

Data mining and Machine learning

Part 5. Kernel Methods

Input and Feature Space

For mining and analysis, it is important to find a suitable data representation. For example, for complex data such as text, sequences, images, and so on, we must typically extract or construct a set of attributes or features, so that we can represent the data instances as multivariate vectors.

Given a data instance \mathbf{x} (e.g., a sequence), we need to find a mapping ϕ , so that $\phi(\mathbf{x})$ is the vector representation of \mathbf{x} .

Even when the input data is a numeric data matrix a nonlinear mapping ϕ may be used to discover nonlinear relationships.

The term *input space* refers to the data space for the input data \mathbf{x} and *feature space* refers to the space of mapped vectors $\phi(\mathbf{x})$.

Sequence-based Features

Consider a dataset of DNA sequences over the alphabet $\Sigma = \{A, C, G, T\}$.

One simple feature space is to represent each sequence in terms of the probability distribution over symbols in Σ . That is, given a sequence x with length $|x| = m$, the mapping into feature space is given as

$$\phi(x) = \{P(A), P(C), P(G), P(T)\}$$

where $P(s) = \frac{n_s}{m}$ is the probability of observing symbol $s \in \Sigma$, and n_s is the number of times s appears in sequence x .

For example, if $x = ACAGCAGTA$, with $m = |x| = 9$, since A occurs four times, C and G occur twice, and T occurs once, we have

$$\phi(x) = (4/9, 2/9, 2/9, 1/9) = (0.44, 0.22, 0.22, 0.11)$$

We can compute larger feature spaces by considering, for example, the probability distribution over all substrings or words of size up to k over the alphabet Σ .

Nonlinear Features

Consider the mapping ϕ that takes as input a vector $\mathbf{x} = (x_1, x_2)^T \in \mathbb{R}^2$ and maps it to a “quadratic” feature space via the nonlinear mapping

$$\phi(\mathbf{x}) = (x_1^2, x_2^2, \sqrt{2}x_1x_2)^T \in \mathbb{R}^3$$

For example, the point $\mathbf{x} = (5.9, 3)^T$ is mapped to the vector

$$\phi(\mathbf{x}) = (5.9^2, 3^2, \sqrt{2} \cdot 5.9 \cdot 3)^T = (34.81, 9, 25.03)^T$$

We can then apply well-known linear analysis methods in the feature space.

Kernel Method

Let \mathcal{I} denote the input space, which can comprise any arbitrary set of objects, and let $\mathcal{D} = \{\mathbf{x}_i\}_{i=1}^n \subset \mathcal{I}$ be a dataset comprising n objects in the input space. Let $\phi: \mathcal{I} \rightarrow \mathcal{F}$ be a mapping from the input space \mathcal{I} to the feature space \mathcal{F} .

Kernel methods avoid explicitly transforming each point \mathbf{x} in the input space into the mapped point $\phi(\mathbf{x})$ in the feature space. Instead, the input objects are represented via their pairwise similarity values comprising the $n \times n$ *kernel matrix*, defined as

$$K = \begin{pmatrix} K(\mathbf{x}_1, \mathbf{x}_1) & K(\mathbf{x}_1, \mathbf{x}_2) & \cdots & K(\mathbf{x}_1, \mathbf{x}_n) \\ K(\mathbf{x}_2, \mathbf{x}_1) & K(\mathbf{x}_2, \mathbf{x}_2) & \cdots & K(\mathbf{x}_2, \mathbf{x}_n) \\ \vdots & \vdots & \ddots & \vdots \\ K(\mathbf{x}_n, \mathbf{x}_1) & K(\mathbf{x}_n, \mathbf{x}_2) & \cdots & K(\mathbf{x}_n, \mathbf{x}_n) \end{pmatrix}$$

$K: \mathcal{I} \times \mathcal{I} \rightarrow \mathbb{R}$ is a *kernel function* on any two points in input space, which should satisfy the condition

$$K(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$$

Intuitively, we need to be able to compute the value of the dot product using the original input representation \mathbf{x} , without having recourse to the mapping $\phi(\mathbf{x})$.

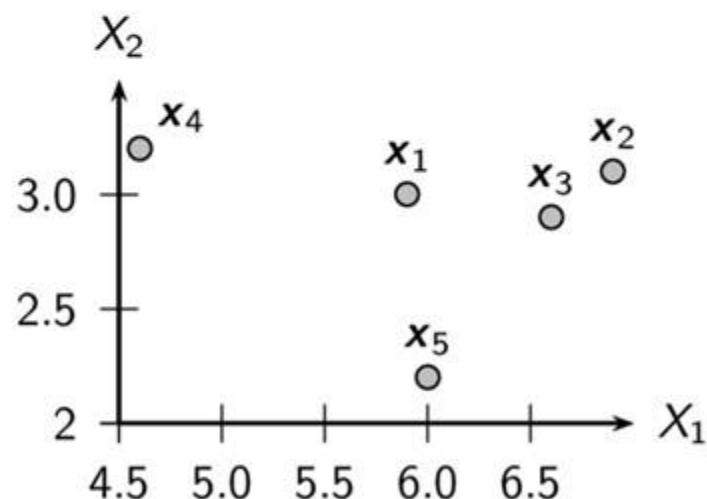
Linear Kernel

Let $\phi(\mathbf{x}) \rightarrow \mathbf{x}$ be the *identity kernel*. This leads to the *linear kernel*, which is simply the dot product between two input vectors:

$$\phi(\mathbf{x})^T \phi(\mathbf{y}) = \mathbf{x}^T \mathbf{y} = K(\mathbf{x}, \mathbf{y})$$

For example, if $\mathbf{x}_1 = (5.9 \quad 3)^T$ and $\mathbf{x}_2 = (6.9 \quad 3.1)^T$, then we have

$$K(\mathbf{x}_1, \mathbf{x}_2) = \mathbf{x}_1^T \mathbf{x}_2 = 5.9 \times 6.9 + 3 \times 3.1 = 40.71 + 9.3 = 50.01$$



K	x_1	x_2	x_3	x_4	x_5
x_1	43.81	50.01	47.64	36.74	42.00
x_2	50.01	57.22	54.53	41.66	48.22
x_3	47.64	54.53	51.97	39.64	45.98
x_4	36.74	41.66	39.64	31.40	34.64
x_5	42.00	48.22	45.98	34.64	40.84

Kernel Trick

Many data mining methods can be *kernelized* that is, instead of mapping the input points into feature space, the data can be represented via the $n \times n$ kernel matrix \mathbf{K} , and all relevant analysis can be performed over \mathbf{K} .

This is done via the *kernel trick*, that is, show that the analysis task requires only dot products $\phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$ in feature space, which can be replaced by the corresponding kernel $K(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$ that can be computed efficiently in input space.

Once the kernel matrix has been computed, we no longer even need the input points \mathbf{x}_i , as all operations involving only dot products in the feature space can be performed over the $n \times n$ kernel matrix \mathbf{K} .

Kernel Matrix

A function K is called a **positive semidefinite kernel** if and only if it is symmetric:

$$K(\mathbf{x}_i, \mathbf{x}_j) = K(\mathbf{x}_j, \mathbf{x}_i)$$

and the corresponding kernel matrix \mathbf{K} for any subset $D \subset \mathcal{I}$ is positive semidefinite, that is,

$$\mathbf{a}^T \mathbf{K} \mathbf{a} \geq 0, \text{ for all vectors } \mathbf{a} \in \mathbb{R}^n$$

which implies that

$$\sum_{i=1}^n \sum_{j=1}^n a_i a_j K(\mathbf{x}_i, \mathbf{x}_j) \geq 0, \text{ for all } a_i \in \mathbb{R}, i \in [1, n]$$

Dot Products and Positive Semi-definite Kernels

Positive Semidefinite Kernel

If $K(x_i, x_j)$ represents the dot product $\phi(x_i)^T \phi(x_j)$ in some feature space, then K is a positive semidefinite kernel.

First, K is symmetric since the dot product is symmetric, which also implies that K is symmetric.

Second, K is positive semidefinite because

$$\begin{aligned} \mathbf{a}^T \mathbf{K} \mathbf{a} &= \sum_{i=1}^n \sum_{j=1}^n a_i a_j K(x_i, x_j) \\ &= \sum_{i=1}^n \sum_{j=1}^n a_i a_j \phi(x_i)^T \phi(x_j) \\ &= \left(\sum_{i=1}^n a_i \phi(x_i) \right)^T \left(\sum_{j=1}^n a_j \phi(x_j) \right) \\ &= \left\| \sum_{i=1}^n a_i \phi(x_i) \right\|^2 \geq 0 \end{aligned}$$

Empirical Kernel Map

We now show that if we are given a positive semidefinite kernel $K: \mathcal{I} \times \mathcal{I} \rightarrow \mathbb{R}$, then it corresponds to a dot product in some feature space \mathcal{F} .

Define the map ϕ as follows:

$$\phi(\mathbf{x}) = \left((K(\mathbf{x}_1, \mathbf{x}), K(\mathbf{x}_2, \mathbf{x}), \dots, K(\mathbf{x}_n, \mathbf{x})) \right)^T \in \mathbb{R}^n$$

The *empirical kernel map* is defined as

$$\phi(\mathbf{x}) = \mathbf{K}^{-1/2} \cdot \left((K(\mathbf{x}_1, \mathbf{x}), K(\mathbf{x}_2, \mathbf{x}), \dots, K(\mathbf{x}_n, \mathbf{x})) \right)^T \in \mathbb{R}^n$$

so that the dot product yields

$$\begin{aligned}\phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j) &= \left(\mathbf{K}^{-1/2} \mathbf{K}_i \right)^T \left(\mathbf{K}^{-1/2} \mathbf{K}_j \right) \\ &= \mathbf{K}_i^T (\mathbf{K}^{-1/2} \mathbf{K}^{-1/2}) \mathbf{K}_j \\ &= \mathbf{K}_i^T \mathbf{K}^{-1} \mathbf{K}_j\end{aligned}$$

where \mathbf{K}_i is the i th column of \mathbf{K} . Over all pairs of mapped points, we have

$$\{\mathbf{K}_i^T \mathbf{K}^{-1} \mathbf{K}_j\}_{i,j=1}^n = \mathbf{K} \mathbf{K}^{-1} \mathbf{K} = \mathbf{K}$$

Data-specific Mercer Kernel Map

The Mercer kernel map also corresponds to a dot product in feature space.

Since K is a symmetric positive semidefinite matrix, it has real and non-negative eigenvalues. It can be decomposed as follows:

$$K = U \Lambda U^T$$

where U is the orthonormal matrix of eigenvectors $u_i = (u_{i1}, u_{i2}, \dots, u_{in})^T \in \mathbb{R}^n$ (for $i = 1, \dots, n$), and Λ is the diagonal matrix of eigenvalues, with both arranged in non-increasing order of the eigenvalues $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n \geq 0$:

The Mercer map ϕ is given as

$$\phi(x_i) = \sqrt{\Lambda} U_i$$

where U_i is the i th row of U .

The kernel value is simply the dot product between scaled rows of U :

$$\phi(x_i)^T \phi(x_j) = (\sqrt{\Lambda} U_i)^T (\sqrt{\Lambda} U_j) = U_i^T \Lambda U_j$$

Polynomial Kernel

Polynomial kernels are of two types: homogeneous or inhomogeneous.

Let $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$. The (inhomogeneous) *polynomial kernel* is defined as

$$K_q(\mathbf{x}, \mathbf{y}) = \phi(\mathbf{x})^T \phi(\mathbf{y}) = (c + \mathbf{x}^T \mathbf{y})^q$$

where q is the degree of the polynomial, and $c \geq 0$ is some constant. When $c = 0$ we obtain the homogeneous kernel, comprising only degree q terms. When $c > 0$, the feature space is spanned by all products of at most q attributes. This can be seen from the binomial expansion

$$K_q(\mathbf{x}, \mathbf{y}) = (c + \mathbf{x}^T \mathbf{y})^q = \sum_{k=1}^q \binom{q}{k} c^{q-k} (\mathbf{x}^T \mathbf{y})^k$$

The most typical cases are the *linear* (with $q = 1$) and *quadratic* (with $q = 2$) kernels, given as

$$K_1(\mathbf{x}, \mathbf{y}) = c + \mathbf{x}^T \mathbf{y}$$

$$K_2(\mathbf{x}, \mathbf{y}) = (c + \mathbf{x}^T \mathbf{y})^2$$

Gaussian Kernel

The Gaussian kernel, also called the Gaussian radial basis function (RBF) kernel, is defined as

$$K(\mathbf{x}, \mathbf{y}) = \exp \left\{ -\frac{\|\mathbf{x} - \mathbf{y}\|^2}{2\sigma^2} \right\}$$

where $\sigma > 0$ is the spread parameter that plays the same role as the standard deviation in a normal density function.

Note that $K(\mathbf{x}, \mathbf{x}) = 1$, and further that the kernel value is inversely related to the distance between the two points \mathbf{x} and \mathbf{y} .

A feature space for the Gaussian kernel has infinite dimensionality.

Basic Kernel Operations in Feature Space

Basic data analysis tasks that can be performed solely via kernels, without instantiating $\phi(\mathbf{x})$.

Norm of a Point: We can compute the norm of a point $\phi(\mathbf{x})$ in feature space as follows:

$$\|\phi(\mathbf{x})\|^2 = \phi(\mathbf{x})^T \phi(\mathbf{x}) = K(\mathbf{x}, \mathbf{x})$$

which implies that $\|\phi(\mathbf{x})\| = \sqrt{K(\mathbf{x}, \mathbf{x})}$.

Distance between Points: The distance between $\phi(\mathbf{x}_i)$ and $\phi(\mathbf{x}_j)$ is

$$\begin{aligned}\|\phi(\mathbf{x}_i) - \phi(\mathbf{x}_j)\|^2 &= \|\phi(\mathbf{x}_i)\|^2 + \|\phi(\mathbf{x}_j)\|^2 - 2\phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j) \\ &= K(\mathbf{x}_i, \mathbf{x}_i) + K(\mathbf{x}_j, \mathbf{x}_j) - 2K(\mathbf{x}_i, \mathbf{x}_j)\end{aligned}$$

which implies that

$$\|\phi(\mathbf{x}_i) - \phi(\mathbf{x}_j)\| = \sqrt{K(\mathbf{x}_i, \mathbf{x}_i) + K(\mathbf{x}_j, \mathbf{x}_j) - 2K(\mathbf{x}_i, \mathbf{x}_j)}$$

Basic Kernel Operations in Feature Space

Kernel Value as Similarity: We can rearrange the terms in

$$\|\phi(\mathbf{x}_i) - \phi(\mathbf{x}_j)\|^2 = K(\mathbf{x}_i, \mathbf{x}_i) + K(\mathbf{x}_j, \mathbf{x}_j) - 2K(\mathbf{x}_i, \mathbf{x}_j)$$

to obtain

$$\frac{1}{2} (\|\phi(\mathbf{x}_i)\|^2 + \|\phi(\mathbf{x}_j)\|^2 - \|\phi(\mathbf{x}_i) - \phi(\mathbf{x}_j)\|^2) = K(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$$

The more the distance $\|\phi(\mathbf{x}_i) - \phi(\mathbf{x}_j)\|$ between the two points in feature space, the less the kernel value, that is, the less the similarity.

Mean in Feature Space: The mean of the points in feature space is given as $\mu_\phi = 1/n \sum_{i=1}^n \phi(\mathbf{x}_i)$. Thus, we cannot compute it explicitly. However, the squared norm of the mean is:

$$\|\mu_\phi\|^2 = \mu_\phi^T \mu_\phi = \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n K(\mathbf{x}_i, \mathbf{x}_j)$$

The squared norm of the mean in feature space is simply the average of the values in the kernel matrix \mathbf{K} .

Basic Kernel Operations in Feature Space

Total Variance in Feature Space: The total variance in feature space is obtained by taking the average squared deviation of points from the mean in feature space:

$$\sigma_{\phi}^2 = \frac{1}{n} \sum_{i=1}^n \|\phi(\mathbf{x}_i) - \boldsymbol{\mu}_{\phi}\|^2 = \frac{1}{n} \sum_{i=1}^n K(\mathbf{x}_i, \mathbf{x}_i) - \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n K(\mathbf{x}_i, \mathbf{x}_j)$$

Centering in Feature Space We can center each point in feature space by subtracting the mean from it, as follows:

$$\hat{\phi}(\mathbf{x}_i) = \phi(\mathbf{x}_i) - \boldsymbol{\mu}_{\phi}$$

The kernel between centered points is given as

$$\begin{aligned}\hat{K}(\mathbf{x}_i, \mathbf{x}_j) &= \hat{\phi}(\mathbf{x}_i)^T \hat{\phi}(\mathbf{x}_j) \\ &= K(\mathbf{x}_i, \mathbf{x}_j) - \frac{1}{n} \sum_{k=1}^n K(\mathbf{x}_i, \mathbf{x}_k) - \frac{1}{n} \sum_{k=1}^n K(\mathbf{x}_j, \mathbf{x}_k) + \frac{1}{n^2} \sum_{a=1}^n \sum_{b=1}^n K(\mathbf{x}_a, \mathbf{x}_b)\end{aligned}$$

More compactly, we have:

$$\hat{K} = \left(\mathbf{I} - \frac{1}{n} \mathbf{1}_{n \times n} \right) K \left(\mathbf{I} - \frac{1}{n} \mathbf{1}_{n \times n} \right)$$

where $\mathbf{1}_{n \times n}$ is the $n \times n$ matrix of ones.

Basic Kernel Operations in Feature Space

Normalizing in Feature Space: The dot product between normalized points in feature space corresponds to the cosine of the angle between them

$$\phi_n(\mathbf{x}_i)^T \phi_n(\mathbf{x}_j) = \frac{\phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)}{\|\phi(\mathbf{x}_i)\| \cdot \|\phi(\mathbf{x}_j)\|} = \cos \theta$$

If the mapped points are both centered and normalized, then a dot product corresponds to the correlation between the two points in feature space.

The normalized kernel matrix, \mathbf{K}_n , can be computed using only the kernel function K , as

$$\mathbf{K}_n(\mathbf{x}_i, \mathbf{x}_j) = \frac{\phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)}{\|\phi(\mathbf{x}_i)\| \cdot \|\phi(\mathbf{x}_j)\|} = \frac{K(\mathbf{x}_i, \mathbf{x}_j)}{\sqrt{K(\mathbf{x}_i, \mathbf{x}_i) \cdot K(\mathbf{x}_j, \mathbf{x}_j)}}$$

\mathbf{K}_n has all diagonal elements as 1.

Spectrum Kernel for Strings

Given alphabet Σ , the l -spectrum feature map is the mapping $\phi: \Sigma^* \rightarrow \mathbb{R}^{|\Sigma|^l}$ from the set of substrings over Σ to the $|\Sigma|^l$ -dimensional space representing the number of occurrences of all possible substrings of length l , defined as

$$\phi(x) = \left(\dots, \#(\alpha), \dots \right)_{\alpha \in \Sigma^l}^T$$

where $\#(\alpha)$ is the number of occurrences of the l -length string α in x .

The (full) spectrum map considers all lengths from $l=0$ to $l=\infty$, leading to an infinite dimensional feature map $\phi: \Sigma^* \rightarrow \mathbb{R}^\infty$:

$$\phi(x) = \left(\dots, \#(\alpha), \dots \right)_{\alpha \in \Sigma^*}^T$$

where $\#(\alpha)$ is the number of occurrences of the string α in x .

The (l -)spectrum kernel between two strings x_i, x_j is simply the dot product between their (l -)spectrum maps:

$$K(x_i, x_j) = \phi(x_i)^T \phi(x_j)$$

The (full) spectrum kernel can be computed efficiently via suffix trees in $O(n+m)$ time for two strings of length n and m .

Diffusion Kernels on Graph Nodes

Let \mathbf{S} be some symmetric similarity matrix between nodes of a graph $G = (V, E)$. For instance, \mathbf{S} can be the (weighted) adjacency matrix \mathbf{A} or the Laplacian matrix $\mathbf{L} = \mathbf{A} - \Delta$ (or its negation), where Δ is the degree matrix for an undirected graph G , defined as $\Delta(i, i) = d_i$ and $\Delta(i, j) = 0$ for all $i \neq j$, and d_i is the degree of node i .

Power Kernels: Summing up the product of the base similarities over all l -length paths between two nodes, we obtain the l -length similarity matrix $\mathbf{S}^{(l)}$, which is simply the l th power of \mathbf{S} , that is,

$$\mathbf{S}^{(l)} = \mathbf{S}^l$$

Even path lengths lead to positive semidefinite kernels, but odd path lengths are not guaranteed to do so, unless the base matrix \mathbf{S} is itself a positive semidefinite matrix.

Power kernel \mathbf{K} can be obtained via the eigen-decomposition of \mathbf{S}^l :

$$\mathbf{K} = \mathbf{S}^l = (\mathbf{U} \Lambda \mathbf{U}^T)^l = \mathbf{U} (\Lambda^l) \mathbf{U}^T$$

Exponential Diffusion Kernel

The exponential diffusion kernel we can obtain a new kernel between nodes of a graph by paths of all possible lengths, but damps the contribution of longer paths

$$\begin{aligned} K &= \sum_{l=0}^{\infty} \frac{1}{l!} \beta^l S^l \\ &= I + \beta S + \frac{1}{2!} \beta^2 S^2 + \frac{1}{3!} \beta^3 S^3 + \dots \\ &= \exp\{\beta S\} \end{aligned}$$

where β is a damping factor, and $\exp\{\beta S\}$ is the matrix exponential. The series on the right hand side above converges for all $\beta \geq 0$.

Substituting $S = U \Lambda U^T$ the kernel can be computed as

$$\begin{aligned} K &= I + \beta S + \frac{1}{2!} \beta^2 S^2 + \dots \\ &= U \begin{pmatrix} \exp\{\beta \lambda_1\} & 0 & \cdots & 0 \\ 0 & \exp\{\beta \lambda_2\} & \cdots & 0 \\ \vdots & \vdots & \ddots & 0 \\ 0 & 0 & \cdots & \exp\{\beta \lambda_n\} \end{pmatrix} U^T \end{aligned}$$

where λ_i is an eigenvalue of S .

Von Neumann Diffusion Kernel

The *von Neumann diffusion kernel* is defined as

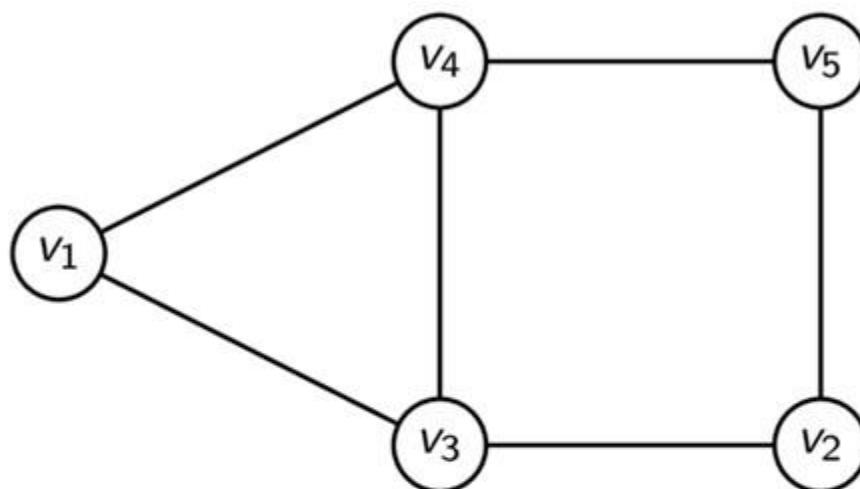
$$\boldsymbol{\mathcal{K}} = \sum_{l=0}^{\infty} \beta^l \boldsymbol{S}^l$$

where $\beta \geq 0$. Expanding and rearranging the terms, we obtain

$$\boldsymbol{\mathcal{K}} = (\boldsymbol{I} - \beta \boldsymbol{S})^{-1}$$

The kernel is guaranteed to be positive semidefinite if $|\beta| < 1/\rho(\boldsymbol{S})$, where $\rho(\boldsymbol{S}) = \max_i\{|\lambda_i|\}$ is called the *spectral radius* of \boldsymbol{S} , defined as the largest eigenvalue of \boldsymbol{S} in absolute value.

Graph Diffusion Kernel: Example



Adjacency and degree matrices are given as

$$\mathbf{A} = \begin{pmatrix} 0 & 0 & 1 & 1 & 0 \\ 0 & 0 & 1 & 0 & 1 \\ 1 & 1 & 0 & 1 & 0 \\ 1 & 0 & 1 & 0 & 1 \\ 0 & 1 & 0 & 1 & 0 \end{pmatrix} \quad \Delta = \begin{pmatrix} 2 & 0 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 & 0 \\ 0 & 0 & 3 & 0 & 0 \\ 0 & 0 & 0 & 3 & 0 \\ 0 & 0 & 0 & 0 & 2 \end{pmatrix}$$

Graph Diffusion Kernel: Example

Let the base similarity matrix \mathbf{S} be the negated Laplacian matrix

$$\mathbf{S} = -\mathbf{L} = \mathbf{A} - \mathbf{D} = \begin{pmatrix} -2 & 0 & 1 & 1 & 0 \\ 0 & -2 & 1 & 0 & 1 \\ 1 & 1 & -3 & 1 & 0 \\ 1 & 0 & 1 & -3 & 1 \\ 0 & 1 & 0 & 1 & -2 \end{pmatrix}$$

The eigenvalues of \mathbf{S} are as follows:

$$\lambda_1 = 0 \quad \lambda_2 = -1.38 \quad \lambda_3 = -2.38 \quad \lambda_4 = -3.62 \quad \lambda_5 = -4.62$$

and the eigenvectors of \mathbf{S} are

$$\mathbf{U} = \begin{pmatrix} \mathbf{u}_1 & \mathbf{u}_2 & \mathbf{u}_3 & \mathbf{u}_4 & \mathbf{u}_5 \end{pmatrix} = \begin{pmatrix} 0.45 & -0.63 & 0.00 & 0.63 & 0.00 \\ 0.45 & 0.51 & -0.60 & 0.20 & -0.37 \\ 0.45 & -0.20 & -0.37 & -0.51 & 0.60 \\ 0.45 & -0.20 & 0.37 & -0.51 & -0.60 \\ 0.45 & 0.51 & 0.60 & 0.20 & 0.37 \end{pmatrix}$$

Graph Diffusion Kernel: Example

Assuming $\beta = 0.2$, the exponential diffusion kernel matrix is given as

$$\begin{aligned}\boldsymbol{\kappa} &= \exp\{0.2\mathbf{S}\} = \mathbf{U} \begin{pmatrix} \exp\{0.2\lambda_1\} & 0 & \cdots & 0 \\ 0 & \exp\{0.2\lambda_2\} & \cdots & 0 \\ \vdots & \vdots & \ddots & 0 \\ 0 & 0 & \cdots & \exp\{0.2\lambda_n\} \end{pmatrix} \mathbf{U}^T \\ &= \begin{pmatrix} 0.70 & 0.01 & 0.14 & 0.14 & 0.01 \\ 0.01 & 0.70 & 0.13 & 0.03 & 0.14 \\ 0.14 & 0.13 & 0.59 & 0.13 & 0.03 \\ 0.14 & 0.03 & 0.13 & 0.59 & 0.13 \\ 0.01 & 0.14 & 0.03 & 0.13 & 0.70 \end{pmatrix}\end{aligned}$$

Assuming $\beta = 0.2$, the von Neumann kernel is given as

$$\boldsymbol{\kappa} = \mathbf{U}(\mathbf{I} - 0.2\Lambda)^{-1} \mathbf{U}^T = \begin{pmatrix} 0.75 & 0.02 & 0.11 & 0.11 & 0.02 \\ 0.02 & 0.74 & 0.10 & 0.03 & 0.11 \\ 0.11 & 0.10 & 0.66 & 0.10 & 0.03 \\ 0.11 & 0.03 & 0.10 & 0.66 & 0.10 \\ 0.02 & 0.11 & 0.03 & 0.10 & 0.74 \end{pmatrix}$$

Data mining and Machine learning

Part 6. High-Dimensional Data

High-dimensional Space

Let D be a $n \times d$ data matrix. In data mining typically the data is very high dimensional. Understanding the nature of high-dimensional space, or *hyperspace*, is very important, especially because it does not behave like the more familiar geometry in two or three dimensions.

Hyper-rectangle: The data space is a d -dimensional *hyper-rectangle*

$$R_d = \prod_{j=1}^d [\min(X_j), \max(X_j)]$$

where $\min(X_j)$ and $\max(X_j)$ specify the range of X_j .

Hypercube: Assume the data is centered, and let m denote the maximum attribute value

$$m = \max_{j=1}^d \max_{i=1}^n \{ |x_{ij}| \}$$

The data hyperspace can be represented as a *hypercube*, centered at 0, with all sides of length $l = 2m$, given as

$$H_d(l) = \left\{ \mathbf{x} = (x_1, x_2, \dots, x_d)^T \mid \forall i, x_i \in [-l/2, l/2] \right\}$$

The *unit hypercube* has all sides of length $l = 1$, and is denoted as $H_d(1)$.

Hypersphere

Assume that the data has been centered, so that $\mu = 0$. Let r denote the largest magnitude among all points:

$$r = \max_i \{ \|x_i\| \}$$

The data hyperspace can be represented as a d -dimensional *hyperball* centered at 0 with radius r , defined as

$$B_d(r) = \{x \mid \|x\| \leq r\} \text{ or } B_d(r) = \left\{ x = (x_1, x_2, \dots, x_d) \mid \sum_{j=1}^d x_j^2 \leq r^2 \right\}$$

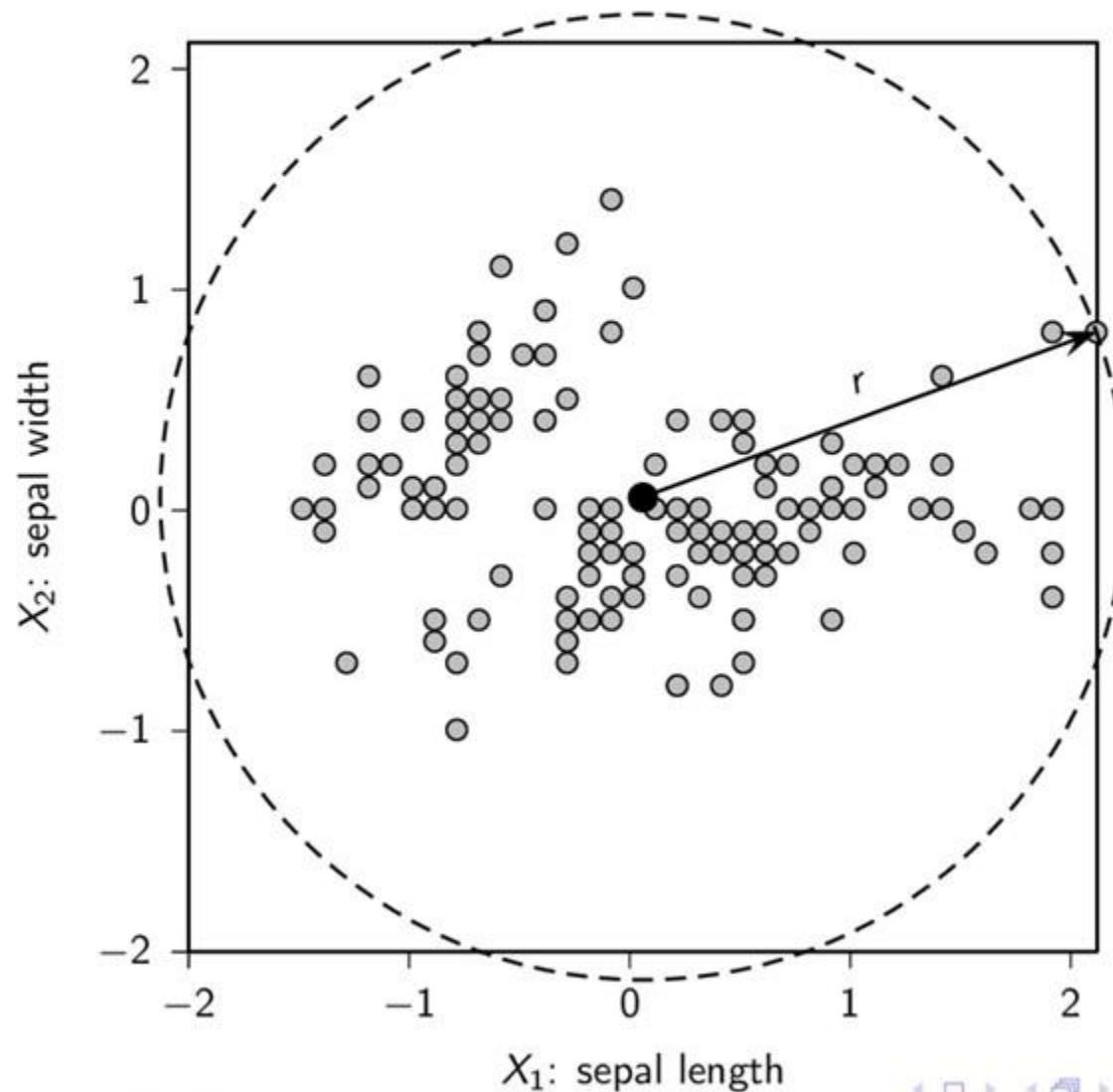
The surface of the hyperball is called a *hypersphere*, and it consists of all the points exactly at distance r from the center of the hyperball

$$S_d(r) = \{x \mid \|x\| = r\}$$

$$\text{or } S_d(r) = \left\{ x = (x_1, x_2, \dots, x_d) \mid \sum_{j=1}^d (x_j)^2 = r^2 \right\}$$

Iris Data Hyperspace: Hypercube and Hypersphere

$l = 4.12$ and $r = 2.19$



High-dimensional Volumes

Hypercube: The volume of a hypercube with edge length l is given as

$$\text{vol}(H_d(l)) = l^d$$

Hypersphere The volume of a hyperball and its corresponding hypersphere is identical
The volume of a hypersphere is given as

$$\begin{aligned} \text{In 1D: } \text{vol}(S_1(r)) &= 2r & \text{In 2D: } \text{vol}(S_2(r)) &= \pi r^2 & \text{In 3D: } \text{vol}(S_3(r)) &= \frac{4}{3}\pi r^3 \end{aligned}$$

In d -dimensions:

$$\text{vol}(S_d(r)) = K_d r^d = \left(\frac{\pi^{\frac{d}{2}}}{\Gamma(\frac{d}{2} + 1)} \right) r^d$$

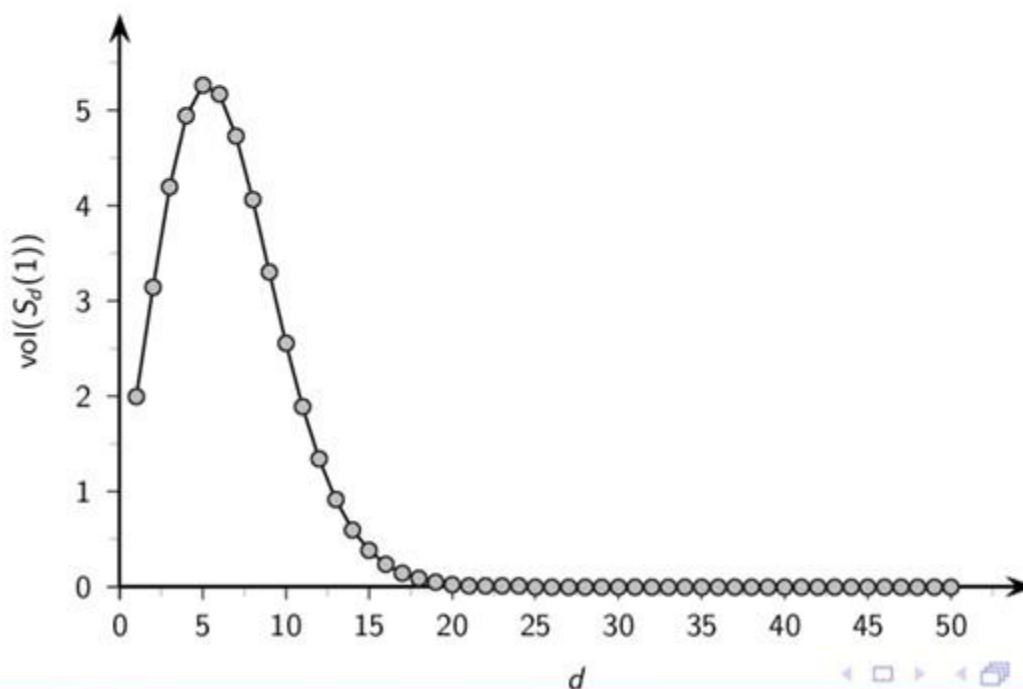
where

$$\Gamma\left(\frac{d}{2} + 1\right) = \begin{cases} \left(\frac{d}{2}\right)! & \text{if } d \text{ is even} \\ \sqrt{\pi} \left(\frac{d!!}{2^{(d+1)/2}}\right) & \text{if } d \text{ is odd} \end{cases}$$

Volume of Unit Hypersphere

With increasing dimensionality the hypersphere volume first increases up to a point, and then starts to decrease, and ultimately vanishes. In particular, for the unit hypersphere with $r = 1$,

$$\lim_{d \rightarrow \infty} \text{vol}(S_d(1)) = \lim_{d \rightarrow \infty} \frac{\pi^{\frac{d}{2}}}{\Gamma(\frac{d}{2} + 1)} \rightarrow 0$$



Hypersphere Inscribed within Hypercube

Consider the space enclosed within the largest hypersphere that can be accommodated within a hypercube (which represents the dataspace).

The ratio of the volume of the hypersphere of radius r to the hypercube with side length $l = 2r$ is given as

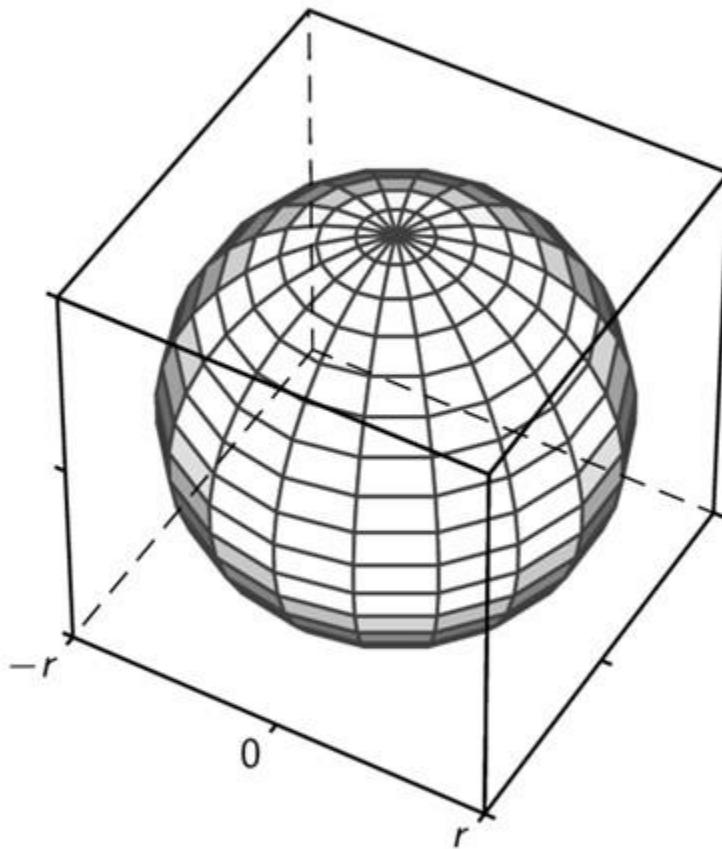
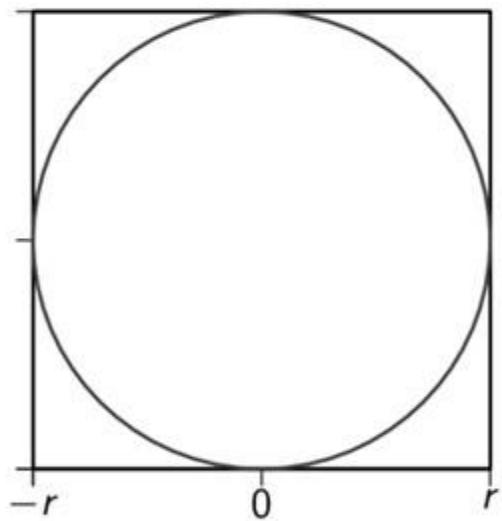
$$\text{In 2 dimensions: } \frac{\text{vol}(S_2(r))}{\text{vol}(H_2(2r))} = \frac{\pi r^2}{4r^2} = \frac{\pi}{4} = 78.5\%$$

$$\text{In 3 dimensions: } \frac{\text{vol}(S_3(r))}{\text{vol}(H_3(2r))} = \frac{\frac{4}{3}\pi r^3}{8r^3} = \frac{\pi}{6} = 52.4\%$$

$$\text{In } d \text{ dimensions: } \lim_{d \rightarrow \infty} \frac{\text{vol}(S_d(r))}{\text{vol}(H_d(2r))} = \lim_{d \rightarrow \infty} \frac{\pi^{d/2}}{2^d \Gamma(\frac{d}{2} + 1)} \rightarrow 0$$

As the dimensionality increases, most of the volume of the hypercube is in the "corners," whereas the center is essentially empty.

Hypersphere Inscribed inside a Hypercube

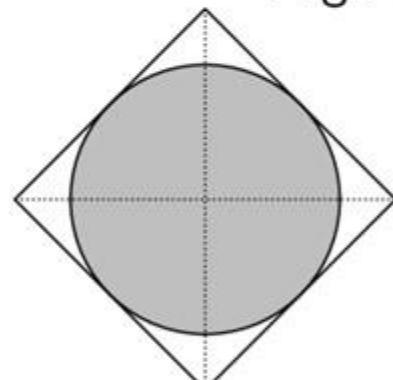


Conceptual View of High-dimensional Space

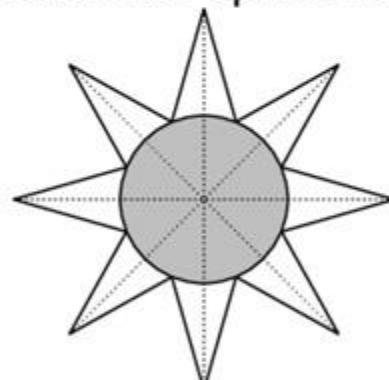
Two, three, four, and higher dimensions

All the volume of the hyperspace is in the corners, with the center being essentially empty.

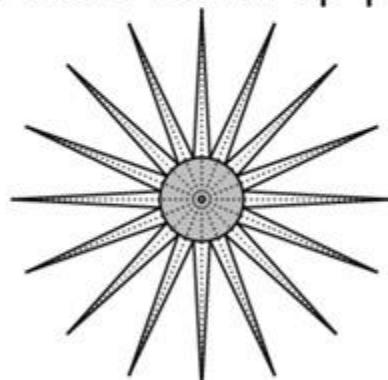
High-dimensional space looks like a rolled-up porcupine!



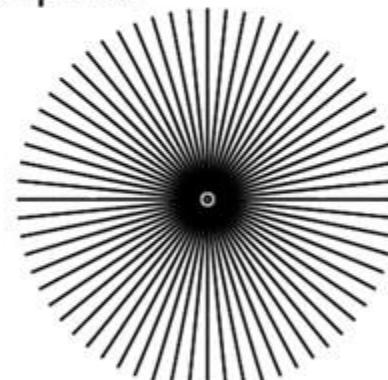
(a) 2D



(b) 3D



(c) 4D



(d) dD

Volume of a Thin Shell

The volume of a thin hypershell of width ϵ is given as

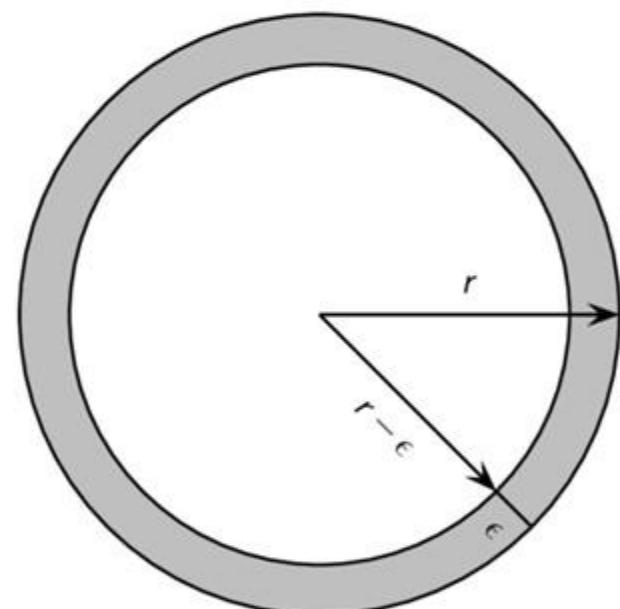
$$\begin{aligned}\text{vol}(S_d(r, \epsilon)) &= \text{vol}(S_d(r)) - \text{vol}(S_d(r - \epsilon)) \\ &= K_d r^d - K_d (r - \epsilon)^d.\end{aligned}$$

The ratio of volume of the thin shell to the volume of the outer sphere:

$$\frac{\text{vol}(S_d(r, \epsilon))}{\text{vol}(S_d(r))} = \frac{K_d r^d - K_d (r - \epsilon)^d}{K_d r^d} = 1 - \left(1 - \frac{\epsilon}{r}\right)^d$$

As d increases, we have

$$\lim_{d \rightarrow \infty} \frac{\text{vol}(S_d(r, \epsilon))}{\text{vol}(S_d(r))} = \lim_{d \rightarrow \infty} 1 - \left(1 - \frac{\epsilon}{r}\right)^d \rightarrow 1$$



Diagonals in Hyperspace

Consider a d -dimensional hypercube, with origin $0_d = (0_1, 0_2, \dots, 0_d)$, and bounded in each dimension in the range $[-1, 1]$. Each “corner” of the hyperspace is a d -dimensional vector of the form $(\pm 1_1, \pm 1_2, \dots, \pm 1_d)^T$.

Let $\mathbf{e}_i = (0_1, \dots, 1_i, \dots, 0_d)^T$ denote the d -dimensional canonical unit vector in dimension i , and let 1 denote the d -dimensional diagonal vector $(1_1, 1_2, \dots, 1_d)^T$.

Consider the angle θ_d between the diagonal vector 1 and the first axis \mathbf{e}_1 , in d dimensions:

$$\cos \theta_d = \frac{\mathbf{e}_1^T 1}{\|\mathbf{e}_1\| \|1\|} = \frac{\mathbf{e}_1^T 1}{\sqrt{\mathbf{e}_1^T \mathbf{e}_1} \sqrt{1^T 1}} = \frac{1}{\sqrt{1} \sqrt{d}} = \frac{1}{\sqrt{d}}$$

Diagonals in Hyperspace

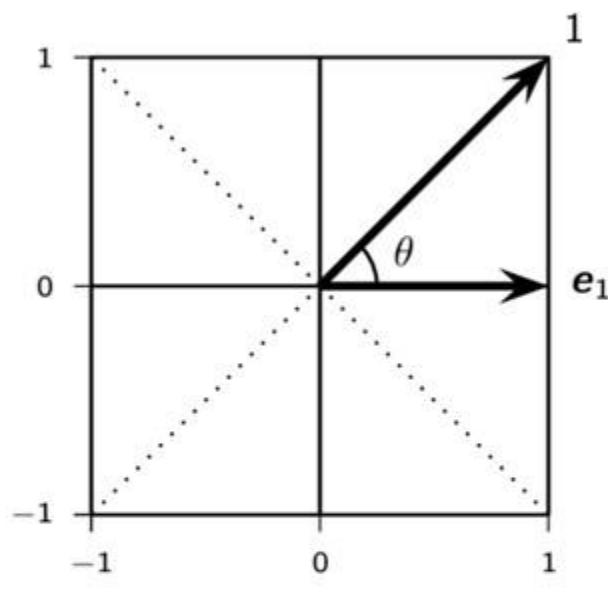
As d increases, we have

$$\lim_{d \rightarrow \infty} \cos \theta_d = \lim_{d \rightarrow \infty} \frac{1}{\sqrt{d}} \rightarrow 0$$

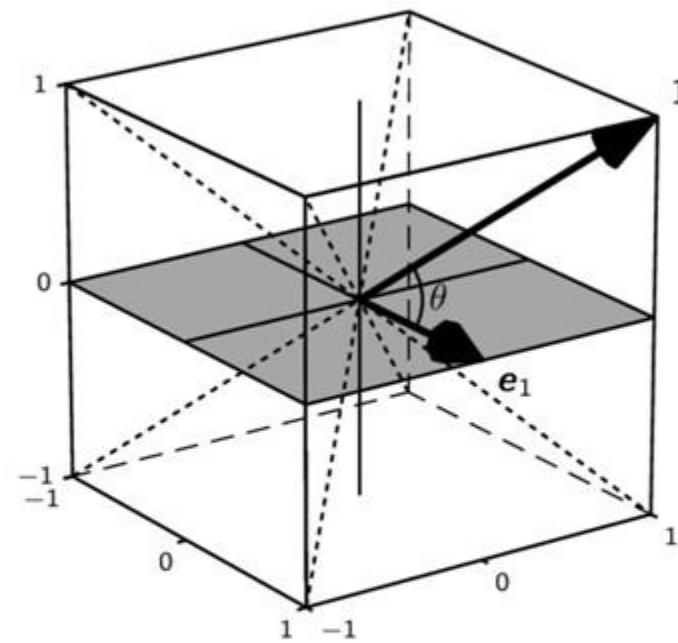
which implies that

$$\lim_{d \rightarrow \infty} \theta_d \rightarrow \frac{\pi}{2} = 90^\circ$$

Angle between Diagonal Vector 1 and e_1



(a) In 2D



(b) In 3D

In high dimensions all of the diagonal vectors are perpendicular (or orthogonal) to all the coordinates axes! Each of the 2^{d-1} new axes connecting pairs of 2^d corners are essentially orthogonal to all of the d principal coordinate axes! Thus, in effect, high-dimensional space has an exponential number of orthogonal “axes.”

Density of the Multivariate Normal

Consider the standard multivariate normal distribution with $\mu = 0$, and $\Sigma = I$

$$f(\mathbf{x}) = \frac{1}{(\sqrt{2\pi})^d} \exp\left\{-\frac{\mathbf{x}^T \mathbf{x}}{2}\right\}$$

The peak of the density is at the mean. Consider the set of points \mathbf{x} with density at least α fraction of the density at the mean

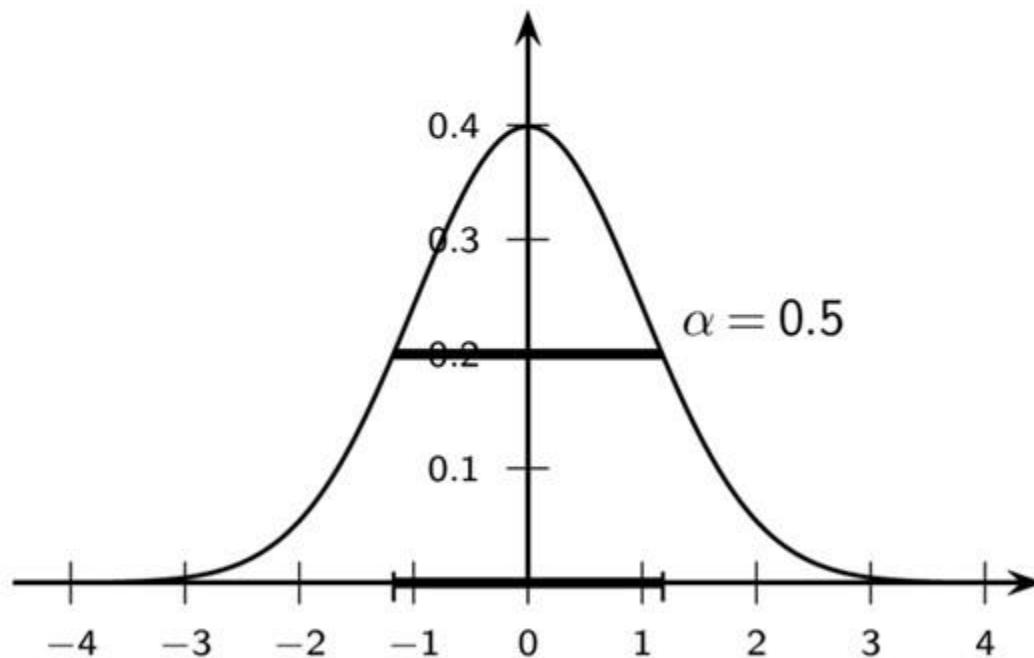
$$\begin{aligned}\frac{f(\mathbf{x})}{f(0)} &\geq \alpha \\ \exp\left\{-\frac{\mathbf{x}^T \mathbf{x}}{2}\right\} &\geq \alpha \\ \mathbf{x}^T \mathbf{x} &\leq -2 \ln(\alpha) \\ \sum_{i=1}^d (x_i)^2 &\leq -2 \ln(\alpha)\end{aligned}$$

The sum of squared IID random variables follows a chi-squared distribution χ_d^2 . Thus,

$$P\left(\frac{f(\mathbf{x})}{f(0)} \geq \alpha\right) = F_{\chi_d^2}(-2 \ln(\alpha))$$

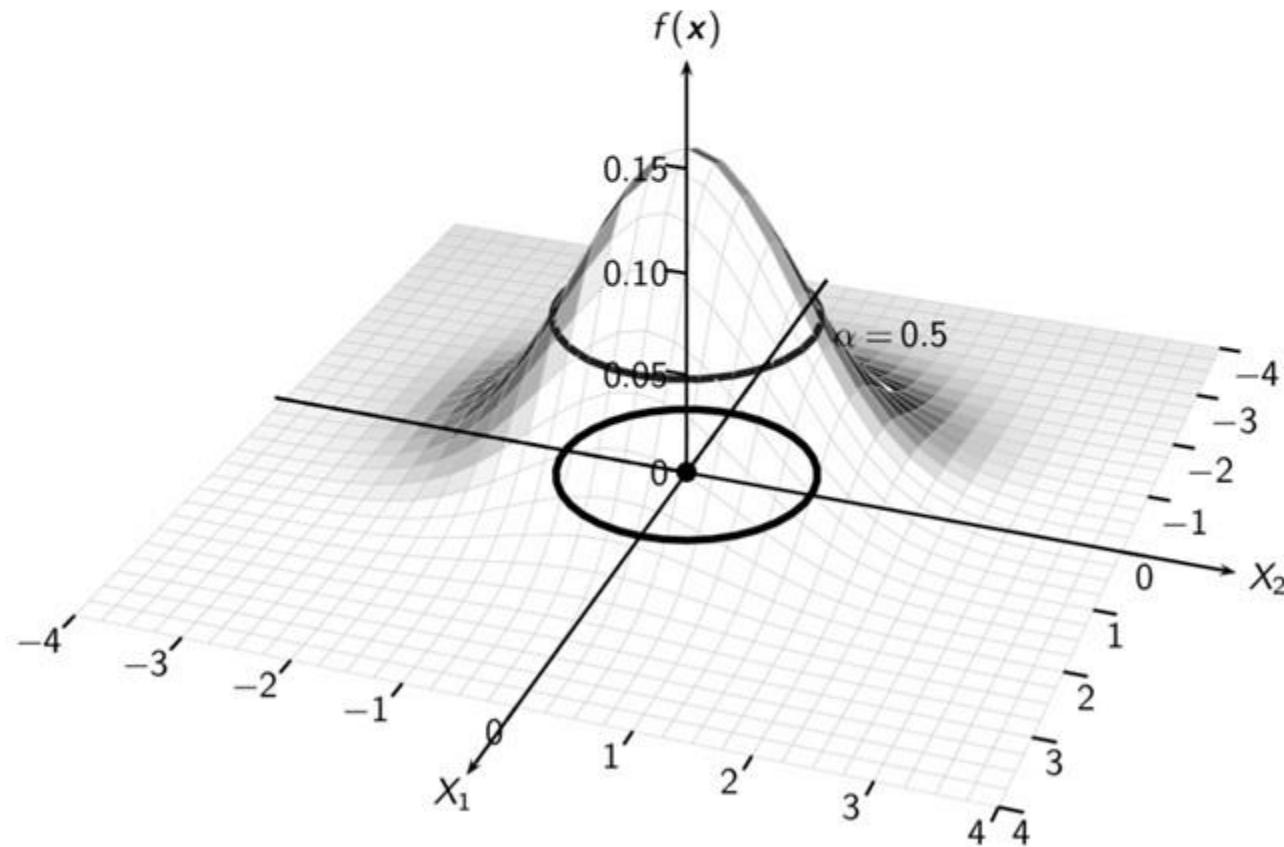
Density Contour for α Fraction of the Density at the Mean: One Dimension

Let $\alpha = 0.5$, then $-2\ln(0.5) = 1.386$ and $F_{\chi^2_1}(1.386) = 0.76$. Thus, 24% of the density is in the tail regions.



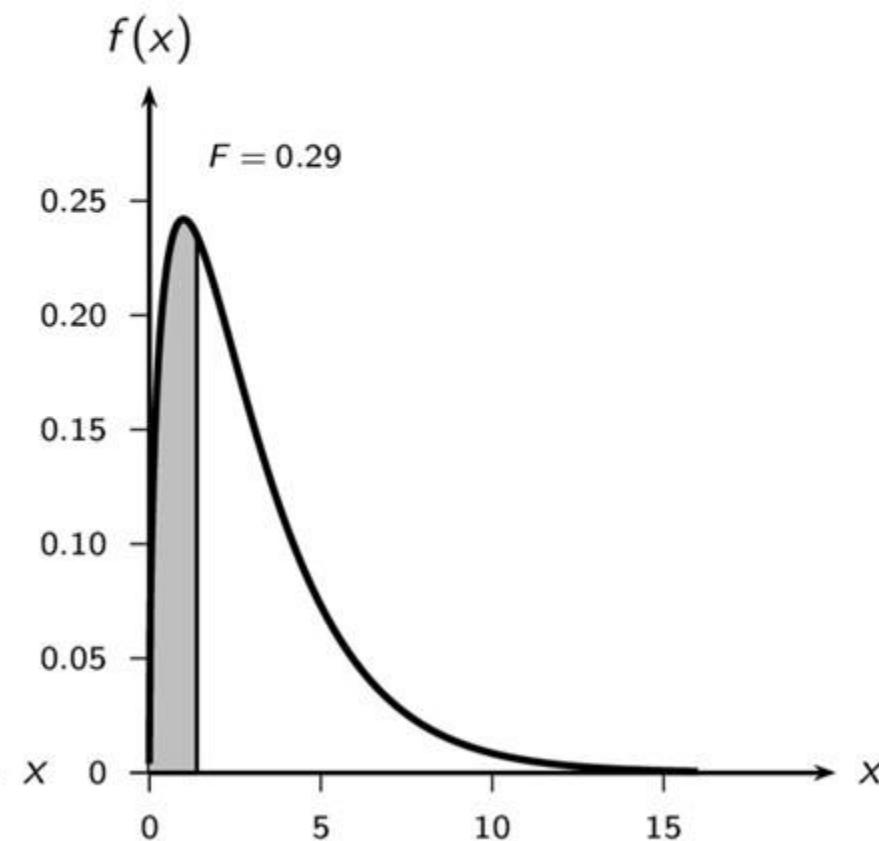
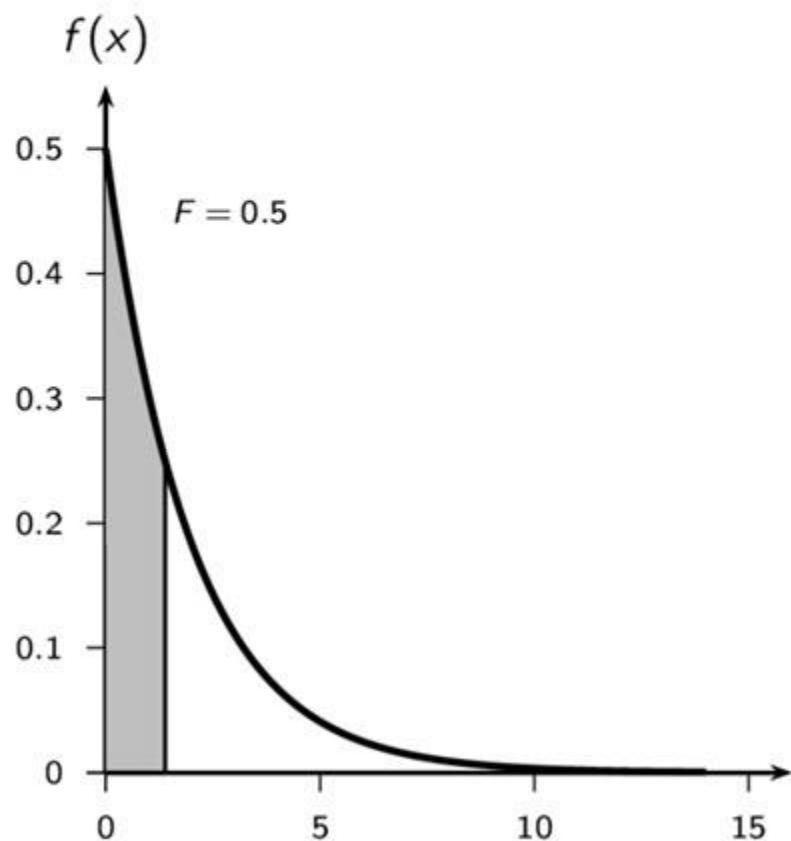
Density Contour for α Fraction of the Density at the Mean: Two Dimensions

Let $\alpha = 0.5$, then $-2\ln(0.5) = 1.386$ and $F_{\chi^2_2}(1.386) = 0.50$. Thus, 50% of the density is in the tail regions.

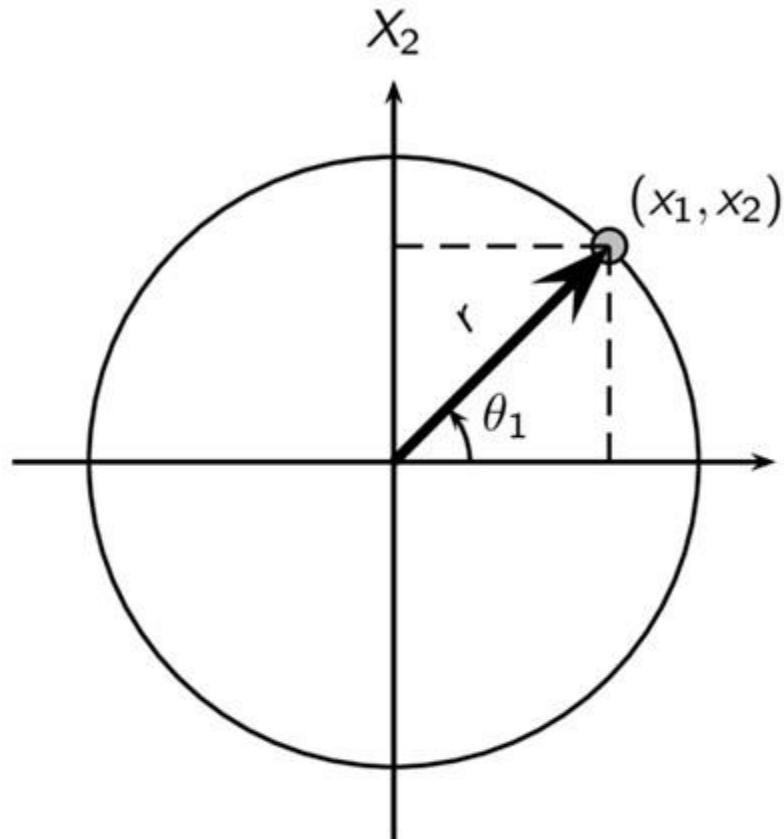


Chi-Squared Distribution: $P(f(x)/f(0) \geq \alpha)$

This probability decreases rapidly with dimensionality. For 2D, it is 0.5. For 3D it is 0.29, ie., 71% of the density is in the tails. By $d = 10$, it decreases to 0.075%, that is, 99.925% of the points lie in the extreme or tail regions.



Hypersphere Volume: Polar Coordinates in 2D



The point $x = (x_1, x_2)$ in polar coordinates

$$x_1 = r \cos \theta_1 = rc_1$$

$$x_2 = r \sin \theta_1 = rs_1$$

where $r = \|x\|$, and $\cos \theta_1 = c_1$ and $\sin \theta_1 = s_1$.

The *Jacobian matrix* for this transformation is given as

$$J(\theta_1) = \begin{pmatrix} \frac{\partial x_1}{\partial r} & \frac{\partial x_1}{\partial \theta_1} \\ \frac{\partial x_2}{\partial r} & \frac{\partial x_2}{\partial \theta_1} \end{pmatrix} = \begin{pmatrix} c_1 & -rs_1 \\ s_1 & rc_1 \end{pmatrix}$$

Hypersphere volume is obtained by integration over r and θ_1 (with $r > 0$, and $0 \leq \theta_1 \leq 2\pi$):

$$\begin{aligned} \text{vol}(S_2(r)) &= \int_r \int_{\theta_1} |\det(J(\theta_1))| dr d\theta_1 \\ &= \int_0^r \int_0^{2\pi} r dr d\theta_1 = \int_0^r r dr \int_0^{2\pi} d\theta_1 \\ &= \frac{r^2}{2} \left. \theta_1 \right|_0^{2\pi} = \pi r^2 \end{aligned}$$

Hypersphere Volume: Polar Coordinates in 3D

$x = (x_1, x_2, x_3)$ in polar coordinates

$$x_1 = r \cos \theta_1 \cos \theta_2 = r c_1 c_2$$

$$x_2 = r \cos \theta_1 \sin \theta_2 = r c_1 s_2$$

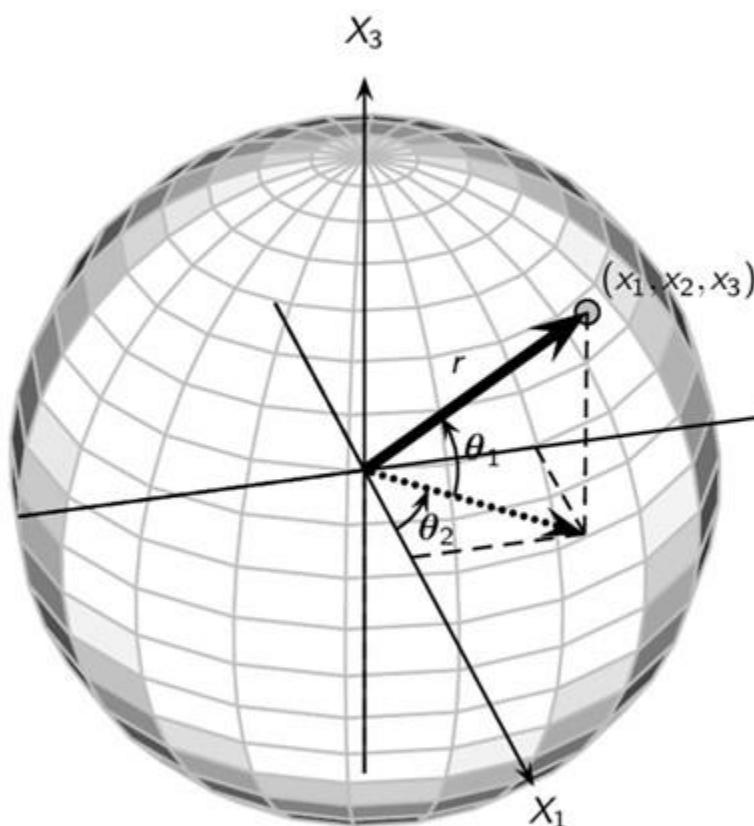
$$x_3 = r \sin \theta_1 = r s_1$$

The Jacobian matrix is given as

$$J(\theta_1, \theta_2) = \begin{pmatrix} c_1 c_2 & -r s_1 c_2 & -r c_1 s_2 \\ c_1 s_2 & -r s_1 s_2 & r c_1 c_2 \\ s_1 & r c_1 & 0 \end{pmatrix}$$

The volume of the hypersphere for $d = 3$ is obtained via a triple integral with $r > 0$, $-\pi/2 \leq \theta_1 \leq \pi/2$, and $0 \leq \theta_2 \leq 2\pi$

$$\begin{aligned} \text{vol}(S_3(r)) &= \int_r \int_{\theta_1} \int_{\theta_2} |\det(J(\theta_1, \theta_2))| dr d\theta_1 d\theta_2 \\ &= \frac{4}{3} \pi r^3 \end{aligned}$$



Hypersphere Volume in d Dimensions

The determinant of the d -dimensional Jacobian matrix is

$$\det(J(\theta_1, \theta_2, \dots, \theta_{d-1})) = (-1)^d r^{d-1} c_1^{d-2} c_2^{d-3} \dots c_{d-2}$$

The volume of the hypersphere is given by the d -dimensional integral with $r > 0$, $-\pi/2 \leq \theta_i \leq \pi/2$ for all $i = 1, \dots, d-2$, and $0 \leq \theta_{d-1} \leq 2\pi$:

$$\begin{aligned}\text{vol}(S_d(r)) &= \int_r \int_{\theta_1} \int_{\theta_2} \cdots \int_{\theta_{d-1}} \left| \det(J(\theta_1, \theta_2, \dots, \theta_{d-1})) \right| dr d\theta_1 d\theta_2 \dots d\theta_{d-1} \\ &= \int_0^r r^{d-1} dr \int_{-\pi/2}^{\pi/2} c_1^{d-2} d\theta_1 \cdots \int_{-\pi/2}^{\pi/2} c_{d-2} d\theta_{d-2} \int_0^{2\pi} d\theta_{d-1} \\ &= \frac{r^d}{d} \frac{\Gamma(\frac{d-1}{2}) \Gamma(\frac{1}{2})}{\Gamma(\frac{d}{2})} \frac{\Gamma(\frac{d-2}{2}) \Gamma(\frac{1}{2})}{\Gamma(\frac{d-1}{2})} \cdots \frac{\Gamma(1) \Gamma(\frac{1}{2})}{\Gamma(\frac{3}{2})} 2\pi \\ &= \frac{\pi \Gamma(\frac{1}{2})^{d/2-1} r^d}{\frac{d}{2} \Gamma(\frac{d}{2})} \\ &= \left(\frac{\pi^{d/2}}{\Gamma(\frac{d}{2} + 1)} \right) r^d\end{aligned}$$

Data mining and Machine learning

Part 7. Dimensionality Reduction

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

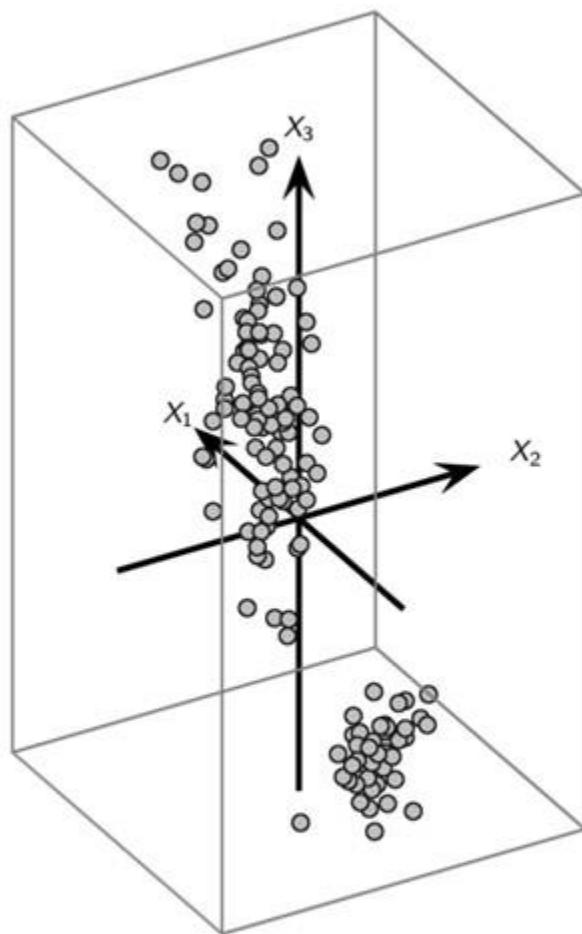
Given the projected vector $\mathbf{x}' = \mathbf{P}_r \mathbf{x}$, the corresponding *error vector*, is the projection onto the remaining $d - r$ basis vectors

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

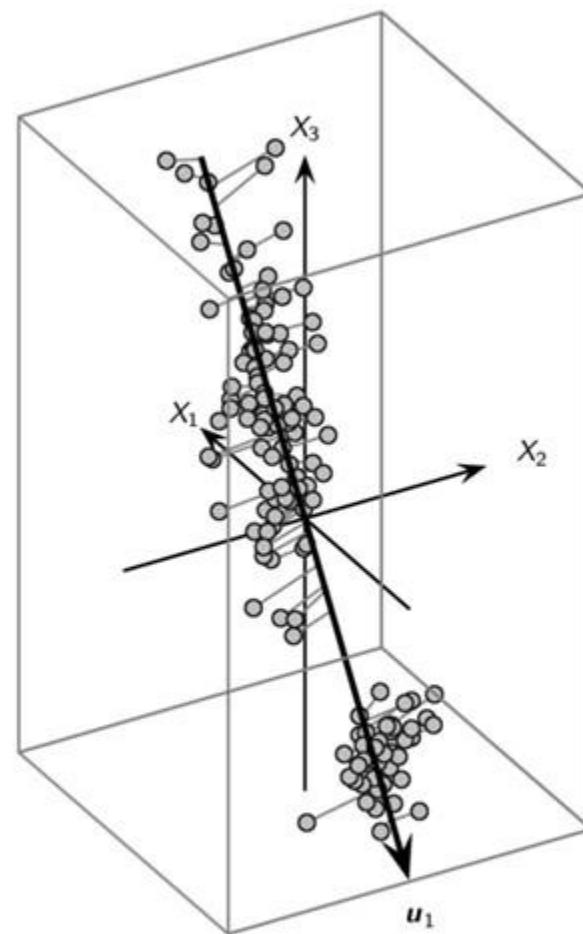
The error vector ϵ is orthogonal to \mathbf{x}' .

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

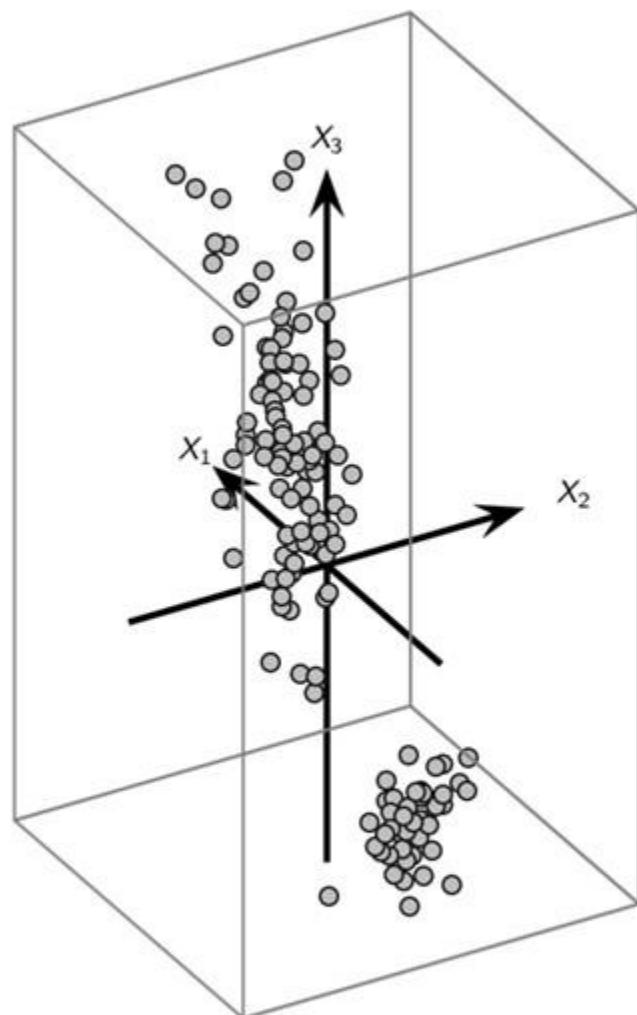


Iris Data: 3D

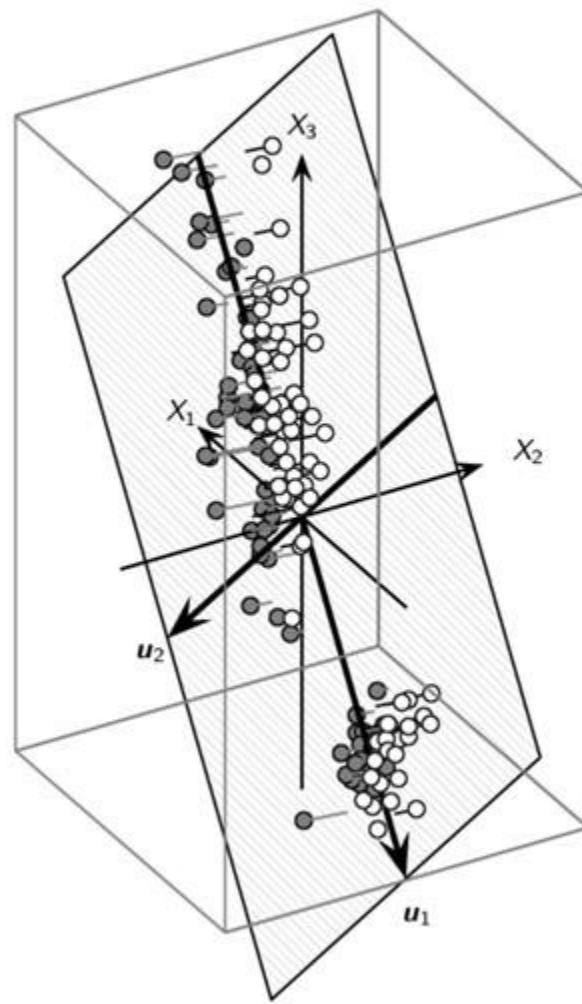


Optimal 1D Basis

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 Σ be its covariance matrix.

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

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

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

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

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

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

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

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

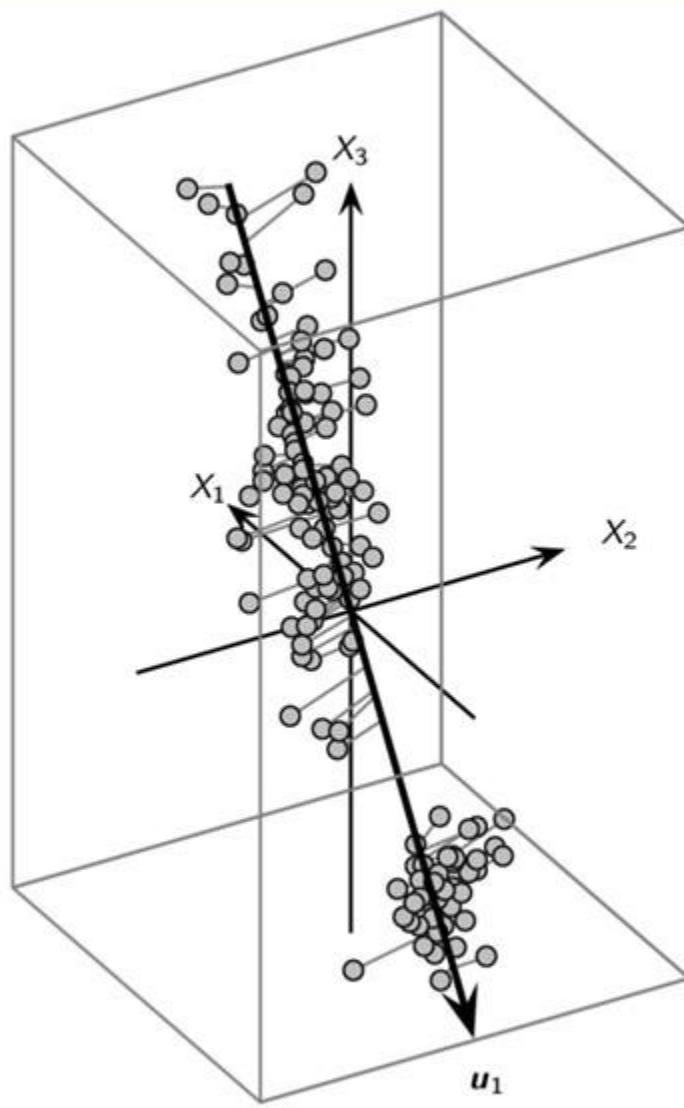
Thus α is an eigenvalue of the covariance matrix Σ , with the associated eigenvector \mathbf{u} .

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 λ_1 of Σ , 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

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

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 λ_1 and λ_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

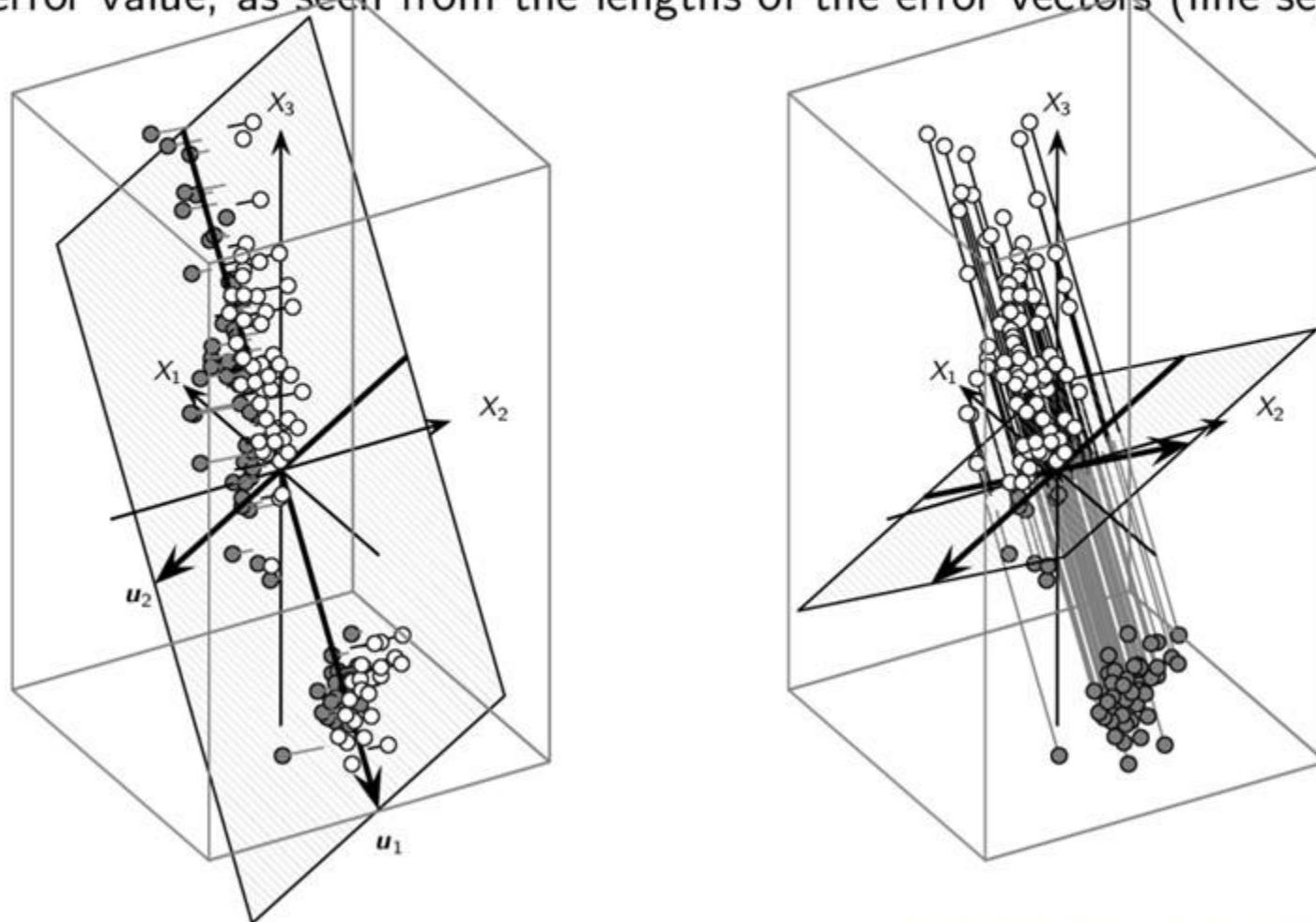
$$var(\mathbf{A}) = \mathbf{u}_1^T \Sigma \mathbf{u}_1 + \mathbf{u}_2^T \Sigma \mathbf{u}_2 = \mathbf{u}_1^T \lambda_1 \mathbf{u}_1 + \mathbf{u}_2^T \lambda_2 \mathbf{u}_2 = \lambda_1 + \lambda_2$$

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 Σ . Because Σ 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 α , 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 α fraction (say 0.9) of the total variance.

Principal Component Analysis: Algorithm

PCA (D, α):

- 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 Σ

$$\lambda_1 = 3.662$$

$$\lambda_2 = 0.239$$

$$\lambda_3 = 0.059$$

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

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

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

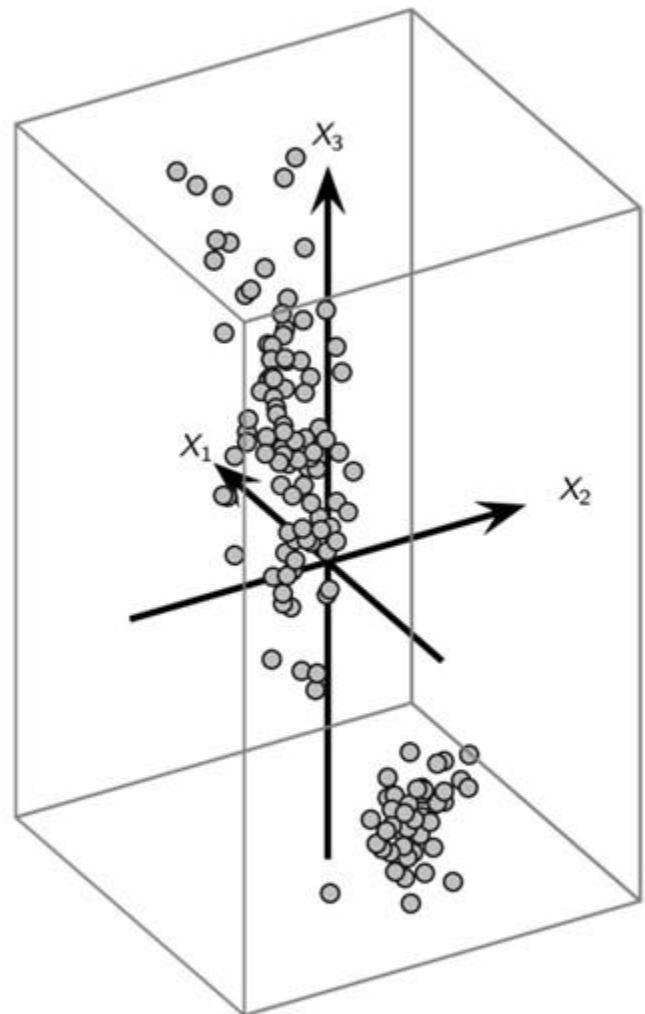
The total variance is therefore $\lambda_1 + \lambda_2 + \lambda_3 = 3.662 + 0.239 + 0.059 = 3.96$.

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

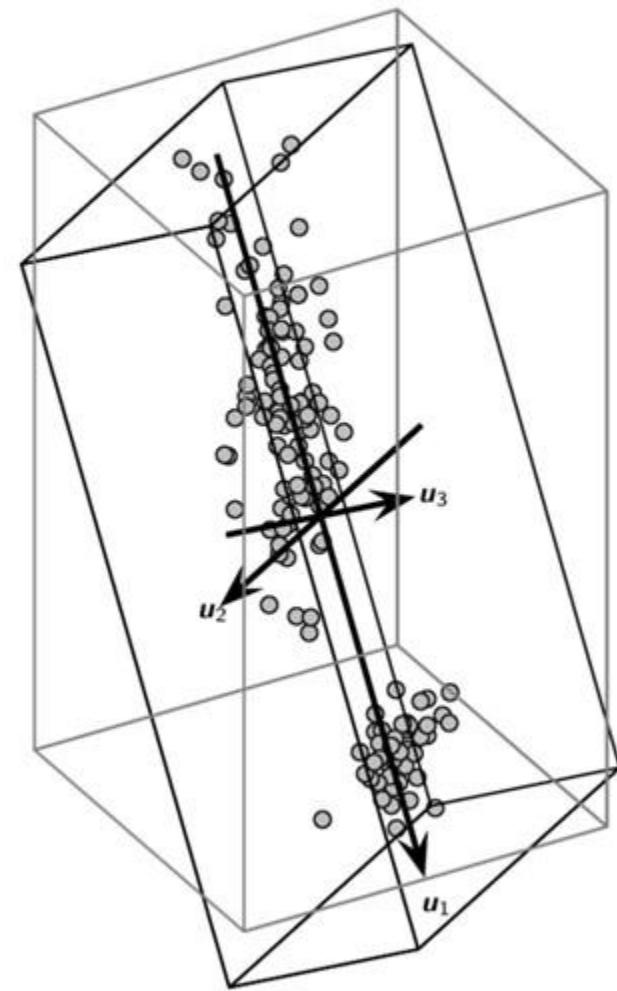
r	1	2	3
$f(r)$	0.925	0.985	1.0

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

Iris Data: Optimal 3D PC Basis

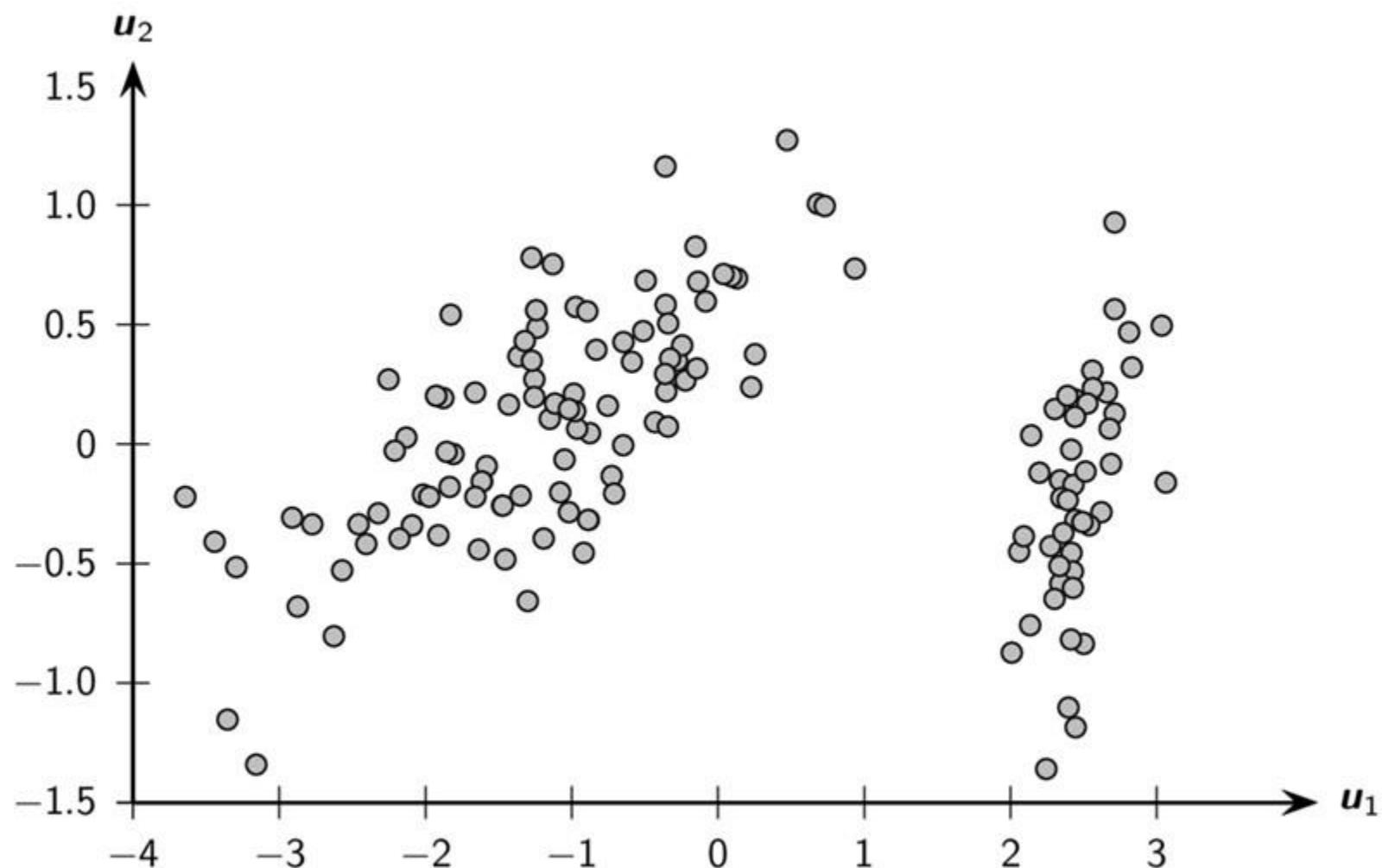


Iris Data (3D)



Optimal 3D Basis

Iris Principal Components: Projected Data (2D)



Geometry of PCA

Geometrically, when $r = d$, PCA corresponds to an 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 Σ

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

which can be written compactly in matrix notation:

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

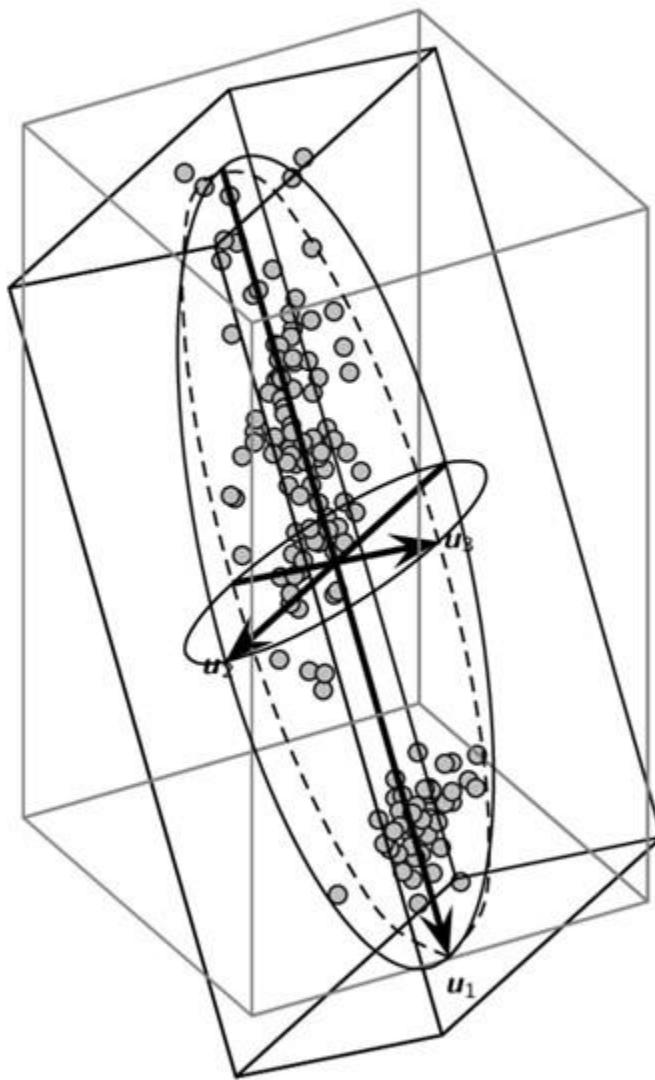
Thus, Λ represents the covariance matrix in the new PC basis.

In the new PC basis, the equation

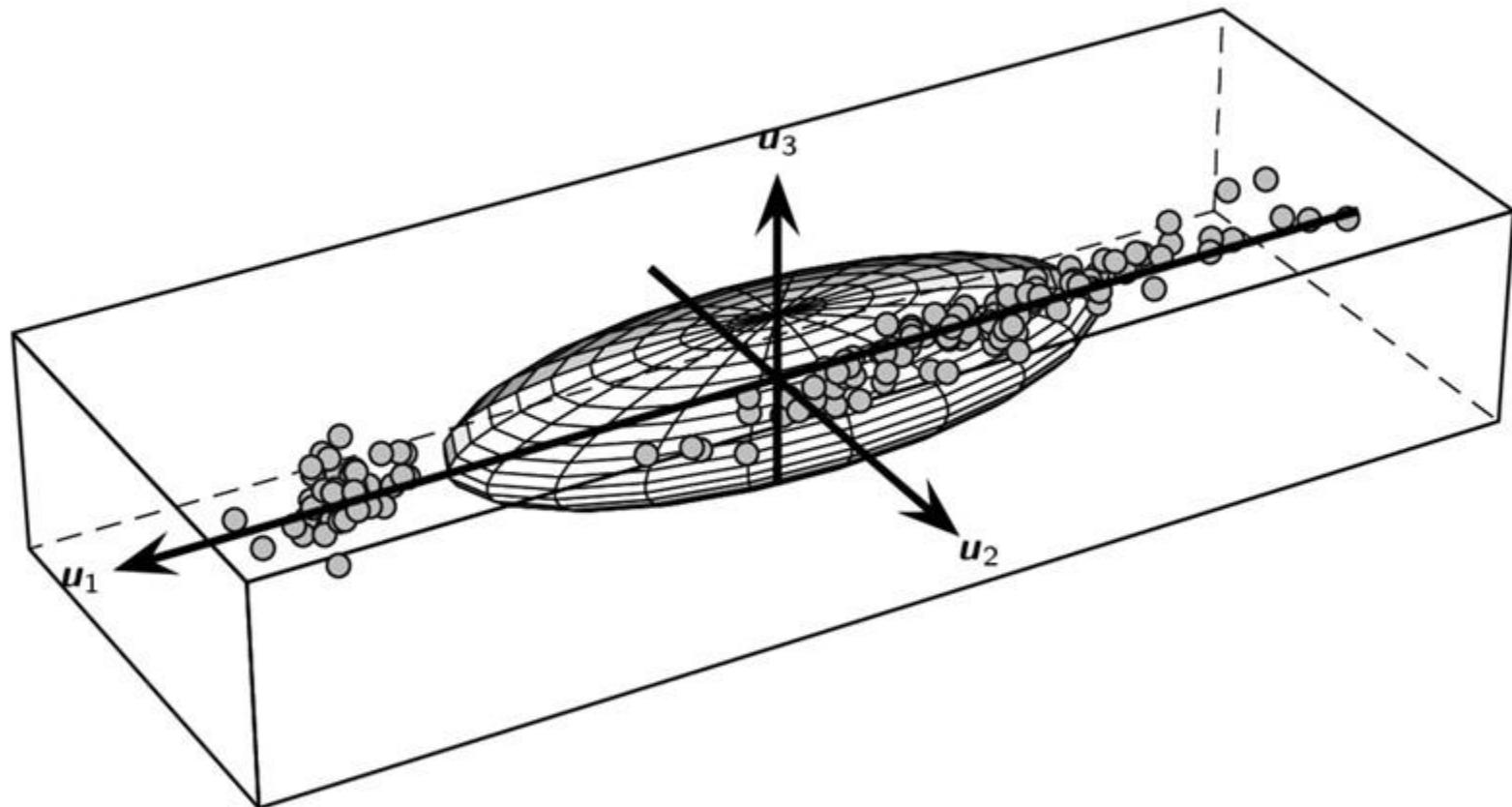
$$\mathbf{x}^T \Sigma^{-1} \mathbf{x} = 1$$

defines a d -dimensional ellipsoid (or hyper-ellipse). The eigenvectors \mathbf{u}_i of Σ , 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



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 ϕ be a function that maps a point x in input space to its image $\phi(x_i)$ in feature space. Let the points in feature space be centered and let Σ_ϕ 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 Σ_ϕ 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 η_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

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

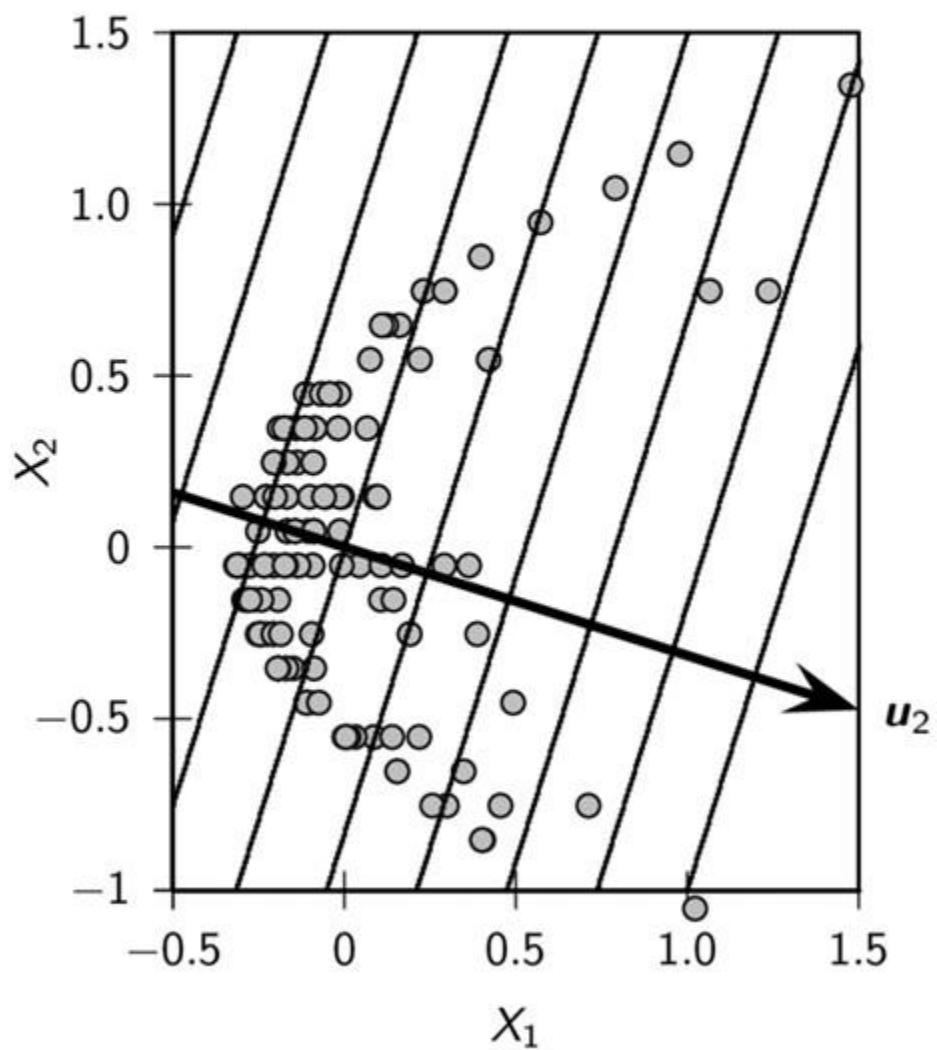
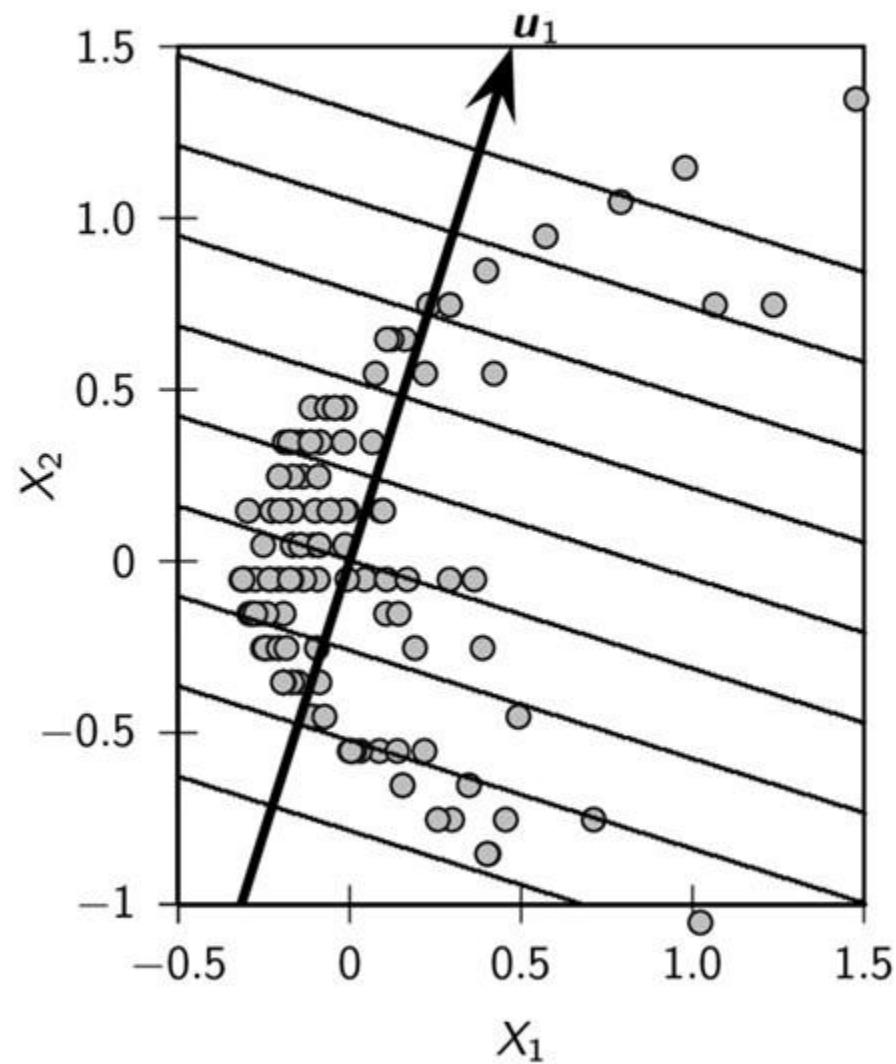
If we sort the eigenvalues of \mathbf{K} in decreasing order $\eta_1 \geq \eta_2 \geq \dots \geq \eta_n \geq 0$, we can obtain the j th principal component as the corresponding eigenvector \mathbf{c}_j . The variance along the j th principal component is given as $\lambda_j = \frac{\eta_j}{n}$.

Kernel PCA Algorithm

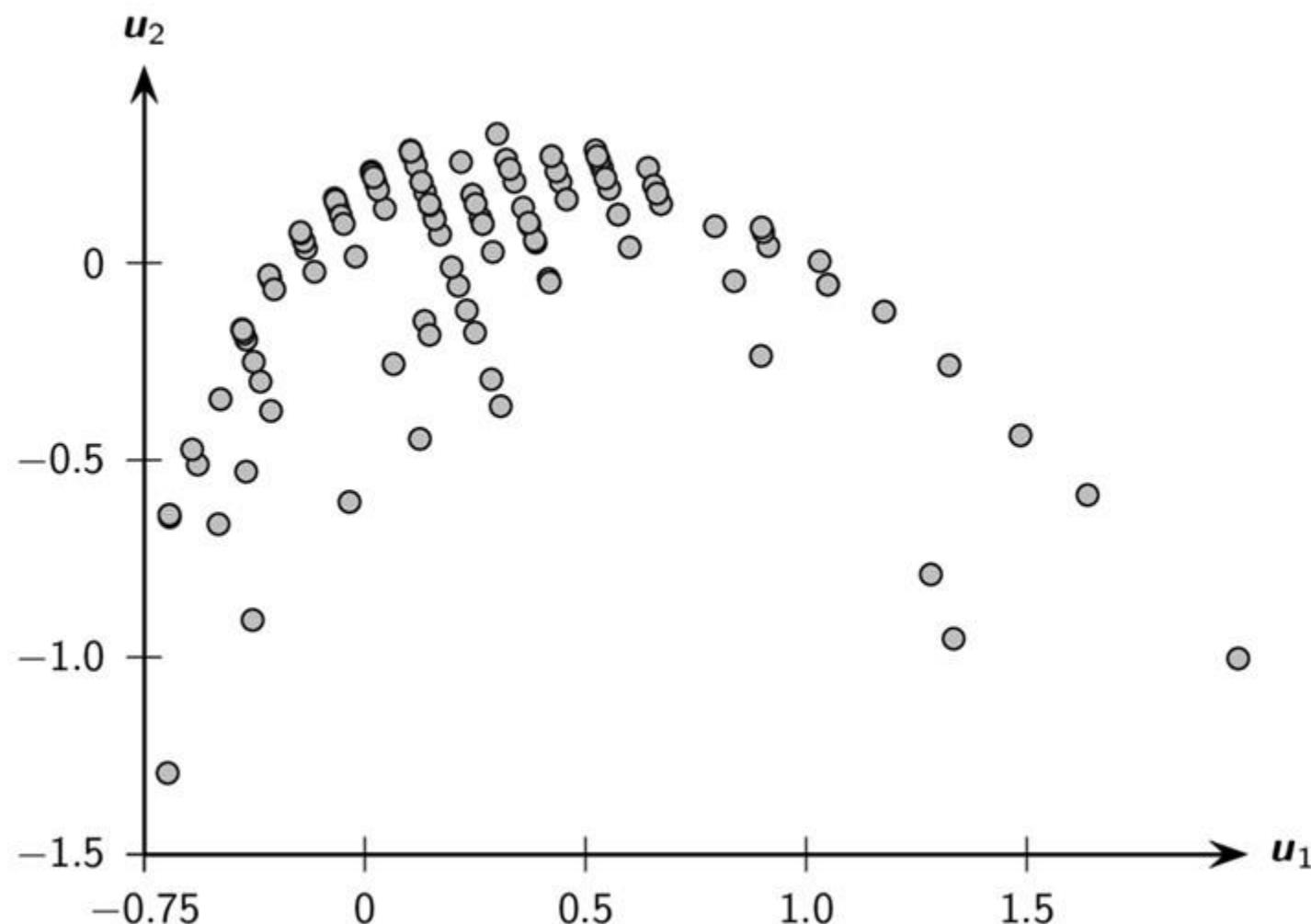
KernelPCA (D, K, α):

- 1 $K = \{K(x_i, x_j)\}_{i,j=1,\dots,n}$ // compute $n \times n$ kernel matrix
- 2 $K = (I - \frac{1}{n}1_{n \times n})K(I - \frac{1}{n}1_{n \times n})$ // center the kernel matrix
- 3 $(\eta_1, \eta_2, \dots, \eta_d) = \text{eigenvalues}(K)$ // compute eigenvalues
- 4 $(c_1 \ c_2 \ \dots \ c_n) = \text{eigenvectors}(K)$ // compute eigenvectors
- 5 $\lambda_i = \frac{\eta_i}{n}$ for all $i = 1, \dots, n$ // compute variance for each component
- 6 $c_i = \sqrt{\frac{1}{\eta_i}} \cdot c_i$ for all $i = 1, \dots, n$ // ensure that $u_i^T u_i = 1$
- 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 $C_r = (c_1 \ c_2 \ \dots \ c_r)$ // reduced basis
- 10 $A = \{a_i \mid a_i = C_r^T K_i, \text{for } i = 1, \dots, n\}$ // reduced dimensionality data

Nonlinear Iris Data: PCA in Input Space

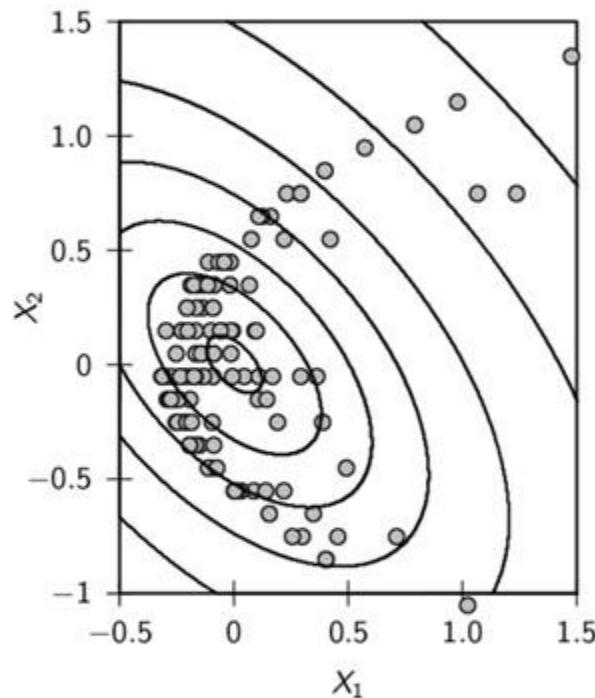


Nonlinear Iris Data: Projection onto PCs

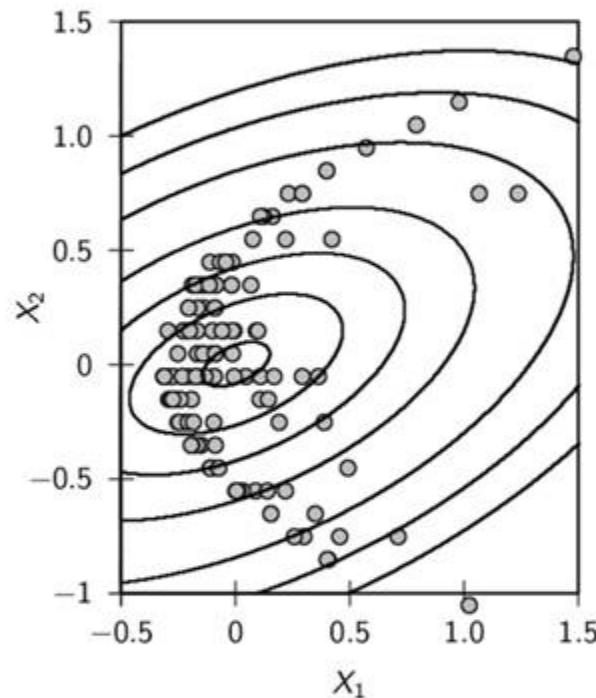


Kernel PCA: 3 PCs (Contours of Constant Projection)

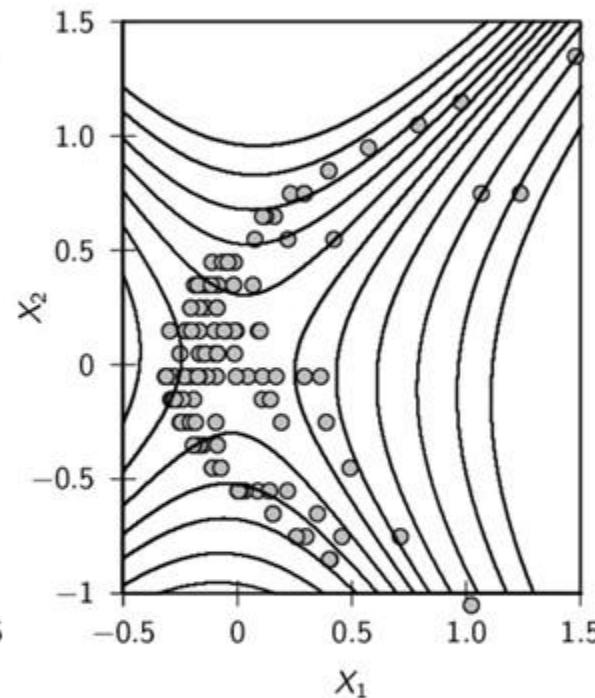
Homogeneous Quadratic Kernel: $K(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i^T \mathbf{x}_j)^2$



(a) $\lambda_1 = 0.2067$



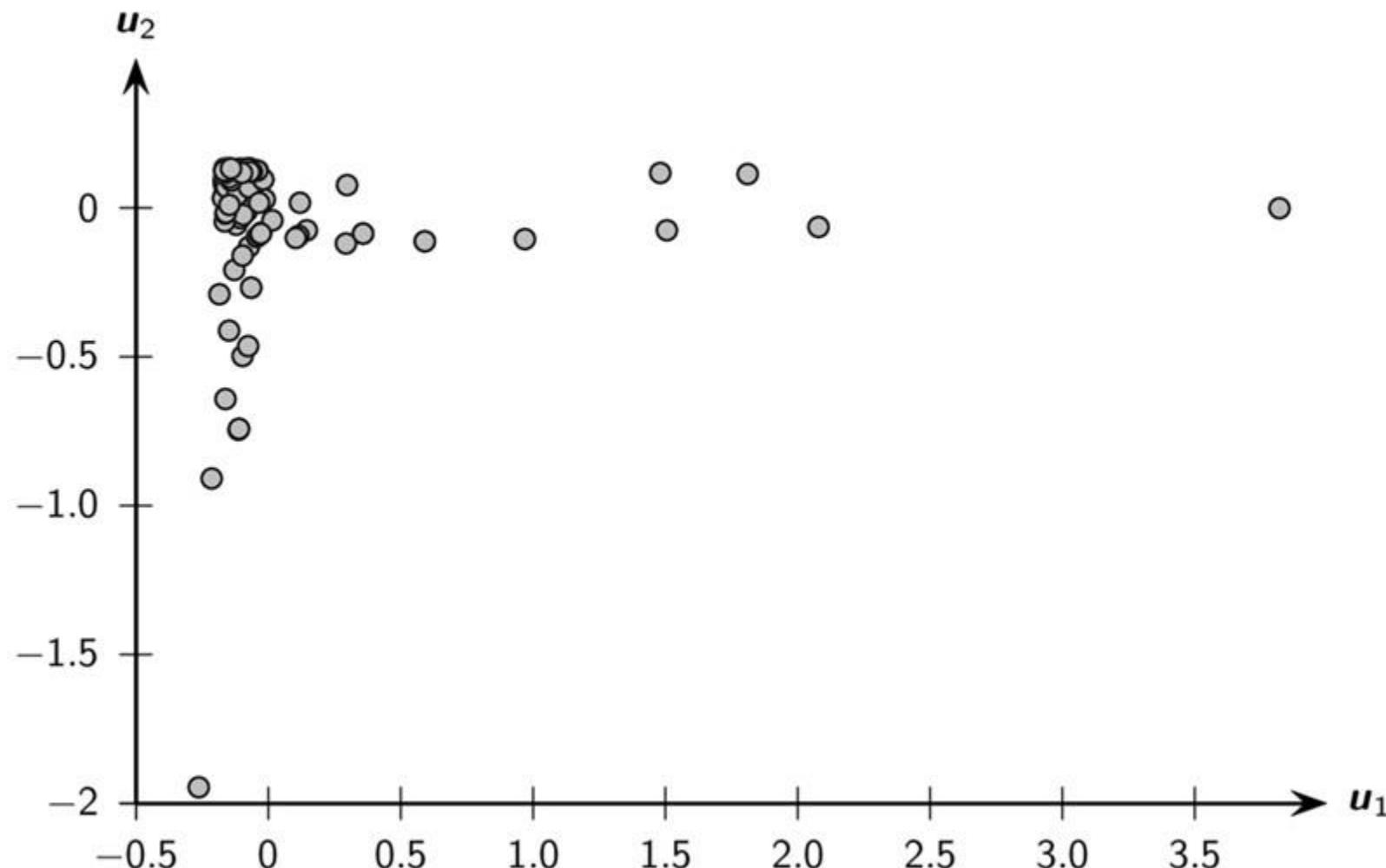
(b) $\lambda_2 = 0.0596$



(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 Λ 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 Δ 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 δ_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 Δ_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 Δ_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 Σ . The singular values of \mathbf{D} are related to the eigenvalues of Σ 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 δ_i^2 .

Q&A

