

# Statistical Machine Learning

## Lecture 6: Singular Value Decomposition (SVD) and Principal Component Analysis (PCA)

2022-23

# SVD

- ▶ Common data pre-processing technique to reduce high-dimensional data.
- ▶ Provides a robust justification for matrix (data arrays) approximation (reduction).
- ▶ Provides a systematic way to determine a low-dimensional approximation to a high-dimensional data in terms of dominant patterns

# SVD Theorem

- *Theorem:* Let  $\mathbf{A} \in \mathbb{R}^{m \times n}$  be matrix of rank  $r \in \{0, 1, \dots, \min\{m, n\}\}$ . There exists an  $m \times m$  real orthogonal matrix  $\mathbf{U}$  and an  $n \times n$  real orthogonal matrix  $\mathbf{V}$  such that:

$$\mathbf{A} = \mathbf{U} \mathbf{D} \mathbf{V}^T,$$

where

$$\mathbf{D} = \begin{pmatrix} \mathbf{\Sigma} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} \quad \mathbf{\Sigma} = \begin{pmatrix} \sigma_1 & 0 & \dots & 0 \\ 0 & \sigma_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & 0 & \sigma_r \end{pmatrix},$$

where  $\mathbf{D} \in \mathbb{R}^{m \times n}$ ,  $\mathbf{\Sigma} \in \mathbb{R}^{r \times r}$  and diagonal, and  $\sigma_i \in \mathbb{R}$  such that  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0$ .

# SVD Theorem

- The decomposition is also expressed using the following partition:

$$\mathbf{A} = \mathbf{U}\mathbf{D}\mathbf{V}^T = (\mathbf{U}_1, \mathbf{U}_2) \begin{pmatrix} \mathbf{\Sigma} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} \begin{pmatrix} \mathbf{V}_1^T \\ \mathbf{V}_2^T \end{pmatrix} = \mathbf{U}_1 \mathbf{\Sigma} \mathbf{V}_1^T,$$

where  $\mathbf{U}_1$  and  $\mathbf{V}_1$  are  $m \times r$  and  $n \times r$  matrices, respectively, with orthonormal columns and the  $\mathbf{0}$  submatrices have compatible dimensions for the above partition to be sensible.

# SVD Theorem

- *Proof:* The matrix  $\mathbf{A}^T \mathbf{A}$  is an  $n \times n$  symmetric matrix. Those, its eigenvalues  $\lambda$  are all real. Also:

$$\begin{aligned}\mathbf{A}^T \mathbf{A} \mathbf{v}_i &= \lambda_i \mathbf{v}_i \\ \mathbf{v}_i^T \mathbf{A}^T \mathbf{A} \mathbf{v}_i &= \lambda_i \mathbf{v}_i^T \mathbf{v}_i \\ (\mathbf{A} \mathbf{v}_i)^T \mathbf{A} \mathbf{v}_i &= \lambda_i \|\mathbf{v}_i\|^2 \\ 0 \leq \frac{\|\mathbf{A} \mathbf{v}_i\|^2}{\|\mathbf{v}_i\|^2} &= \lambda_i.\end{aligned}$$

Let  $\{\lambda_1, \dots, \lambda_n\}$  denote the set of eigenvalues of  $\mathbf{A}^T \mathbf{A}$ .

# SVD Theorem

- The spectral decomposition of  $\mathbf{A}^T \mathbf{A}$  is then:

$$\mathbf{A}^T \mathbf{A} = \mathbf{V} \mathbf{\Lambda} \mathbf{V}^T,$$

where  $[\mathbf{\Lambda}]_{i,i} = \lambda_i$  and  $[\mathbf{\Lambda}]_{i,j} = 0$  for all  $i \neq j$ . Also  $\mathbf{V}$  is an  $n \times n$  real orthonormal matrix. We can assume, without loss of generality, that  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n \geq 0$ . Since the rank of the matrix  $\mathbf{A}$  is  $\rho(\mathbf{A}) = r$ , and by using the fundamental theorem of ranks  $\rho(\mathbf{A}) = \rho(\mathbf{A}) = \rho(\mathbf{A}\mathbf{A}^T) = \rho(\mathbf{A}^T \mathbf{A})$ . It follows that

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_r > \lambda_{r+1} = \dots = \lambda_n = 0.$$

# SVD Theorem

- ▶ The  $n$  columns of  $\mathbf{V} = (\mathbf{v}_1, \dots, \mathbf{v}_n)$  are the correspondent eigenvectors of  $\mathbf{A}^T \mathbf{A}$ . Partition  $\mathbf{V} = (\mathbf{V}_1, \mathbf{V}_2)$  such that  $\mathbf{V}_1 = (\mathbf{v}_1, \dots, \mathbf{v}_r)$  is  $n \times r$  and  $\mathbf{V}_1^T \mathbf{V}_1 = \mathbf{I}_{r \times r}$ . Let  $\mathbf{\Lambda}_1$  be the  $r \times r$  diagonal matrix with  $\lambda_1 \geq \dots \geq \lambda_r > 0$  elements in the diagonal. The spectral decomposition produces:

$$\begin{aligned}\mathbf{A}^T \mathbf{A} &= \mathbf{V} \mathbf{\Lambda} \mathbf{V}^T \\ \mathbf{A}^T \mathbf{A} \mathbf{V} &= \mathbf{V} \mathbf{\Lambda} \\ &= (\mathbf{V}_1, \mathbf{V}_2) \mathbf{\Lambda} \\ &= (\mathbf{V}_1 \mathbf{\Lambda}_1, \mathbf{0}).\end{aligned}$$

# SVD Theorem

► Thus

$$\begin{aligned}\mathbf{A}^T \mathbf{A} \mathbf{V}_2 &= \mathbf{0} \\ (\mathbf{A} \mathbf{V}_2)^T \mathbf{A} \mathbf{V}_2 &= \mathbf{0} \\ \mathbf{A} \mathbf{V}_2 &= \mathbf{0},\end{aligned}$$

thus each column of  $\mathbf{V}_2$  belongs to the null space of  $\mathbf{A}$ . Let us define  $\mathbf{\Sigma} = \mathbf{\Lambda}_1^{\frac{1}{2}}$ , and define the  $m \times r$  matrix  $\mathbf{U}_1$ :

$$\mathbf{U}_1 = \mathbf{A} \mathbf{V}_1 \mathbf{\Sigma}^{-1}.$$



# SVD Theorem

- Observe that from the spectral decomposition we also have that

$$\mathbf{A}^T \mathbf{A} \mathbf{V}_1 = \mathbf{V}_1 \mathbf{\Lambda}_1 = \mathbf{V}_1 \mathbf{\Sigma}^2$$

$$(\mathbf{A} \mathbf{V}_1)^T \mathbf{A} \mathbf{V}_1 = \mathbf{V}_1^T \mathbf{V}_1 \mathbf{\Sigma}^2 = \mathbf{\Sigma}^2$$

$$\mathbf{\Sigma}^{-1} (\mathbf{A} \mathbf{V}_1)^T \mathbf{A} \mathbf{V}_1 \mathbf{\Sigma}^{-1} = \mathbf{I}_{r \times r}$$

$$(\mathbf{A} \mathbf{V}_1 \mathbf{\Sigma}^{-1})^T \mathbf{A} \mathbf{V}_1 \mathbf{\Sigma}^{-1} = \mathbf{U}_1^T \mathbf{U}_1 = \mathbf{I}_{r \times r}.$$

Therefore the columns of  $\mathbf{U}_1$  are orthonormal.

# SVD Theorem

- We also observe that:

$$\begin{aligned}U_1 &= AV_1 \Sigma^{-1} \\ I_{r \times r} &= U_1^T U_1 = U_1^T AV_1 \Sigma^{-1} \\ \Sigma &= U_1^T AV_1 \\ U_1 \Sigma V_1^T &= A.\end{aligned}$$

Let us choose  $(m - r)$  vectors  $\mathbf{u}_k$ ,  $k = r + 1, \dots, m$ , such that these vectors are perpendicular to the columns of  $U_1$  and  $\mathbf{u}_k^T \mathbf{u}_{k'} = \delta_{k,k'}$ .

# SVD Theorem

► So

$$\begin{aligned}U_2 &= (\mathbf{u}_{r+1}, \dots, \mathbf{u}_m) \\U_2^T \mathbf{A} \mathbf{V}_1 &= U_2^T U_1 \boldsymbol{\Sigma} = \mathbf{0}.\end{aligned}$$

This implies that

$$\begin{aligned}U^T \mathbf{A} \mathbf{V} &= \begin{pmatrix} U_1^T \\ U_2^T \end{pmatrix} \mathbf{A} (\mathbf{V}_1, \mathbf{V}_2) \\&= \begin{pmatrix} U_1^T \mathbf{A} \mathbf{V}_1 & U_1^T \mathbf{A} \mathbf{V}_2 \\ U_2^T \mathbf{A} \mathbf{V}_1 & U_2^T \mathbf{A} \mathbf{V}_2 \end{pmatrix}\end{aligned}$$

# SVD Theorem



$$\begin{aligned} U^T AV &= \begin{pmatrix} U_1^T AV_1 & 0 \\ U_2^T AV_1 & 0 \end{pmatrix} & AV_2 &= 0 \\ &= \begin{pmatrix} U_1^T AV_1 & 0 \\ 0 & 0 \end{pmatrix}, \end{aligned}$$

thus

$$\begin{aligned} U^T AV &= \begin{pmatrix} \Sigma & 0 \\ 0 & 0 \end{pmatrix} = D \\ A &= UDV^T. \spadesuit \end{aligned}$$

# Approximation Theorem

- *Theorem:* The optimal rank- $r$  approximation to  $\mathbf{X}$ , in a least-squares sense, is given by the rank- $r$  SVD truncation  $\tilde{\mathbf{X}}$  :

$$\operatorname{argmin}_{\tilde{\mathbf{X}}: \rho(\tilde{\mathbf{X}})=r} \left\| \mathbf{X} - \tilde{\mathbf{X}} \right\|_F = \tilde{\mathbf{U}} \tilde{\mathbf{\Sigma}} \tilde{\mathbf{V}}^T,$$

where  $\tilde{\mathbf{U}}$  and  $\tilde{\mathbf{V}}$  denote the first  $r$  leading columns of the matrices  $\mathbf{U}$  and  $\mathbf{V}$ , and  $\tilde{\mathbf{\Sigma}}$  contains the leading  $r \times r$  sub-block of  $\mathbf{\Sigma}$ .  $\|\cdot\|_F$  is the Frobenius norm:

$$\forall \mathbf{M} \in \mathbb{R}^{n \times m}, \|\mathbf{M}\|_F = \sqrt{\sum_{i=1}^n \sum_{j=1}^m M_{i,j}^2}.$$

# Approximation Theorem

- ▶ Thus, if we consider  $r' < r$

$$\tilde{\mathbf{X}}' = \sum_{k=1}^{r'} \sigma_k \mathbf{u}_k \mathbf{v}_k^T$$

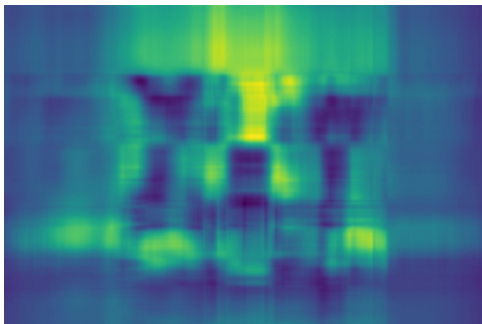
$\tilde{\mathbf{X}}'$  is the optimal approximation to  $\mathbf{X}$  with  $r'$  components.

# Approximation Theorem



# Approximation Theorem

►  $r = 5$





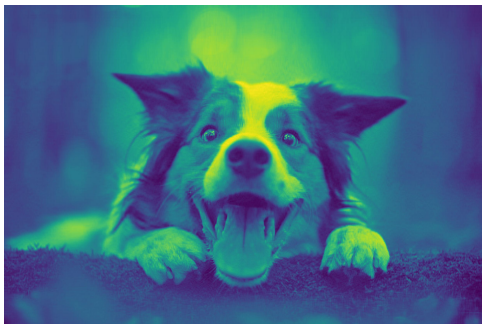
# Approximation Theorem

- ▶  $r = 20$



# Approximation Theorem

►  $r = 100$



# Systems of Linear Equations

- The *pseudo-inverse* of a matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$  is defined as the matrix  $\mathbf{A}^+ \in \mathbb{R}^{n \times m}$  such that if

$$\mathbf{A} = \mathbf{U}\mathbf{D}\mathbf{V}^T$$

then

$$\mathbf{A}^+ \equiv \mathbf{V}\mathbf{D}^+\mathbf{U}^T$$

with

$$\mathbf{D} = \begin{cases} \begin{pmatrix} \mathbf{\Sigma} & \mathbf{0}_{r \times (n-r)} \\ \mathbf{0}_{(m-r) \times r} & \mathbf{0}_{(m-r) \times (n-r)} \end{pmatrix} & r < \min\{m, n\} \\ \begin{pmatrix} \mathbf{\Sigma} \\ \mathbf{0}_{(m-r) \times r} \end{pmatrix} & n = r < m \\ \begin{pmatrix} \mathbf{\Sigma} & \mathbf{0}_{r \times (n-r)} \end{pmatrix} & m = r < n \\ \mathbf{\Sigma} & m = n = r \end{cases}$$

# Systems of Linear Equations

► Then

$$D^+ = \begin{cases} \begin{pmatrix} \Sigma^{-1} & \mathbf{0}_{r \times (m-r)} \\ \mathbf{0}_{(n-r) \times r} & \mathbf{0}_{(n-r) \times (m-r)} \end{pmatrix} & r < \min\{m, n\} \\ \begin{pmatrix} \Sigma^{-1} & \mathbf{0}_{r \times (m-r)} \end{pmatrix} & n = r < m \\ \begin{pmatrix} \Sigma^{-1} \\ \mathbf{0}_{(n-r) \times r} \end{pmatrix} & m = r < n \\ \Sigma^{-1} & m = n = r \end{cases}.$$

# Systems of Linear Equations

- ▶ Consider the system  $\mathbf{Ax} = \mathbf{b}$  for  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $\mathbf{b} \in \mathbb{R}^m$  and a suitable  $\mathbf{x} \in \mathbb{R}^n$ .
- ▶ Let us explore the properties of the vector  $\tilde{\mathbf{x}} \equiv \mathbf{A}^+\mathbf{b}$ .
- ▶ Let us consider the minimum norm defined as

$$\theta \equiv \min\{\|\mathbf{Ax} - \mathbf{b}\|_2 : \mathbf{x} \in \mathbb{R}^n\},$$

where

$$\|\mathbf{v}\|_2 = \sqrt{\mathbf{v}^T \mathbf{v}}$$

is the Euclidean norm.

# Systems of Linear Equations

- ▶ The Euclidean norm is indifferent under changes of coordinates. Given a unitary matrix  $\mathbf{P} \in \mathbb{R}^{m \times m}$  such that  $\mathbf{P}\mathbf{P}^T = \mathbf{P}^T\mathbf{P} = \mathbf{I}_{m \times m}$  we have that for all  $\mathbf{y} \in \mathbb{R}^m$   
 $\|\mathbf{y}\|_2 = \|\mathbf{P}^T\mathbf{y}\|_2$ .
- ▶ Then

$$\begin{aligned}\tilde{\theta} &\equiv \|\mathbf{A}\tilde{\mathbf{x}} - \mathbf{b}\|_2 \\ &= \|\mathbf{A}\mathbf{A}^+\mathbf{b} - \mathbf{b}\|_2 \\ &= \|\mathbf{U}\mathbf{D}\mathbf{V}^T\mathbf{V}\mathbf{D}^+\mathbf{U}^T\mathbf{b} - \mathbf{b}\|_2 \\ &= \|\mathbf{U}\mathbf{D}\mathbf{D}^+\mathbf{U}^T\mathbf{b} - \mathbf{b}\|_2 \\ &= \|\mathbf{D}\mathbf{D}^+\mathbf{U}^T\mathbf{b} - \mathbf{U}^T\mathbf{b}\|_2.\end{aligned}$$

# Systems of Linear Equations

- Observe that

$$DD^+ = \begin{cases} \begin{pmatrix} I_{r \times r} & \mathbf{0}_{r \times (m-r)} \\ \mathbf{0}_{(m-r) \times r} & \mathbf{0}_{(m-r) \times (m-r)} \end{pmatrix} & r < \min\{m, n\} \\ I_{r \times r} & m = r \end{cases} .$$

# Systems of Linear Equations

► Therefore

$$\begin{aligned}(I_{m \times m} - DD^+)(U^T b) &= \begin{pmatrix} \mathbf{0}_{r \times r} & \mathbf{0}_{r \times (m-r)} \\ \mathbf{0}_{(m-r) \times r} & I_{(m-r) \times (m-r)} \end{pmatrix} (U^T b) \\ &= \begin{pmatrix} 0 \\ \vdots \\ 0 \\ [U^T b]_r \\ \vdots \\ [U^T b]_m \end{pmatrix},\end{aligned}$$



# Systems of Linear Equations

- Then

$$\tilde{\theta} = \begin{cases} \sqrt{\sum_{\ell=r+1}^m [\mathbf{U}^T \mathbf{b}]_{\ell}^2} & r < \min\{m, n\} \\ 0 & r = m \end{cases}.$$

- If  $r = \rho(\mathbf{A}) = m$  and therefore  $m \leq n$ , then  $\tilde{\theta} = 0$  which must be the minimum  $\theta$ . Therefore  $\mathbf{A}^+ \mathbf{b}$  is a solution to the linear system (with minimum norm).

# Systems of Linear Equations

- ▶ If  $r = \rho(\mathbf{A}) = m$  and therefore  $m \leq n$ , then  $\tilde{\theta} = 0$  which must be the minimum  $\theta$ . Therefore  $\mathbf{A}^+ \mathbf{b}$  is a solution to the linear system (with minimum norm).
- ▶ If  $r < \min\{m, n\}$  we have that  $\rho(\mathbf{A}) < \dim(\mathbf{b})$  and the linear system may be undetermined (therefore no solution exists). In such a case observe that for any  $\mathbf{z} = \tilde{\mathbf{x}} + \mathbf{v}$  we have that

$$\begin{aligned}\theta(\mathbf{z}) &= \|\mathbf{A}(\tilde{\mathbf{x}} + \mathbf{v}) - \mathbf{b}\|_2 \\ &= \left\| (\mathbf{D}\mathbf{D}^+ - \mathbf{I}_{m \times m}) \mathbf{U}^T \mathbf{b} + \mathbf{U}^T \mathbf{A} \mathbf{v} \right\|_2\end{aligned}$$

# Systems of Linear Equations

- Observe that

$$\begin{aligned}U^T A \mathbf{v} &= D \mathbf{V}^T \mathbf{v} \\&= \begin{pmatrix} \mathbf{\Sigma} & \mathbf{0}_{r \times (n-r)} \\ \mathbf{0}_{(m-r) \times r} & \mathbf{0}_{(m-r) \times (n-r)} \end{pmatrix} \begin{pmatrix} \{\mathbf{V}^T \mathbf{v}\}_r \\ \{\mathbf{V}^T \mathbf{v}\}_{n-r} \end{pmatrix} \\&= \begin{pmatrix} \mathbf{\Sigma} \{\mathbf{V}^T \mathbf{v}\}_r \\ \mathbf{0}_{m-r} \end{pmatrix}.\end{aligned}$$

# Systems of Linear Equations

► Thus

$$\begin{aligned}\theta(z) &= \left\| \begin{pmatrix} \mathbf{\Sigma}\{\mathbf{V}^T \mathbf{v}\}_r \\ \mathbf{0}_{m-r} \end{pmatrix} - \begin{pmatrix} \mathbf{0}_r \\ \{\mathbf{U}^T \mathbf{b}\}_{m-r} \end{pmatrix} \right\|_2 \\ &= \sqrt{\sum_{\ell=1}^r (\sigma_\ell [\mathbf{V}^T \mathbf{v}]_\ell)^2 + \sum_{\ell=r+1}^m [\mathbf{U}^T \mathbf{b}]_\ell^2} \\ &\geq \tilde{\theta}_{\clubsuit}.\end{aligned}$$

► Therefore the only  $\mathbf{v}$  that satisfies the equal sign is  $\mathbf{0}_n$ .

# Dominant Correlations

- There are two correlation matrices that can be computed from  $\mathbf{X}$  :  $\mathbf{X}\mathbf{X}^T$  and  $\mathbf{X}^T\mathbf{X}$ :

$$\mathbf{X}\mathbf{X}^T = \mathbf{U} \begin{pmatrix} \mathbf{\Sigma} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} \mathbf{V}^T \mathbf{V} \begin{pmatrix} \mathbf{\Sigma} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} \mathbf{U}^T = \mathbf{U} \begin{pmatrix} \mathbf{\Sigma}^2 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} \mathbf{U}^T$$

$$\mathbf{X}^T\mathbf{X} = \mathbf{V} \begin{pmatrix} \mathbf{\Sigma} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} \mathbf{U}^T \mathbf{U} \begin{pmatrix} \mathbf{\Sigma} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} \mathbf{V}^T = \mathbf{V}\mathbf{\Sigma}^2\mathbf{V}^T.$$

# Dominant Correlations

- ▶  $\mathbf{X}\mathbf{X}^T$  is usually much larger than  $\mathbf{X}^T\mathbf{X}$ . Observe that we can find

$$\mathbf{X}^T\mathbf{X}\mathbf{V} = \mathbf{V}\mathbf{\Sigma}^2$$

and from here approximate  $\tilde{\mathbf{U}}$  (the first  $r$  columns of  $\mathbf{U}$ ) as

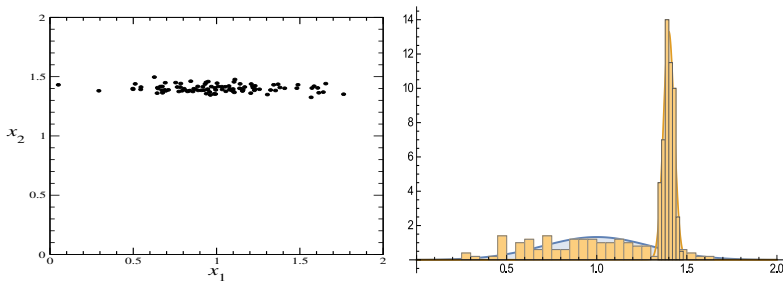
$$\tilde{\mathbf{U}} = \mathbf{X}\tilde{\mathbf{V}}\tilde{\mathbf{\Sigma}}^{-1}.$$

# PCA

- ▶ PCA provides a hierarchical coordinate system to represent high-dimensional correlated data.
- ▶ Preprocessing of data by subtraction of the mean and reducing the variance to unity before performing SVD.
- ▶ The new coordinates are uncorrelated (orthogonal) but guard a maximum correlation with the measurements.
- ▶ Measurements are collected into a row vector. (Feature vector associated to an observable).

# PCA

- Suppose we have a data set composed by 2-dimensional features such that:



- If we filter out the variations on  $x_2$  the 'loss in information' is less than filtering out  $x_1$ .



# PCA

- ▶ Given  $\mathcal{D} = \{\mathbf{x}_n\}_{n=1}^N$ ,  $\mathbf{x} \in \mathbb{R}^d$ , then

$$\boldsymbol{\mu} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n,$$

there is a matrix  $\mathbf{B} \in \mathbb{R}^{d \times q}$  with  $d \geq q > 0$  such that  $\mathbf{z} \in \mathbb{R}^q$  defined as

$$\mathbf{z} = \mathbf{B}^T(\mathbf{x} - \boldsymbol{\mu}).$$

- ▶ Observe that  $\mathbf{B}$  implements a reduction of dimensionality of the elements of the data set.

# PCA

- Observe that

$$\begin{aligned}\mathbb{V}_z[z] &= \mathbb{V}_x[\mathbf{B}^T(\mathbf{x} - \boldsymbol{\mu})] \\&= \mathbb{E}_x[(\mathbf{B}^T(\mathbf{x} - \boldsymbol{\mu}))^2] - \mathbb{E}_x[\mathbf{B}^T(\mathbf{x} - \boldsymbol{\mu})]^2 \\&= \mathbb{E}_x[\mathbf{x}^T \mathbf{B} \mathbf{B}^T \mathbf{x} - 2\mathbf{x}^T \mathbf{B} \mathbf{B}^T \boldsymbol{\mu} + \boldsymbol{\mu}^T \mathbf{B} \mathbf{B}^T \boldsymbol{\mu}] \\&= \mathbb{E}_x[\mathbf{x}^T \mathbf{B} \mathbf{B}^T \mathbf{x}] - \boldsymbol{\mu}^T \mathbf{B} \mathbf{B}^T \boldsymbol{\mu} \\&= \mathbb{V}_x[\mathbf{B}^T \mathbf{x}].\end{aligned}$$

- We want to find  $\mathbf{B}$  that maximizes the variance of  $z$ .

# PCA

- ▶ We proceed sequentially. Let us define

$$V_1 = \mathbb{V}[z_1] = \frac{1}{N} \sum_{n=1}^N z_{1,n}^2.$$

- ▶ We need to find the direction in space  $\mathbb{R}^d$ , represented by the unit vector  $\mathbf{b}_1$ , such that  $z_{1,n} = \mathbf{b}_1^T (\mathbf{x}_n - \boldsymbol{\mu})$  maximizes  $V_1$ . (Observe that the variances in  $\mathbf{z}$  and in  $\mathbf{x}$  without the shift in  $\boldsymbol{\mu}$  are identical)

► Therefore

$$\begin{aligned}V_1 &= \frac{1}{N} \sum_{n=1}^N \left( \mathbf{b}_1^T \mathbf{x}_n \right)^2 \\&= \mathbf{b}_1^T \left( \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^T \right) \mathbf{b}_1 \\&= \mathbf{b}_1^T \mathbf{S} \mathbf{b}_1 \\ \mathbf{S} &= \mathbf{b}_1^T \left( \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^T \right) \mathbf{b}_1,\end{aligned}$$

where  $\mathbf{S}$  is the data covariance matrix.

# PCA

- ▶ The problem can be stated as:

$$\max_{\mathbf{b}_1} \mathbf{b}_1^T \mathbf{S} \mathbf{b}_1, \text{ subject to } \mathbf{b}_1^T \mathbf{b}_1 = 1.$$

- ▶ The constraint maximization problem can be defined with the aid of a Lagrange multiplier  $\lambda_1$ , such that:

$$\mathcal{L}_1(\mathbf{b}_1, \lambda_1) = \mathbf{b}_1^T \mathbf{S} \mathbf{b}_1 + \lambda_1(1 - \mathbf{b}_1^T \mathbf{b}_1).$$

- ▶ The equations to be solved are:

$$\nabla_{\mathbf{b}_1} \mathcal{L}_1 = 2\mathbf{S} \mathbf{b}_1 - 2\lambda_1 \mathbf{b}_1 = \mathbf{0}$$

$$\frac{\partial \mathcal{L}_1}{\partial \lambda_1} = 1 - \mathbf{b}_1^T \mathbf{b}_1 = 0 \quad (\text{normalization})$$

# PCA

- ▶ The solutions to the problem are the eigenvectors of  $\mathbf{S}$ . We keep the one correspondent to the largest eigenvalue.
- ▶ Subsequent projections can be obtained by increasing the number of orthogonality constraints with previously found vectors, i.e.

$$\mathcal{L}_2(\mathbf{b}_2, \lambda_1, \gamma_{1,2}) = \mathbf{b}_2^T \mathbf{S} \mathbf{b}_2 + \lambda_2(1 - \mathbf{b}_2^T \mathbf{b}_2) + \gamma_{1,2} \mathbf{b}_1^T \mathbf{b}_2$$

$$\nabla_{\mathbf{b}_2} \mathcal{L}_2 = 2\mathbf{S} \mathbf{b}_2 - 2\lambda_2 \mathbf{b}_2 + \gamma_{1,2} \mathbf{b}_1 = \mathbf{0}$$

$$\frac{\partial \mathcal{L}_2}{\partial \lambda_2} = 1 - \mathbf{b}_2^T \mathbf{b}_2 = 0 \quad (\text{normalization})$$

$$\frac{\partial \mathcal{L}_2}{\partial \gamma_{1,2}} = \mathbf{b}_1^T \mathbf{b}_2 = 0 \quad (\text{orthogonality}).$$

# PCA

- From the first equation we obtain:

$$\begin{aligned}\mathbf{S}\mathbf{b}_2 &= \lambda_2\mathbf{b}_2 - \frac{\gamma_{1,2}}{2}\mathbf{b}_1 \\ \mathbf{b}_1^T\mathbf{S}\mathbf{b}_2 &= \lambda_2\mathbf{b}_1^T\mathbf{b}_2 - \frac{\gamma_{1,2}}{2} \\ &= (\mathbf{b}_2^T\mathbf{S}\mathbf{b}_1)^T = \lambda_1\mathbf{b}_2^T\mathbf{b}_1\end{aligned}$$

and by the orthogonality  $\mathbf{b}_1^T\mathbf{b}_2 = 0$  we have that,  $\gamma_{1,2} = 0$  and

$$\mathbf{S}\mathbf{b}_2 = \lambda_2\mathbf{b}_2.$$

- $\mathbf{b}_2$  must be the eigenvector of  $\mathbf{S}$  corresponding to the second largest eigenvalue.

# PCA

- ▶ Given the spectrum  $\{\lambda_1, \dots, \lambda_d\}$  of  $\mathbf{S}$ , with  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d > 0$ , the matrix  $\mathbf{B}$  is then  $\mathbf{B} = (\mathbf{b}_1, \dots, \mathbf{b}_q)$ .
- ▶ The vector  $\mathbf{z}_n$  becomes

$$z_{n,j} = \mathbf{b}_j^T (\mathbf{x}_n - \boldsymbol{\mu}).$$

- ▶ The total fraction of the variance that is carried into the  $\mathbf{z}$ -representation is

$$\frac{\sum_{j=1}^q \lambda_j}{\sum_{j=1}^d \lambda_j}.$$