

# Anomaly Detection

Dimensionality reduction: PCA, robust PCA

Paul Irofti  
Andrei Pătrașcu  
Cristian Rusu

Computer Science Department  
Faculty of Mathematics and Computer Science  
University of Bucharest  
Email: `first.last@fmi.unibuc.ro`



- ▶ eigenvalue and singular value decomposition
- ▶ Principal Component Analysis
- ▶ Robust PCA
- ▶ Matrix Factorization

The course references are Aggarwal 2017, Ch.3 with papers for Robust PCA by Candès et al. 2011 and Netrapalli et al. 2014. For a thorough recap of eigen and singular values see Golub and Van Loan 2013.



# Preliminaries



# Eigenvalues and Eigenvalue Decomposition (EVD)

Given square matrix  $A \in \mathbb{R}^{n \times n}$  then its eigenvalues  $\lambda$  and associated eigenvectors  $\mathbf{v}$  follow:

$$\mathbf{A}\mathbf{v} = \lambda\mathbf{v} \quad (1)$$

which can be obtained through  $(\mathbf{A} - \lambda\mathbf{I}_n)\mathbf{v} = 0$  or the factorization of the polynomial  $\det(\mathbf{A} - \lambda\mathbf{I}_n) = 0$



# Eigenvalues and Eigenvalue Decomposition (EVD)

Given square matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$  then its eigenvalues  $\lambda$  and associated eigenvectors  $\mathbf{v}$  follow:

$$\mathbf{A}\mathbf{v} = \lambda\mathbf{v} \quad (1)$$

which can be obtained through  $(\mathbf{A} - \lambda\mathbf{I}_n)\mathbf{v} = 0$  or the factorization of the polynomial  $\det(\mathbf{A} - \lambda\mathbf{I}_n) = 0$

There are  $n$  eigenvalues and eigenvectors, thus we can write the eigenvalue decomposition (EVD):

$$\mathbf{A}\mathbf{V} = \mathbf{V}\mathbf{\Lambda} \iff \mathbf{\Lambda} = \mathbf{V}^{-1}\mathbf{A}\mathbf{V} \quad (2)$$



# Eigenvalues and Eigenvalue Decomposition (EVD)

Given square matrix  $A \in \mathbb{R}^{n \times n}$  then its eigenvalues  $\lambda$  and associated eigenvectors  $\mathbf{v}$  follow:

$$\mathbf{A}\mathbf{v} = \lambda\mathbf{v} \quad (1)$$

which can be obtained through  $(\mathbf{A} - \lambda\mathbf{I}_n)\mathbf{v} = 0$  or the factorization of the polynomial  $\det(\mathbf{A} - \lambda\mathbf{I}_n) = 0$

There are  $n$  eigenvalues and eigenvectors, thus we can write the eigenvalue decomposition (EVD):

$$\mathbf{A}\mathbf{V} = \mathbf{V}\mathbf{\Lambda} \iff \mathbf{\Lambda} = \mathbf{V}^{-1}\mathbf{A}\mathbf{V} \quad (2)$$

**Remark:** For symmetric matrices  $\mathbf{A}^\top = \mathbf{A}$  we have  $\mathbf{V}^{-1} = \mathbf{V}^\top$  such that  $\mathbf{A} = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^\top$ .



# Singular Values and Singular Value Decomposition (SVD)

Given rectangular matrix  $A \in \mathbb{R}^{n \times m}$  the singular values  $\sigma$ , the associated left-hand side singular vectors  $\mathbf{u}$ , and associated right-hand side singular vectors  $\mathbf{v}$

$$\mathbf{A}\mathbf{v} = \sigma\mathbf{v}\mathbf{A}^T\mathbf{u} = \sigma\mathbf{u} \quad (3)$$



# Singular Values and Singular Value Decomposition (SVD)

Given rectangular matrix  $A \in \mathbb{R}^{n \times m}$  the singular values  $\sigma$ , the associated left-hand side singular vectors  $\mathbf{u}$ , and associated right-hand side singular vectors  $\mathbf{v}$

$$\mathbf{A}\mathbf{v} = \sigma\mathbf{v}\mathbf{A}^\top\mathbf{u} = \sigma\mathbf{u} \quad (3)$$

There are  $\min(n, m)$  singular values,  $n$  left singular vectors, and  $m$  right singular vectors, thus we can write the singular value decomposition (SVD):

$$\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top = \sum_{i=1}^{\min\{n,m\}} \sigma_i \mathbf{u}_i \mathbf{v}_i^\top \quad (4)$$





# Singular Values and Singular Value Decomposition (SVD)

Given rectangular matrix  $A \in \mathbb{R}^{n \times m}$  the singular values  $\sigma$ , the associated left-hand side singular vectors  $\mathbf{u}$ , and associated right-hand side singular vectors  $\mathbf{v}$

$$\mathbf{A}\mathbf{v} = \sigma\mathbf{v}\mathbf{A}^\top\mathbf{u} = \sigma\mathbf{u} \quad (3)$$

There are  $\min(n, m)$  singular values,  $n$  left singular vectors, and  $m$  right singular vectors, thus we can write the singular value decomposition (SVD):

$$\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top = \sum_{i=1}^{\min\{n,m\}} \sigma_i \mathbf{u}_i \mathbf{v}_i^\top \quad (4)$$

**Theorem:** The optimal low-rank matrix  $\mathbf{L} \in \mathbb{R}^{n \times m}$  with rank  $k$  that approximates  $\mathbf{A}$  is  $\sum_{i=1}^k \sigma_i \mathbf{u}_i \mathbf{v}_i^\top$



# Least Squares (LS)

Given data  $\mathbf{X} \in \mathbb{R}^{N \times d}$ , where  $d$  is the data dimension and  $N$  the number of samples, the least-squares problem solves the following:

$$\min_{\beta} \|\mathbf{X}\beta - \mathbf{y}\|_F^2 \quad (5)$$

- ▶ when  $N = d$  we have  $\beta^* = \mathbf{X}^{-1}\mathbf{y}$
- ▶ when  $N > d$  we have  $\beta^* = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y}$
- ▶ when  $N < d$  we have  $\beta^* = \mathbf{X}^T(\mathbf{X}\mathbf{X}^T)^{-1}\mathbf{y}$

See Lecture 2 for more details.



# LS: line fit

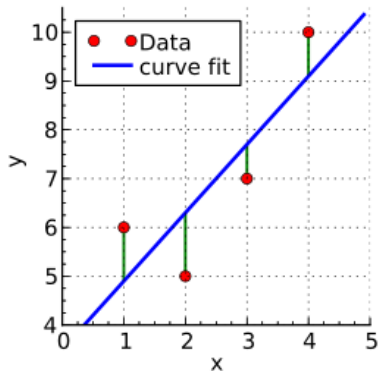


Figure: LS fits 2D points on a line

Source: [https://en.wikipedia.org/wiki/Linear\\_least\\_squares](https://en.wikipedia.org/wiki/Linear_least_squares)



# LS: projection

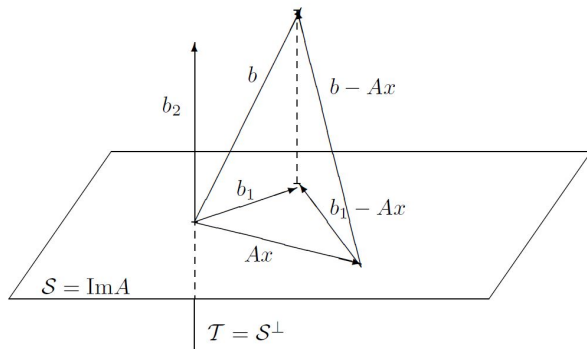


Figure: LS projects vectors on  $\text{Im} A$



Least-squares properties:

- ▶ finds  $d - 1$  subspace or hyperplane
- ▶ the hyperplane is an optimum fit to data
- ▶ anomaly score: based on length on orthogonal direction



Least-squares properties:

- ▶ finds  $d - 1$  subspace or hyperplane
- ▶ the hyperplane is an optimum fit to data
- ▶ anomaly score: based on length on orthogonal direction

Generalization:

- ▶ what is the  $k < d$  subspace or hyperplane?
- ▶ what is the anomaly score then?
- ▶ what is an optimum fit on any  $k$ -dimensional subspace?



# Nonlinear PCA



# Principal Component Analysis (PCA)

PCA starts from the covariance matrix of the mean-centered data matrix  $\mathbf{X} \in \mathbb{R}^{N \times d}$

$$\Sigma = \frac{\mathbf{X}^T \mathbf{X}}{N} \quad (6)$$

such that  $\Sigma \in \mathbb{R}^{d \times d}$  where element  $\Sigma_{ij}$  is the covariance between data dimension  $i$  and  $j$ .





# Principal Component Analysis (PCA)

PCA starts from the covariance matrix of the mean-centered data matrix  $\mathbf{X} \in \mathbb{R}^{N \times d}$

$$\Sigma = \frac{\mathbf{X}^T \mathbf{X}}{N} \quad (6)$$

such that  $\Sigma \in \mathbb{R}^{d \times d}$  where element  $\Sigma_{ij}$  is the covariance between data dimension  $i$  and  $j$ .

Properties:

- ▶ the covariance matrix is symmetric and positive definite



# Principal Component Analysis (PCA)

PCA starts from the covariance matrix of the mean-centered data matrix  $\mathbf{X} \in \mathbb{R}^{N \times d}$

$$\Sigma = \frac{\mathbf{X}^T \mathbf{X}}{N} \quad (6)$$

such that  $\Sigma \in \mathbb{R}^{d \times d}$  where element  $\Sigma_{ij}$  is the covariance between data dimension  $i$  and  $j$ .

Properties:

- ▶ the covariance matrix is symmetric and positive definite
- ▶ the EVD of  $\Sigma = \mathbf{P} \Delta \mathbf{P}^T$



# Principal Component Analysis (PCA)

PCA starts from the covariance matrix of the mean-centered data matrix  $\mathbf{X} \in \mathbb{R}^{N \times d}$

$$\Sigma = \frac{\mathbf{X}^T \mathbf{X}}{N} \quad (6)$$

such that  $\Sigma \in \mathbb{R}^{d \times d}$  where element  $\Sigma_{ij}$  is the covariance between data dimension  $i$  and  $j$ .

Properties:

- ▶ the covariance matrix is symmetric and positive definite
- ▶ the EVD of  $\Sigma = \mathbf{P} \Delta \mathbf{P}^T$
- ▶  $\Delta$  is diagonal and contains the eigenvalues between  $\lambda_{max}$  and  $\lambda_{min}$



# Principal Component Analysis (PCA)

PCA starts from the covariance matrix of the mean-centered data matrix  $\mathbf{X} \in \mathbb{R}^{N \times d}$

$$\Sigma = \frac{\mathbf{X}^T \mathbf{X}}{N} \quad (6)$$

such that  $\Sigma \in \mathbb{R}^{d \times d}$  where element  $\Sigma_{ij}$  is the covariance between data dimension  $i$  and  $j$ .

Properties:

- ▶ the covariance matrix is symmetric and positive definite
- ▶ the EVD of  $\Sigma = \mathbf{P} \Delta \mathbf{P}^T$
- ▶  $\Delta$  is diagonal and contains the eigenvalues between  $\lambda_{max}$  and  $\lambda_{min}$
- ▶  $\mathbf{P} \in \mathbb{R}^{d \times d}$  represents the orthonormal eigenvectors of the covariance corresponding to  $\Delta$



# Principal Component Analysis (PCA)

PCA starts from the covariance matrix of the mean-centered data matrix  $\mathbf{X} \in \mathbb{R}^{N \times d}$

$$\Sigma = \frac{\mathbf{X}^T \mathbf{X}}{N} \quad (6)$$

such that  $\Sigma \in \mathbb{R}^{d \times d}$  where element  $\Sigma_{ij}$  is the covariance between data dimension  $i$  and  $j$ .

Properties:

- ▶ the covariance matrix is symmetric and positive definite
- ▶ the EVD of  $\Sigma = \mathbf{P} \Delta \mathbf{P}^T$
- ▶  $\Delta$  is diagonal and contains the eigenvalues between  $\lambda_{max}$  and  $\lambda_{min}$
- ▶  $\mathbf{P} \in \mathbb{R}^{d \times d}$  represents the orthonormal eigenvectors of the covariance corresponding to  $\Delta$
- ▶ the normal hyperplane to  $\mathbf{p}_{min} \in \mathbf{P}$  is the *LS*-hyperplane of dimension  $k = d - 1$



# PCA as Generalized LS

PCA starts from the covariance matrix of the mean-centered data matrix  $\mathbf{X} \in \mathbb{R}^{N \times d}$

$$\Sigma = \frac{\mathbf{X}^T \mathbf{X}}{N} = \mathbf{P} \Delta \mathbf{P}^T \quad (7)$$

such that  $\Sigma \in \mathbb{R}^{d \times d}$ .



# PCA as Generalized LS

PCA starts from the covariance matrix of the mean-centered data matrix  $\mathbf{X} \in \mathbb{R}^{N \times d}$

$$\Sigma = \frac{\mathbf{X}^T \mathbf{X}}{N} = \mathbf{P} \Delta \mathbf{P}^T \quad (7)$$

such that  $\Sigma \in \mathbb{R}^{d \times d}$ .

**Remark:** the normal hyperplane to  $\mathbf{p}_{min} \in \mathbf{P}$  is the *LS*-hyperplane of dimension  $k = d - 1$ .



# PCA as Generalized LS

PCA starts from the covariance matrix of the mean-centered data matrix  $\mathbf{X} \in \mathbb{R}^{N \times d}$

$$\Sigma = \frac{\mathbf{X}^T \mathbf{X}}{N} = \mathbf{P} \Delta \mathbf{P}^T \quad (7)$$

such that  $\Sigma \in \mathbb{R}^{d \times d}$ .

**Remark:** the normal hyperplane to  $\mathbf{p}_{min} \in \mathbf{P}$  is the *LS*-hyperplane of dimension  $k = d - 1$ .

This implies that:

- ▶ the eigenvectors subspace corresponding to the largest  $d - 1$  eigenvalues provides a good data approximation





# PCA as Generalized LS

PCA starts from the covariance matrix of the mean-centered data matrix  $\mathbf{X} \in \mathbb{R}^{N \times d}$

$$\Sigma = \frac{\mathbf{X}^T \mathbf{X}}{N} = \mathbf{P} \Delta \mathbf{P}^T \quad (7)$$

such that  $\Sigma \in \mathbb{R}^{d \times d}$ .

**Remark:** the normal hyperplane to  $\mathbf{p}_{min} \in \mathbf{P}$  is the *LS*-hyperplane of dimension  $k = d - 1$ .

This implies that:

- ▶ the eigenvectors subspace corresponding to the largest  $d - 1$  eigenvalues provides a good data approximation
- ▶ what about  $k = d - 2$  or any  $k \in [d - 1] = \{1, \dots, d - 1\}$ ?



# PCA as Generalized LS

PCA starts from the covariance matrix of the mean-centered data matrix  $\mathbf{X} \in \mathbb{R}^{N \times d}$

$$\Sigma = \frac{\mathbf{X}^T \mathbf{X}}{N} = \mathbf{P} \Delta \mathbf{P}^T \quad (7)$$

such that  $\Sigma \in \mathbb{R}^{d \times d}$ .

**Remark:** the normal hyperplane to  $\mathbf{p}_{min} \in \mathbf{P}$  is the *LS*-hyperplane of dimension  $k = d - 1$ .

This implies that:

- ▶ the eigenvectors subspace corresponding to the largest  $d - 1$  eigenvalues provides a good data approximation
- ▶ what about  $k = d - 2$  or any  $k \in [d - 1] = \{1, \dots, d - 1\}$ ?
- ▶ the subspace will be generated by the largest  $k$  eigenvectors such that the residual is minimal



# PCA as Generalized LS

PCA starts from the covariance matrix of the mean-centered data matrix  $\mathbf{X} \in \mathbb{R}^{N \times d}$

$$\Sigma = \frac{\mathbf{X}^T \mathbf{X}}{N} = \mathbf{P} \Delta \mathbf{P}^T \quad (7)$$

such that  $\Sigma \in \mathbb{R}^{d \times d}$ .

**Remark:** the normal hyperplane to  $\mathbf{p}_{min} \in \mathbf{P}$  is the *LS*-hyperplane of dimension  $k = d - 1$ .

This implies that:

- ▶ the eigenvectors subspace corresponding to the largest  $d - 1$  eigenvalues provides a good data approximation
- ▶ what about  $k = d - 2$  or any  $k \in [d - 1] = \{1, \dots, d - 1\}$ ?
- ▶ the subspace will be generated by the largest  $k$  eigenvectors such that the residual is minimal
- ▶ anomalies are the data whose error is high in this new subspace



# PCA as Generalized LS

PCA starts from the covariance matrix of the mean-centered data matrix  $\mathbf{X} \in \mathbb{R}^{N \times d}$

$$\Sigma = \frac{\mathbf{X}^T \mathbf{X}}{N} = \mathbf{P} \Delta \mathbf{P}^T \quad (7)$$

such that  $\Sigma \in \mathbb{R}^{d \times d}$ .

**Remark:** the normal hyperplane to  $\mathbf{p}_{min} \in \mathbf{P}$  is the *LS*-hyperplane of dimension  $k = d - 1$ .

This implies that:

- ▶ the eigenvectors subspace corresponding to the largest  $d - 1$  eigenvalues provides a good data approximation
- ▶ what about  $k = d - 2$  or any  $k \in [d - 1] = \{1, \dots, d - 1\}$ ?
- ▶ the subspace will be generated by the largest  $k$  eigenvectors such that the residual is minimal
- ▶ anomalies are the data whose error is high in this new subspace
- ▶ anomalies have a large normal component



# Example: Principal Eigenvectors

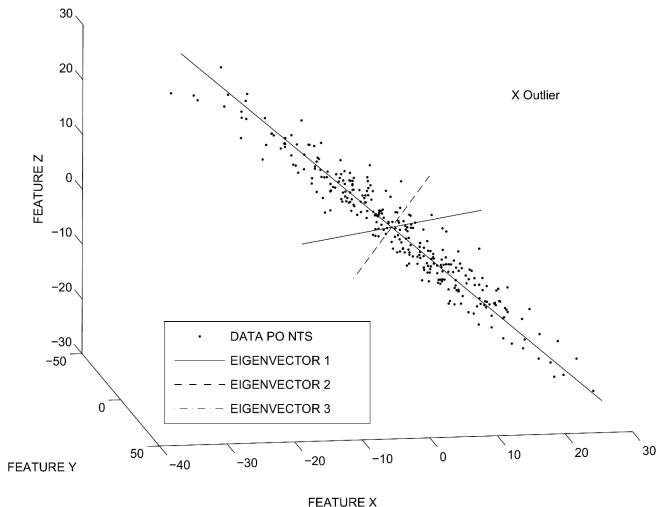


Figure: Distribution along first  $k = 3$  eigenvectors (Aggarwal 2017)



**Remark:** the subspace will be generated by the largest  $k$  eigenvectors such that the residual is minimal.



**Remark:** the subspace will be generated by the largest  $k$  eigenvectors such that the residual is minimal.

Implications:

- ▶ the new axis of the data representation are the  $k$  orthonormal eigenvectors



**Remark:** the subspace will be generated by the largest  $k$  eigenvectors such that the residual is minimal.

Implications:

- ▶ the new axis of the data representation are the  $k$  orthonormal eigenvectors
- ▶ there is no covariance in this new subspace as the eigenvectors are orthogonal





**Remark:** the subspace will be generated by the largest  $k$  eigenvectors such that the residual is minimal.

Implications:

- ▶ the new axis of the data representation are the  $k$  orthonormal eigenvectors
- ▶ there is no covariance in this new subspace as the eigenvectors are orthogonal
- ▶ the variance along each axis is the eigenvalue



**Remark:** the subspace will be generated by the largest  $k$  eigenvectors such that the residual is minimal.

Implications:

- ▶ the new axis of the data representation are the  $k$  orthonormal eigenvectors
- ▶ there is no covariance in this new subspace as the eigenvectors are orthogonal
- ▶ the variance along each axis is the eigenvalue
- ▶ a small eigenvalue implies a low variance



**Remark:** the subspace will be generated by the largest  $k$  eigenvectors such that the residual is minimal.

Implications:

- ▶ the new axis of the data representation are the  $k$  orthonormal eigenvectors
- ▶ there is no covariance in this new subspace as the eigenvectors are orthogonal
- ▶ the variance along each axis is the eigenvalue
- ▶ a small eigenvalue implies a low variance
- ▶ can we cancel axis with small eigenvalues? E.g.  $\lambda < 10^{-3}$



**Remark:** the subspace will be generated by the largest  $k$  eigenvectors such that the residual is minimal.

Implications:

- ▶ the new axis of the data representation are the  $k$  orthonormal eigenvectors
- ▶ there is no covariance in this new subspace as the eigenvectors are orthogonal
- ▶ the variance along each axis is the eigenvalue
- ▶ a small eigenvalue implies a low variance
- ▶ can we cancel axis with small eigenvalues? E.g.  $\lambda < 10^{-3}$
- ▶ not if we expect anomalies to have higher variance among low variance axis



# Example: Eigen Histogram

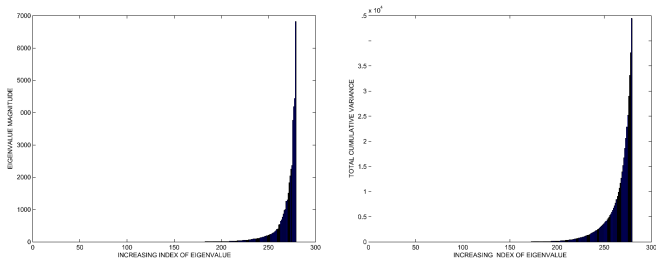
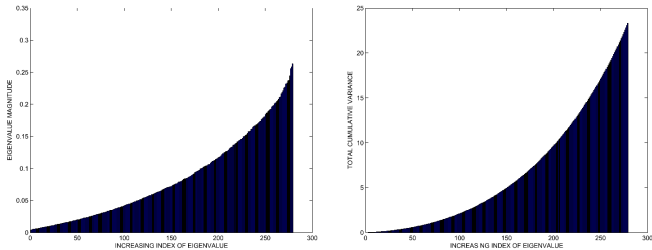


Figure: Eigenvalues magnitude and variance (Aggarwal 2017)



# Example: Eigen Histogram Trimmed



**Figure:** Eigenvalues magnitude and variance after trimming (Aggarwal 2017)



# PCA space

PCA starts from the EVD of the covariance matrix

$$\Sigma = \frac{\mathbf{X}^T \mathbf{X}}{N} = \mathbf{P} \Delta \mathbf{P}^T$$

such that  $\Sigma \in R^{d \times d}$ .



# PCA space

PCA starts from the EVD of the covariance matrix

$$\Sigma = \frac{\mathbf{X}^T \mathbf{X}}{N} = \mathbf{P} \Delta \mathbf{P}^T$$

such that  $\Sigma \in \mathbb{R}^{d \times d}$ .

Then the transformed data space becomes:

$$\mathbf{X}' = \mathbf{X} \mathbf{P} \tag{8}$$

where  $\mathbf{X}' \in \mathbb{R}^{N \times d}$ .





# PCA space

PCA starts from the EVD of the covariance matrix

$$\Sigma = \frac{\mathbf{X}^T \mathbf{X}}{N} = \mathbf{P} \Delta \mathbf{P}^T$$

such that  $\Sigma \in \mathbb{R}^{d \times d}$ .

Then the transformed data space becomes:

$$\mathbf{X}' = \mathbf{X} \mathbf{P} \tag{8}$$

where  $\mathbf{X}' \in \mathbb{R}^{N \times d}$ .

Using the entire space  $k = d$  has the following advantage:



# PCA space

PCA starts from the EVD of the covariance matrix

$$\Sigma = \frac{\mathbf{X}^T \mathbf{X}}{N} = \mathbf{P} \Delta \mathbf{P}^T$$

such that  $\Sigma \in \mathbb{R}^{d \times d}$ .

Then the transformed data space becomes:

$$\mathbf{X}' = \mathbf{X} \mathbf{P} \tag{8}$$

where  $\mathbf{X}' \in \mathbb{R}^{N \times d}$ .

Using the entire space  $k = d$  has the following advantage:

- ▶ take small eigenvector  $\mathbf{p}_j$



# PCA space

PCA starts from the EVD of the covariance matrix

$$\Sigma = \frac{\mathbf{X}^T \mathbf{X}}{N} = \mathbf{P} \Delta \mathbf{P}^T$$

such that  $\Sigma \in \mathbb{R}^{d \times d}$ .

Then the transformed data space becomes:

$$\mathbf{X}' = \mathbf{X} \mathbf{P} \tag{8}$$

where  $\mathbf{X}' \in \mathbb{R}^{N \times d}$ .

Using the entire space  $k = d$  has the following advantage:

- ▶ take small eigenvector  $\mathbf{p}_j$
- ▶ we know that the entries  $x'_{\ell j}$  do not vary much as  $\lambda_j$  is small



# PCA space

PCA starts from the EVD of the covariance matrix

$$\Sigma = \frac{\mathbf{X}^T \mathbf{X}}{N} = \mathbf{P} \Delta \mathbf{P}^T$$

such that  $\Sigma \in \mathbb{R}^{d \times d}$ .

Then the transformed data space becomes:

$$\mathbf{X}' = \mathbf{X} \mathbf{P} \tag{8}$$

where  $\mathbf{X}' \in \mathbb{R}^{N \times d}$ .

Using the entire space  $k = d$  has the following advantage:

- ▶ take small eigenvector  $\mathbf{p}_j$
- ▶ we know that the entries  $x'_{\ell j}$  do not vary much as  $\lambda_j$  is small
- ▶ **outlier:** if  $x'_{ij}$  has a large deviation compared to other  $x'_{\ell j}$  entries



# PCA subspace

PCA starts from the EVD of  $\Sigma = \frac{\mathbf{X}^T \mathbf{X}}{N} = \mathbf{P} \Delta \mathbf{P}^T$ , then the transformed data space becomes  $\mathbf{X}' = \mathbf{X} \mathbf{P}$  where  $\mathbf{X}' \in \mathbb{R}^{N \times d}$ .



PCA starts from the EVD of  $\Sigma = \frac{\mathbf{X}^T \mathbf{X}}{N} = \mathbf{P} \Delta \mathbf{P}^T$ , then the transformed data space becomes  $\mathbf{X}' = \mathbf{X} \mathbf{P}$  where  $\mathbf{X}' \in \mathbb{R}^{N \times d}$ .

The approximation through trimming the smallest  $d - k$  eigenvectors  $\mathbf{X}' = \mathbf{X} \mathbf{P}_k \in \mathbb{R}^{N \times k}$

$$\|\mathbf{X} - \mathbf{X} \mathbf{P}_k\| \tag{9}$$

will contain in each  $x'_{ij}$  with  $j \in [k]$



PCA starts from the EVD of  $\Sigma = \frac{\mathbf{X}^T \mathbf{X}}{N} = \mathbf{P} \Delta \mathbf{P}^T$ , then the transformed data space becomes  $\mathbf{X}' = \mathbf{X} \mathbf{P}$  where  $\mathbf{X}' \in \mathbb{R}^{N \times d}$ .

The approximation through trimming the smallest  $d - k$  eigenvectors  $\mathbf{X}' = \mathbf{X} \mathbf{P}_k \in \mathbb{R}^{N \times k}$

$$\|\mathbf{X} - \mathbf{X} \mathbf{P}_k\| \quad (9)$$

will contain in each  $x'_{ij}$  with  $j \in [k]$

**Hard outlier score:** the residuals representing the distance to the rank- $k$  hyperplane described by  $\mathbf{X} \mathbf{P}_k$ .

**Soft outlier score:** normalize the residuals according to their corresponding variance along the  $d - k$  distances.



**Soft outlier score:** normalize the residuals according to their corresponding variance along the  $d - k$  distances.





**Soft outlier score:** normalize the residuals according to their corresponding variance along the  $d - k$  distances.

Decompose the sum of squares of the  $d - k$  distances and normalize by their corresponding eigenvalue:

$$\text{Score}(\mathbf{x}_\ell) = \sum_{j=k+1}^d \frac{\|\mathbf{x}_{\ell j} - \mathbf{x}_\ell^\top \mathbf{p}_j\|^2}{\lambda_j} \quad (10)$$



**Soft outlier score:** normalize the residuals according to their corresponding variance along the  $d - k$  distances.

Decompose the sum of squares of the  $d - k$  distances and normalize by their corresponding eigenvalue:

$$\text{Score}(\mathbf{x}_\ell) = \sum_{j=k+1}^d \frac{\|\mathbf{x}_{\ell j} - \mathbf{x}_\ell^\top \mathbf{p}_j\|^2}{\lambda_j} \quad (10)$$

**Result:** also reward large deviation along small variance.



**Soft outlier score:** normalize the residuals according to their corresponding variance along the  $d - k$  distances.

Decompose the sum of squares of the  $d - k$  distances and normalize by their corresponding eigenvalue:

$$\text{Score}(\mathbf{x}_\ell) = \sum_{j=k+1}^d \frac{\|\mathbf{x}_{\ell j} - \mathbf{x}_\ell^\top \mathbf{p}_j\|^2}{\lambda_j} \quad (10)$$

**Result:** also reward large deviation along small variance.

**Remark:** both scores focus on representing data in a low-dimensional space which induces parameter  $k$ : selecting the dimensionality



# PCA: Mahalanobis Connection

**Soft outlier score:** normalize the residuals according to their corresponding variance along the  $d - k$  distances.



# PCA: Mahalanobis Connection

**Soft outlier score:** normalize the residuals according to their corresponding variance along the  $d - k$  distances.

Mahalanobis performs the normalization across all  $d$  directions

$$\text{Score}(\mathbf{x}_\ell) = \sum_{j=1}^d \frac{\left\| (\mathbf{x}_\ell - \boldsymbol{\mu})^\top \mathbf{p}_j \right\|^2}{\lambda_j} \quad (11)$$

where  $\boldsymbol{\mu} \in \mathbb{R}^d$  is the data centroid (the mean vector along the data dimension).



# PCA: Mahalanobis Connection

**Soft outlier score:** normalize the residuals according to their corresponding variance along the  $d - k$  distances.

Mahalanobis performs the normalization across all  $d$  directions

$$\text{Score}(\mathbf{x}_\ell) = \sum_{j=1}^d \frac{\|(\mathbf{x}_\ell - \boldsymbol{\mu})^\top \mathbf{p}_j\|^2}{\lambda_j} \quad (11)$$

where  $\boldsymbol{\mu} \in \mathbb{R}^d$  is the data centroid (the mean vector along the data dimension).

**Algorithm:**



# PCA: Mahalanobis Connection

**Soft outlier score:** normalize the residuals according to their corresponding variance along the  $d - k$  distances.

Mahalanobis performs the normalization across all  $d$  directions

$$\text{Score}(\mathbf{x}_\ell) = \sum_{j=1}^d \frac{\|(\mathbf{x}_\ell - \boldsymbol{\mu})^\top \mathbf{p}_j\|^2}{\lambda_j} \quad (11)$$

where  $\boldsymbol{\mu} \in \mathbb{R}^d$  is the data centroid (the mean vector along the data dimension).

**Algorithm:**

1. EVD:  $\boldsymbol{\Sigma} = \mathbf{P}\boldsymbol{\Delta}\mathbf{P}^\top$



# PCA: Mahalanobis Connection

**Soft outlier score:** normalize the residuals according to their corresponding variance along the  $d - k$  distances.

Mahalanobis performs the normalization across all  $d$  directions

$$\text{Score}(\mathbf{x}_\ell) = \sum_{j=1}^d \frac{\|(\mathbf{x}_\ell - \boldsymbol{\mu})^\top \mathbf{p}_j\|^2}{\lambda_j} \quad (11)$$

where  $\boldsymbol{\mu} \in \mathbb{R}^d$  is the data centroid (the mean vector along the data dimension).

**Algorithm:**

1. EVD:  $\boldsymbol{\Sigma} = \mathbf{P}\boldsymbol{\Delta}\mathbf{P}^\top$
2. Transform:  $\mathbf{X}' = \mathbf{X}\mathbf{P}$





# PCA: Mahalanobis Connection

**Soft outlier score:** normalize the residuals according to their corresponding variance along the  $d - k$  distances.

Mahalanobis performs the normalization across all  $d$  directions

$$\text{Score}(\mathbf{x}_\ell) = \sum_{j=1}^d \frac{\|(\mathbf{x}_\ell - \boldsymbol{\mu})^\top \mathbf{p}_j\|^2}{\lambda_j} \quad (11)$$

where  $\boldsymbol{\mu} \in \mathbb{R}^d$  is the data centroid (the mean vector along the data dimension).

**Algorithm:**

1. EVD:  $\boldsymbol{\Sigma} = \mathbf{P}\boldsymbol{\Delta}\mathbf{P}^\top$
2. Transform:  $\mathbf{X}' = \mathbf{X}\mathbf{P}$
3. Normalize:  $\mathbf{X}' = \mathbf{X}'\boldsymbol{\Delta}^{-1}$



# PCA: Mahalanobis Connection

**Soft outlier score:** normalize the residuals according to their corresponding variance along the  $d - k$  distances.

Mahalanobis performs the normalization across all  $d$  directions

$$\text{Score}(\mathbf{x}_\ell) = \sum_{j=1}^d \frac{\|(\mathbf{x}_\ell - \boldsymbol{\mu})^\top \mathbf{p}_j\|^2}{\lambda_j} \quad (11)$$

where  $\boldsymbol{\mu} \in \mathbb{R}^d$  is the data centroid (the mean vector along the data dimension).

**Algorithm:**

1. EVD:  $\boldsymbol{\Sigma} = \mathbf{P}\boldsymbol{\Delta}\mathbf{P}^\top$
2. Transform:  $\mathbf{X}' = \mathbf{X}\mathbf{P}$
3. Normalize:  $\mathbf{X}' = \mathbf{X}'\boldsymbol{\Delta}^{-1}$
4. Anomaly score:  $\text{Score}(\mathbf{x}'_\ell) \quad \forall \ell \in [N]$



# Robust PCA



The lack of data covariance in the eigenspace adds robustness to noise.



# PCA Perturbations

The lack of data covariance in the eigenspace adds robustness to noise.

Robust PCA through iterative pruning:



The lack of data covariance in the eigenspace adds robustness to noise.

Robust PCA through iterative pruning:

1. PCA: compute  $\Sigma = \mathbf{P}\Delta\mathbf{P}^\top$
2. Anomaly score: compute associated scores  $\text{Score}(\mathbf{x}_\ell) \quad \forall \ell \in [N]$



The lack of data covariance in the eigenspace adds robustness to noise.

Robust PCA through iterative pruning:

1. PCA: compute  $\Sigma = \mathbf{P}\Delta\mathbf{P}^\top$
2. Anomaly score: compute associated scores  $\text{Score}(\mathbf{x}_\ell) \quad \forall \ell \in [N]$
3. Prune: remove obvious outliers from the dataset



The lack of data covariance in the eigenspace adds robustness to noise.

Robust PCA through iterative pruning:

1. PCA: compute  $\Sigma = \mathbf{P}\Delta\mathbf{P}^\top$
2. Anomaly score: compute associated scores  $\text{Score}(\mathbf{x}_\ell) \quad \forall \ell \in [N]$
3. Prune: remove obvious outliers from the dataset
4. Reconstruct: compute new covariance matrix





The lack of data covariance in the eigenspace adds robustness to noise.

Robust PCA through iterative pruning:

1. PCA: compute  $\Sigma = \mathbf{P}\Delta\mathbf{P}^\top$
2. Anomaly score: compute associated scores  $\text{Score}(\mathbf{x}_\ell) \quad \forall \ell \in [N]$
3. Prune: remove obvious outliers from the dataset
4. Reconstruct: compute new covariance matrix
5. Goto step 1



# Example: Outlier Perturbation

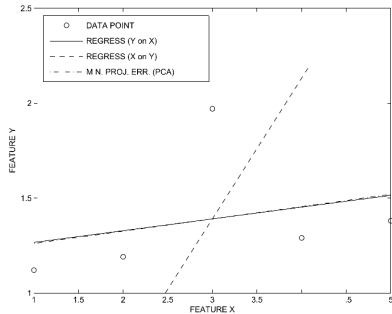
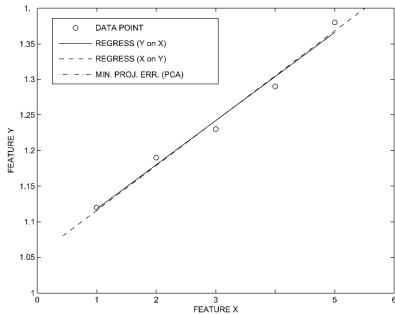


Figure: Sensitivity to outliers (Aggarwal 2017)



**Normalization:** original dimensions scales can vary widely – normalize to unit variance.

**Regularization:**

- ▶ zero variance among some dimensions implies  $\lambda = 0$



**Normalization:** original dimensions scales can vary widely – normalize to unit variance.

**Regularization:**

- ▶ zero variance among some dimensions implies  $\lambda = 0$
- ▶ regularize the covariance matrix  $\Sigma + \alpha I_d$  with  $\alpha > 0$



**Normalization:** original dimensions scales can vary widely – normalize to unit variance.

**Regularization:**

- ▶ zero variance among some dimensions implies  $\lambda = 0$
- ▶ regularize the covariance matrix  $\Sigma + \alpha I_d$  with  $\alpha > 0$
- ▶ shifts all eigenvalues by the constant value  $\alpha$



**Normalization:** original dimensions scales can vary widely – normalize to unit variance.

**Regularization:**

- ▶ zero variance among some dimensions implies  $\lambda = 0$
- ▶ regularize the covariance matrix  $\Sigma + \alpha I_d$  with  $\alpha > 0$
- ▶ shifts all eigenvalues by the constant value  $\alpha$
- ▶ equivalent to adding noise with variance  $\alpha$



**Normalization:** original dimensions scales can vary widely – normalize to unit variance.

**Regularization:**

- ▶ zero variance among some dimensions implies  $\lambda = 0$
- ▶ regularize the covariance matrix  $\Sigma + \alpha I_d$  with  $\alpha > 0$
- ▶ shifts all eigenvalues by the constant value  $\alpha$
- ▶ equivalent to adding noise with variance  $\alpha$

**Score through cross-validation:**

- ▶ split into  $m$ -folds



**Normalization:** original dimensions scales can vary widely – normalize to unit variance.

**Regularization:**

- ▶ zero variance among some dimensions implies  $\lambda = 0$
- ▶ regularize the covariance matrix  $\Sigma + \alpha I_d$  with  $\alpha > 0$
- ▶ shifts all eigenvalues by the constant value  $\alpha$
- ▶ equivalent to adding noise with variance  $\alpha$

**Score through cross-validation:**

- ▶ split into  $m$ -folds
- ▶ perform PCA on  $m - 1$  folds





**Normalization:** original dimensions scales can vary widely – normalize to unit variance.

**Regularization:**

- ▶ zero variance among some dimensions implies  $\lambda = 0$
- ▶ regularize the covariance matrix  $\Sigma + \alpha I_d$  with  $\alpha > 0$
- ▶ shifts all eigenvalues by the constant value  $\alpha$
- ▶ equivalent to adding noise with variance  $\alpha$

**Score through cross-validation:**

- ▶ split into  $m$ -folds
- ▶ perform PCA on  $m - 1$  folds
- ▶ score the data in the  $m^{\text{th}}$ -fold



**Normalization:** original dimensions scales can vary widely – normalize to unit variance.

**Regularization:**

- ▶ zero variance among some dimensions implies  $\lambda = 0$
- ▶ regularize the covariance matrix  $\Sigma + \alpha I_d$  with  $\alpha > 0$
- ▶ shifts all eigenvalues by the constant value  $\alpha$
- ▶ equivalent to adding noise with variance  $\alpha$

**Score through cross-validation:**

- ▶ split into  $m$ -folds
- ▶ perform PCA on  $m - 1$  folds
- ▶ score the data in the  $m^{\text{th}}$ -fold
- ▶ repeat for each fold



**Normalization:** original dimensions scales can vary widely – normalize to unit variance.

**Regularization:**

- ▶ zero variance among some dimensions implies  $\lambda = 0$
- ▶ regularize the covariance matrix  $\Sigma + \alpha I_d$  with  $\alpha > 0$
- ▶ shifts all eigenvalues by the constant value  $\alpha$
- ▶ equivalent to adding noise with variance  $\alpha$

**Score through cross-validation:**

- ▶ split into  $m$ -folds
- ▶ perform PCA on  $m - 1$  folds
- ▶ score the data in the  $m^{\text{th}}$ -fold
- ▶ repeat for each fold
- ▶ alternative: use sub-sampling



# Robust PCA

Treat measurement matrix as the super-position of a low-rank matrix with a sparse noise matrix  $X = L_0 + S_0$ , then recovering  $L_0$  and  $S_0$  involves solving the following optimization problem

$$\arg \min_{L, S} \rho(L) + \lambda \|S\|_0 \quad \text{s.t.} \quad \|X - L - S\|_F^2 = 0 \quad (12)$$

where  $\rho(L)$  is the rank function and  $\|\cdot\|_0$  is the  $\ell_0$ -norm counting the number of non-zeros.



# Robust PCA

Treat measurement matrix as the super-position of a low-rank matrix with a sparse noise matrix  $X = L_0 + S_0$ , then recovering  $L_0$  and  $S_0$  involves solving the following optimization problem

$$\arg \min_{L, S} \rho(L) + \lambda \|S\|_0 \quad \text{s.t.} \quad \|X - L - S\|_F^2 = 0 \quad (12)$$

where  $\rho(L)$  is the rank function and  $\|\cdot\|_0$  is the  $\ell_0$ -norm counting the number of non-zeros.

Candès et al. 2011 show that the convex relaxation of the above can recover  $L_0$  and  $S_0$  under mild assumptions

$$\arg \min_{L, S} \|L\|_{\star} + \lambda \|S\|_1 \quad \text{s.t.} \quad \|X - L - S\|_F^2 = 0 \quad (13)$$

where  $\|\cdot\|_{\star}$  is the nuclear norm summing the singular values.



# Escalator Example: PCA versus Robust PCA



**Figure:** Background separation: truth, PCA and two RPCA implementations (Netrapalli et al. 2014)



# Restaurant Example: PCA versus Robust PCA



**Figure:** Background separation: truth, PCA and two RPCA implementations (Netrapalli et al. 2014)



# Nonlinear PCA





# PCA: sample space versus feature space

What if we used  $\mathbf{S} = \mathbf{X}\mathbf{X}^T \in \mathbb{R}^{N \times N}$  instead of  $\mathbf{\Sigma}$ ?



# PCA: sample space versus feature space

What if we used  $\mathbf{S} = \mathbf{X}\mathbf{X}^\top \in \mathbb{R}^{N \times N}$  instead of  $\mathbf{\Sigma}$ ?

The EVD of  $\mathbf{S}$  becomes:

$$\mathbf{S} = \mathbf{X}\mathbf{X}^\top = \mathbf{Q}\mathbf{\Lambda}^2\mathbf{Q}^\top \quad (14)$$

where  $\mathbf{Q} \in \mathbb{R}^{N \times N}$  are the orthonormal eigenvectors such that only the first  $d$  correspond to non-zero eigenvalues.



# PCA: sample space versus feature space

What if we used  $\mathbf{S} = \mathbf{X}\mathbf{X}^\top \in \mathbb{R}^{N \times N}$  instead of  $\mathbf{\Sigma}$ ?

The EVD of  $\mathbf{S}$  becomes:

$$\mathbf{S} = \mathbf{X}\mathbf{X}^\top = \mathbf{Q}\mathbf{\Lambda}^2\mathbf{Q}^\top \quad (14)$$

where  $\mathbf{Q} \in \mathbb{R}^{N \times N}$  are the orthonormal eigenvectors such that only the first  $d$  correspond to non-zero eigenvalues.

The transform in the sample space is:

$$\mathbf{X}' = \mathbf{X}(\mathbf{Q}\mathbf{\Lambda})_d \quad (15)$$

where we can easily see that  $[\mathbf{X}' \ \mathbf{O}] = \mathbf{X}\mathbf{Q}\mathbf{\Lambda} = [\mathbf{X}(\mathbf{Q}\mathbf{\Lambda})_d \ \mathbf{O}]$ .



## PCA: sample space uses a similarity matrix

**Remark:**  $\mathbf{S} = \mathbf{X}\mathbf{X}^\top \in \mathbb{R}^{N \times N}$  has  $s_{ij}$  as the dot-product between  $\mathbf{x}_i$  and  $\mathbf{x}_j$  thus acting as a similarity matrix.



## PCA: sample space uses a similarity matrix

**Remark:**  $\mathbf{S} = \mathbf{X}\mathbf{X}^\top \in \mathbb{R}^{N \times N}$  has  $s_{ij}$  as the dot-product between  $\mathbf{x}_i$  and  $\mathbf{x}_j$  thus acting as a similarity matrix.

**Remark:** the transform  $\mathbf{S} = \mathbf{Q}\mathbf{\Lambda}^2\mathbf{Q}^\top$  is actually the SVD transform.



# PCA: sample space uses a similarity matrix

**Remark:**  $\mathbf{S} = \mathbf{X}\mathbf{X}^\top \in \mathbb{R}^{N \times N}$  has  $s_{ij}$  as the dot-product between  $\mathbf{x}_i$  and  $\mathbf{x}_j$  thus acting as a similarity matrix.

**Remark:** the transform  $\mathbf{S} = \mathbf{Q}\mathbf{\Lambda}^2\mathbf{Q}^\top$  is actually the SVD transform.

**Remark:** We can use a different similarity matrix instead of the dot-product; we can use a kernel function and employ the kernel trick!



# PCA: sample space uses a similarity matrix

**Remark:**  $\mathbf{S} = \mathbf{X}\mathbf{X}^\top \in \mathbb{R}^{N \times N}$  has  $s_{ij}$  as the dot-product between  $\mathbf{x}_i$  and  $\mathbf{x}_j$  thus acting as a similarity matrix.

**Remark:** the transform  $\mathbf{S} = \mathbf{Q}\mathbf{\Lambda}^2\mathbf{Q}^\top$  is actually the SVD transform.

**Remark:** We can use a different similarity matrix instead of the dot-product; we can use a kernel function and employ the kernel trick!

**Algorithm:**



# PCA: sample space uses a similarity matrix

**Remark:**  $\mathbf{S} = \mathbf{X}\mathbf{X}^\top \in \mathbb{R}^{N \times N}$  has  $s_{ij}$  as the dot-product between  $\mathbf{x}_i$  and  $\mathbf{x}_j$  thus acting as a similarity matrix.

**Remark:** the transform  $\mathbf{S} = \mathbf{Q}\mathbf{\Lambda}^2\mathbf{Q}^\top$  is actually the SVD transform.

**Remark:** We can use a different similarity matrix instead of the dot-product; we can use a kernel function and employ the kernel trick!

## Algorithm:

1. Similarity: Build  $N \times N$  similarity matrix  $\mathbf{S}$  using kernel function  $\kappa$





# PCA: sample space uses a similarity matrix

**Remark:**  $\mathbf{S} = \mathbf{X}\mathbf{X}^\top \in \mathbb{R}^{N \times N}$  has  $s_{ij}$  as the dot-product between  $\mathbf{x}_i$  and  $\mathbf{x}_j$  thus acting as a similarity matrix.

**Remark:** the transform  $\mathbf{S} = \mathbf{Q}\mathbf{\Lambda}^2\mathbf{Q}^\top$  is actually the SVD transform.

**Remark:** We can use a different similarity matrix instead of the dot-product; we can use a kernel function and employ the kernel trick!

## Algorithm:

1. Similarity: Build  $N \times N$  similarity matrix  $\mathbf{S}$  using kernel function  $\kappa$
2. Decomposition:  $\mathbf{S} = \mathbf{Q}\mathbf{\Lambda}^2\mathbf{Q}^\top$



# PCA: sample space uses a similarity matrix

**Remark:**  $\mathbf{S} = \mathbf{X}\mathbf{X}^\top \in \mathbb{R}^{N \times N}$  has  $s_{ij}$  as the dot-product between  $\mathbf{x}_i$  and  $\mathbf{x}_j$  thus acting as a similarity matrix.

**Remark:** the transform  $\mathbf{S} = \mathbf{Q}\mathbf{\Lambda}^2\mathbf{Q}^\top$  is actually the SVD transform.

**Remark:** We can use a different similarity matrix instead of the dot-product; we can use a kernel function and employ the kernel trick!

## Algorithm:

1. Similarity: Build  $N \times N$  similarity matrix  $\mathbf{S}$  using kernel function  $\kappa$
2. Decomposition:  $\mathbf{S} = \mathbf{Q}\mathbf{\Lambda}^2\mathbf{Q}^\top$
3. Transform:  $\mathbf{X}' = \mathbf{X}(\mathbf{Q}\mathbf{\Lambda})_k$  ( $N > k > d$  for some non-linear functions)



# PCA: sample space uses a similarity matrix

**Remark:**  $\mathbf{S} = \mathbf{X}\mathbf{X}^\top \in \mathbb{R}^{N \times N}$  has  $s_{ij}$  as the dot-product between  $\mathbf{x}_i$  and  $\mathbf{x}_j$  thus acting as a similarity matrix.

**Remark:** the transform  $\mathbf{S} = \mathbf{Q}\mathbf{\Lambda}^2\mathbf{Q}^\top$  is actually the SVD transform.

**Remark:** We can use a different similarity matrix instead of the dot-product; we can use a kernel function and employ the kernel trick!

## Algorithm:

1. Similarity: Build  $N \times N$  similarity matrix  $\mathbf{S}$  using kernel function  $\kappa$
2. Decomposition:  $\mathbf{S} = \mathbf{Q}\mathbf{\Lambda}^2\mathbf{Q}^\top$
3. Transform:  $\mathbf{X}' = \mathbf{X}(\mathbf{Q}\mathbf{\Lambda})_k$  ( $N > k > d$  for some non-linear functions)
4. Normalize:  $\mathbf{X}' = \mathbf{X}'\mathbf{\Lambda}^{-1}$



# PCA: sample space uses a similarity matrix

**Remark:**  $\mathbf{S} = \mathbf{X}\mathbf{X}^\top \in \mathbb{R}^{N \times N}$  has  $s_{ij}$  as the dot-product between  $\mathbf{x}_i$  and  $\mathbf{x}_j$  thus acting as a similarity matrix.

**Remark:** the transform  $\mathbf{S} = \mathbf{Q}\mathbf{\Lambda}^2\mathbf{Q}^\top$  is actually the SVD transform.

**Remark:** We can use a different similarity matrix instead of the dot-product; we can use a kernel function and employ the kernel trick!

## Algorithm:

1. Similarity: Build  $N \times N$  similarity matrix  $\mathbf{S}$  using kernel function  $\kappa$
2. Decomposition:  $\mathbf{S} = \mathbf{Q}\mathbf{\Lambda}^2\mathbf{Q}^\top$
3. Transform:  $\mathbf{X}' = \mathbf{X}(\mathbf{Q}\mathbf{\Lambda})_k$  ( $N > k > d$  for some non-linear functions)
4. Normalize:  $\mathbf{X}' = \mathbf{X}'\mathbf{\Lambda}^{-1}$
5. Anomaly score:  $\text{Score}(\mathbf{x}'_\ell) \quad \forall \ell \in [N]$



# Example: Kernel Space

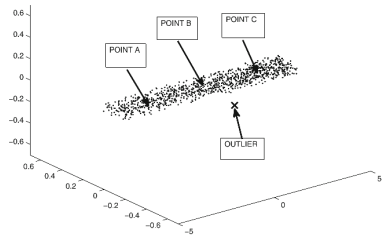
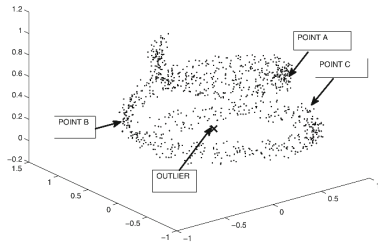


Figure: Sample space to kernel space (Aggarwal 2017)



# Nonlinear PCA: Dealing with large sample spaces

**Remark:** Computing  $\mathbf{S} = \mathbf{Q}\mathbf{\Lambda}^2\mathbf{Q}^\top \in \mathbb{R}^{N \times N}$  for large  $N$  can be prohibitive (e.g.  $N = 100,000$ ).



# Nonlinear PCA: Dealing with large sample spaces

**Remark:** Computing  $\mathbf{S} = \mathbf{Q}\mathbf{\Lambda}^2\mathbf{Q}^\top \in \mathbb{R}^{N \times N}$  for large  $N$  can be prohibitive (e.g.  $N = 100,000$ ).

**Solution:** One can apply sub-sampling using  $s$  samples instead of  $N$  such that  $k < s < N$  thus obtaining an  $s \times s$  similarity matrix.



# Nonlinear PCA: Dealing with large sample spaces

**Remark:** Computing  $\mathbf{S} = \mathbf{Q}\mathbf{\Lambda}^2\mathbf{Q}^\top \in \mathbb{R}^{N \times N}$  for large  $N$  can be prohibitive (e.g.  $N = 100,000$ ).

**Solution:** One can apply sub-sampling using  $s$  samples instead of  $N$  such that  $k < s < N$  thus obtaining an  $s \times s$  similarity matrix.

**Algorithm:**





# Nonlinear PCA: Dealing with large sample spaces

**Remark:** Computing  $\mathbf{S} = \mathbf{Q}\mathbf{\Lambda}^2\mathbf{Q}^\top \in \mathbb{R}^{N \times N}$  for large  $N$  can be prohibitive (e.g.  $N = 100,000$ ).

**Solution:** One can apply sub-sampling using  $s$  samples instead of  $N$  such that  $k < s < N$  thus obtaining an  $s \times s$  similarity matrix.

**Algorithm:**

1. Similarity: Build  $s \times s$  similarity matrix  $\mathbf{S}$  using kernel function  $\kappa$



# Nonlinear PCA: Dealing with large sample spaces

**Remark:** Computing  $\mathbf{S} = \mathbf{Q}\mathbf{\Lambda}^2\mathbf{Q}^\top \in \mathbb{R}^{N \times N}$  for large  $N$  can be prohibitive (e.g.  $N = 100,000$ ).

**Solution:** One can apply sub-sampling using  $s$  samples instead of  $N$  such that  $k < s < N$  thus obtaining an  $s \times s$  similarity matrix.

**Algorithm:**

1. Similarity: Build  $s \times s$  similarity matrix  $\mathbf{S}$  using kernel function  $\kappa$
2. Decomposition:  $\mathbf{S} = \mathbf{Q}\mathbf{\Lambda}^2\mathbf{Q}^\top$



# Nonlinear PCA: Dealing with large sample spaces

**Remark:** Computing  $\mathbf{S} = \mathbf{Q}\mathbf{\Lambda}^2\mathbf{Q}^\top \in \mathbb{R}^{N \times N}$  for large  $N$  can be prohibitive (e.g.  $N = 100,000$ ).

**Solution:** One can apply sub-sampling using  $s$  samples instead of  $N$  such that  $k < s < N$  thus obtaining an  $s \times s$  similarity matrix.

**Algorithm:**

1. Similarity: Build  $s \times s$  similarity matrix  $\mathbf{S}$  using kernel function  $\kappa$
2. Decomposition:  $\mathbf{S} = \mathbf{Q}\mathbf{\Lambda}^2\mathbf{Q}^\top$
3. Transform:  $\mathbf{X}' = \mathbf{X}(\mathbf{Q}\mathbf{\Lambda})_k$  ( $N > k > d$  for some non-linear functions)



# Nonlinear PCA: Dealing with large sample spaces

**Remark:** Computing  $\mathbf{S} = \mathbf{Q}\mathbf{\Lambda}^2\mathbf{Q}^\top \in \mathbb{R}^{N \times N}$  for large  $N$  can be prohibitive (e.g.  $N = 100,000$ ).

**Solution:** One can apply sub-sampling using  $s$  samples instead of  $N$  such that  $k < s < N$  thus obtaining an  $s \times s$  similarity matrix.

## Algorithm:

1. Similarity: Build  $s \times s$  similarity matrix  $\mathbf{S}$  using kernel function  $\kappa$
2. Decomposition:  $\mathbf{S} = \mathbf{Q}\mathbf{\Lambda}^2\mathbf{Q}^\top$
3. Transform:  $\mathbf{X}' = \mathbf{X}(\mathbf{Q}\mathbf{\Lambda})_k$  ( $N > k > d$  for some non-linear functions)
4. Normalize:  $\mathbf{X}' = \mathbf{X}'\mathbf{\Lambda}^{-1}$



# Nonlinear PCA: Dealing with large sample spaces

**Remark:** Computing  $\mathbf{S} = \mathbf{Q}\mathbf{\Lambda}^2\mathbf{Q}^\top \in \mathbb{R}^{N \times N}$  for large  $N$  can be prohibitive (e.g.  $N = 100,000$ ).

**Solution:** One can apply sub-sampling using  $s$  samples instead of  $N$  such that  $k < s < N$  thus obtaining an  $s \times s$  similarity matrix.

## Algorithm:

1. Similarity: Build  $s \times s$  similarity matrix  $\mathbf{S}$  using kernel function  $\kappa$
2. Decomposition:  $\mathbf{S} = \mathbf{Q}\mathbf{\Lambda}^2\mathbf{Q}^\top$
3. Transform:  $\mathbf{X}' = \mathbf{X}(\mathbf{Q}\mathbf{\Lambda})_k$  ( $N > k > d$  for some non-linear functions)
4. Normalize:  $\mathbf{X}' = \mathbf{X}'\mathbf{\Lambda}^{-1}$
5. Anomaly score:  $\text{Score}(\mathbf{x}'_\ell) \quad \forall \ell \in [N]$



# Nonlinear PCA: Dealing with large sample spaces

**Remark:** Computing  $\mathbf{S} = \mathbf{Q}\mathbf{\Lambda}^2\mathbf{Q}^\top \in \mathbb{R}^{N \times N}$  for large  $N$  can be prohibitive (e.g.  $N = 100,000$ ).

**Solution:** One can apply sub-sampling using  $s$  samples instead of  $N$  such that  $k < s < N$  thus obtaining an  $s \times s$  similarity matrix.

**Algorithm:**

1. Similarity: Build  $s \times s$  similarity matrix  $\mathbf{S}$  using kernel function  $\kappa$
2. Decomposition:  $\mathbf{S} = \mathbf{Q}\mathbf{\Lambda}^2\mathbf{Q}^\top$
3. Transform:  $\mathbf{X}' = \mathbf{X}(\mathbf{Q}\mathbf{\Lambda})_k$  ( $N > k > d$  for some non-linear functions)
4. Normalize:  $\mathbf{X}' = \mathbf{X}'\mathbf{\Lambda}^{-1}$
5. Anomaly score:  $\text{Score}(\mathbf{x}'_\ell) \quad \forall \ell \in [N]$
6. Out-of-sample similarity: Build  $(N - s) \times s$  similarity matrix  $\mathbf{S}_0$  using the same kernel



# Nonlinear PCA: Dealing with large sample spaces

**Remark:** Computing  $\mathbf{S} = \mathbf{Q}\mathbf{\Lambda}^2\mathbf{Q}^\top \in \mathbb{R}^{N \times N}$  for large  $N$  can be prohibitive (e.g.  $N = 100,000$ ).

**Solution:** One can apply sub-sampling using  $s$  samples instead of  $N$  such that  $k < s < N$  thus obtaining an  $s \times s$  similarity matrix.

**Algorithm:**

1. Similarity: Build  $s \times s$  similarity matrix  $\mathbf{S}$  using kernel function  $\kappa$
2. Decomposition:  $\mathbf{S} = \mathbf{Q}\mathbf{\Lambda}^2\mathbf{Q}^\top$
3. Transform:  $\mathbf{X}' = \mathbf{X}(\mathbf{Q}\mathbf{\Lambda})_k$  ( $N > k > d$  for some non-linear functions)
4. Normalize:  $\mathbf{X}' = \mathbf{X}'\mathbf{\Lambda}^{-1}$
5. Anomaly score:  $\text{Score}(\mathbf{x}'_\ell) \quad \forall \ell \in [N]$
6. Out-of-sample similarity: Build  $(N - s) \times s$  similarity matrix  $\mathbf{S}_0$  using the same kernel
7. Out-of-sample transform:  $\mathbf{X}'_0 = \mathbf{S}_0(\mathbf{Q}\mathbf{\Lambda}^{-1})_k$



# Nonlinear PCA: Dealing with large sample spaces

**Remark:** Computing  $\mathbf{S} = \mathbf{Q}\mathbf{\Lambda}^2\mathbf{Q}^\top \in \mathbb{R}^{N \times N}$  for large  $N$  can be prohibitive (e.g.  $N = 100,000$ ).

**Solution:** One can apply sub-sampling using  $s$  samples instead of  $N$  such that  $k < s < N$  thus obtaining an  $s \times s$  similarity matrix.

## Algorithm:

1. Similarity: Build  $s \times s$  similarity matrix  $\mathbf{S}$  using kernel function  $\kappa$
2. Decomposition:  $\mathbf{S} = \mathbf{Q}\mathbf{\Lambda}^2\mathbf{Q}^\top$
3. Transform:  $\mathbf{X}' = \mathbf{X}(\mathbf{Q}\mathbf{\Lambda})_k$  ( $N > k > d$  for some non-linear functions)
4. Normalize:  $\mathbf{X}' = \mathbf{X}'\mathbf{\Lambda}^{-1}$
5. Anomaly score:  $\text{Score}(\mathbf{x}'_\ell) \quad \forall \ell \in [N]$
6. Out-of-sample similarity: Build  $(N - s) \times s$  similarity matrix  $\mathbf{S}_0$  using the same kernel
7. Out-of-sample transform:  $\mathbf{X}'_0 = \mathbf{S}_0(\mathbf{Q}\mathbf{\Lambda}^{-1})_k$
8. Complete transform:  $[(\mathbf{Q}\mathbf{\Lambda})_k ; \mathbf{S}_0(\mathbf{Q}\mathbf{\Lambda}^{-1})_k]^\top$





# References

- Aggarwal, Charu C (2017). *An introduction to outlier analysis*. Springer.
- Candès, Emmanuel J et al. (2011). "Robust principal component analysis?" In: *Journal of the ACM (JACM)* 58.3, pp. 1–37.
- Golub, G.H. and C.F. Van Loan (2013). *Matrix Computations*. John Hopkins University Press.
- Netrapalli, Praneeth et al. (2014). "Non-convex robust PCA". In: *Advances in neural information processing systems* 27.

