The best rank  $q$  approximation to the original data  $\mathbf{D}$  is the matrix  $\mathbf{D}_q = \sum_{i=1}^q \delta_i \mathbf{l}_i \mathbf{r}_i^T$  that minimizes the expression  $\|\mathbf{D} - \mathbf{D}_q\|_F$ , where  $\|\mathbf{A}\|_F = \sqrt{\sum_{i=1}^n \sum_{j=1}^d \mathbf{A}(i,j)^2}$  is called the *Frobenius Norm* of  $\mathbf{A}$ .

# Connection Between SVD and PCA

Assume  $\mathbf{D}$  has been centered, and let  $\mathbf{D} = \mathbf{L}\Delta\mathbf{R}^T$  via SVD. Consider the *scatter matrix* for  $\mathbf{D}$ , given as  $\mathbf{D}^T\mathbf{D}$ . We have

$$\mathbf{D}^T\mathbf{D} = (\mathbf{L}\Delta\mathbf{R}^T)^T(\mathbf{L}\Delta\mathbf{R}^T) = \mathbf{R}\Delta^T\mathbf{L}^T\mathbf{L}\Delta\mathbf{R}^T = \mathbf{R}(\Delta^T\Delta)\mathbf{R}^T = \mathbf{R}\Delta_d^2\mathbf{R}^T$$

where  $\Delta_d^2$  is the  $d \times d$  diagonal matrix defined as  $\Delta_d^2(i,i) = \delta_i^2$ , for  $i = 1, \dots, d$ .

The covariance matrix of centered  $\mathbf{D}$  is given as  $\Sigma = \frac{1}{n}\mathbf{D}^T\mathbf{D}$ ; we get

$$\begin{aligned}\mathbf{D}^T\mathbf{D} &= n\Sigma \\ &= n\mathbf{U}\Lambda\mathbf{U}^T \\ &= \mathbf{U}(n\Lambda)\mathbf{U}^T\end{aligned}$$

The right singular vectors  $\mathbf{R}$  are the same as the eigenvectors of  $\Sigma$ . The singular values of  $\mathbf{D}$  are related to the eigenvalues of  $\Sigma$  as

$$n\lambda_i = \delta_i^2, \text{ which implies } \lambda_i = \frac{\delta_i^2}{n}, \text{ for } i = 1, \dots, d$$

Likewise the left singular vectors in  $\mathbf{L}$  are the eigenvectors of the matrix  $n \times n$  matrix  $\mathbf{D}\mathbf{D}^T$ , and the corresponding eigenvalues are given as  $\delta_i^2$ .

# Data Mining and Machine Learning: Fundamental Concepts and Algorithms

dataminingbook.info

Mohammed J. Zaki<sup>1</sup>    Wagner Meira Jr.<sup>2</sup>

<sup>1</sup>Department of Computer Science  
Rensselaer Polytechnic Institute, Troy, NY, USA

<sup>2</sup>Department of Computer Science  
Universidade Federal de Minas Gerais, Belo Horizonte, Brazil

## Chapter 7: Dimensionality Reduction