

# ***Lecture programme***

---

## Lecture 3:

- Dimension reduction with PCA and SVD
- Clustering I (clustering tendency)

**Book:** Sec. 2.4.3, 6.1–6.2

- L4: Clustering II ( $K$ -representatives, hierarchical)
- L5: Clustering III (spectral, validation)

# Why all eigen and singular stuff?

---

**Goal:** nice low-dimensional representation for data, when

- original data high-dimensional
- only pairwise distances/similarities known

**Recall:** Given  $n \times n$  matrix  $\mathbf{A}$ ,  $\lambda$  is eigenvalue and  $\mathbf{v} \neq \mathbf{0}$  corresponding eigenvector, if  $\mathbf{A}\mathbf{v} = \lambda\mathbf{v}$

**Good news:** Everything will be real-valued!

- $k \times k$  covariance matrix and  $n \times n$  Laplacian matrices positive semidefinite <sup>a</sup>  $\Rightarrow$  real  $\lambda$ s and  $\mathbf{v}$ s
- $n \times k$  data matrix real  $\Rightarrow$  singular vectors real (singular values always real)

---

$$^a \mathbf{z}^T \mathbf{A} \mathbf{z} \geq 0 \text{ for all real } \mathbf{z} \neq \mathbf{0}$$

# ***Dimension reduction: motivation***

---

1. **Curse of dimensionality:** hard to distinguish close and far neighbours in high dimensional data! → how to find clusters??
  2. **Redundancy in data:** If features strongly correlated, the same information can be presented with a smaller number of features
    - **intrinsic dimensionality**  $r$  may be  $r \ll d$
- ⇒ Idea: reduce dimensionality by removing this redundancy!

# ***Principal component analysis (PCA) assumptions***

---

1. High **variance** reflects important structures of data.
2. Data can be presented well as a **linear** combination of suitable **orthogonal** basis vectors (PCs).

Idea: Given  $n \times d$  data and suitable  $d \times r$  ( $r < d$ ) matrix  $\mathbf{P}_r$ , new data will be  $n \times r$  matrix  $\mathbf{D}\mathbf{P}_r$

Remember: These assumptions may not always reflect reality!

# PCA intuition

Rotate axes to match highest variance directions + choose the best new axes

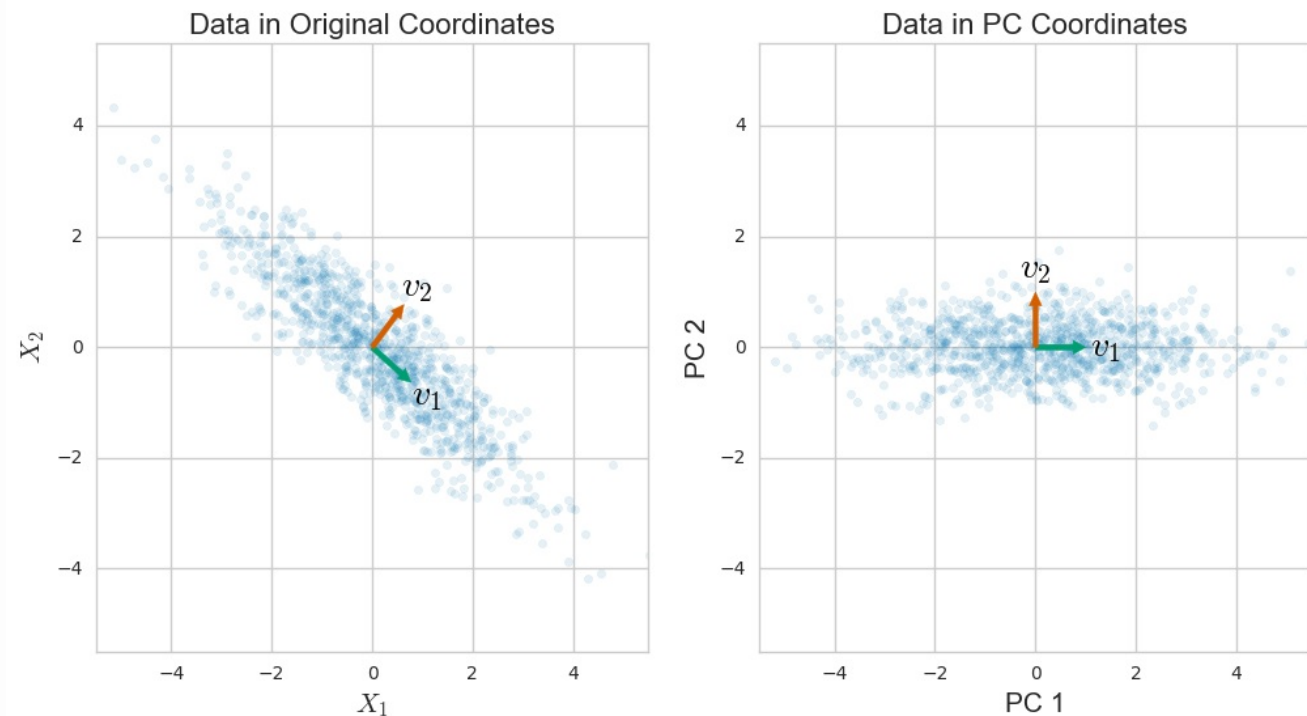


Image source: <https://intoli.com/blog/pca-and-svd>

# PCA: Use eigen-decomposition

---

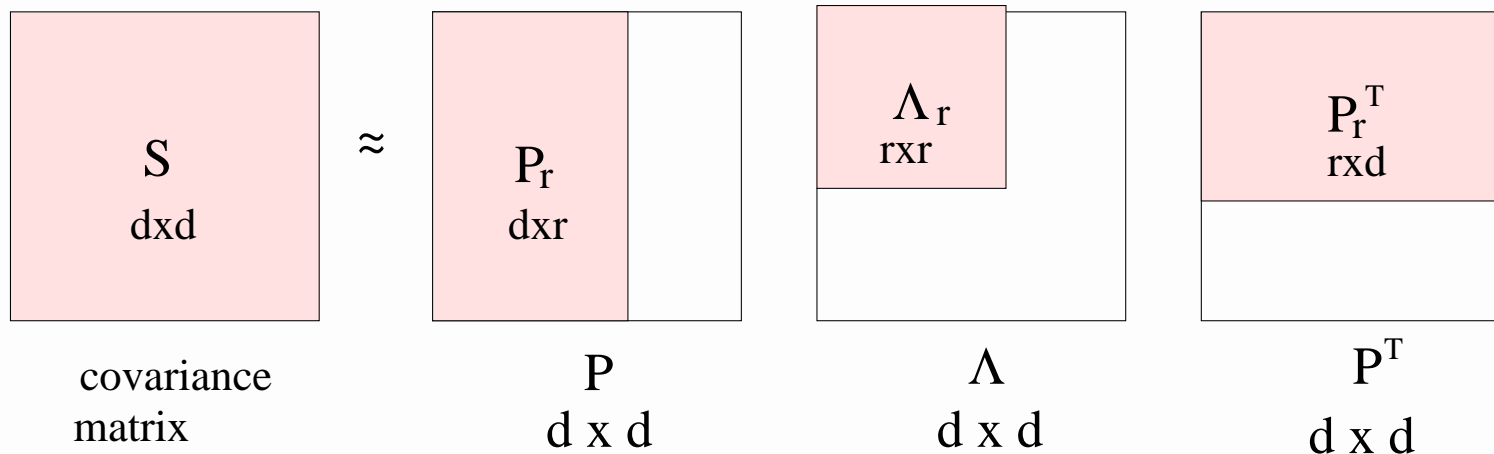
- Assume  $\mathbf{D}$  mean-centered. Covariance matrix  $\mathbf{C} = \frac{1}{n-1} \mathbf{D}^T \mathbf{D}$  <sup>a</sup>
- since  $\mathbf{C}$  positive semidefinite, it can be diagonalized:  
 $\mathbf{C} = \mathbf{P} \mathbf{\Lambda} \mathbf{P}^T$ 
  - $\mathbf{P}$ 's columns= $\mathbf{C}$ 's orthonormal eigenvectors
  - $\mathbf{\Lambda}$  diagonal,  $\Lambda_{ii}$ =eigenvalues
- transformed data  $\mathbf{D}' = \mathbf{D} \mathbf{P}$ 
  - $\mathbf{\Lambda}$  new covariance matrix (diagonal, i.e., no correlations)

---

<sup>a</sup>unbiased estimate; for large  $n$  also  $\frac{1}{n} \mathbf{D}^T \mathbf{D}$  ok

# PCA: Dimension reduction

- Assume  $\Lambda$  ordered into decreasing order by  $\lambda_i = \Lambda_{ii}$
- Keep only  $r$  largest  $(\lambda_1, \dots, \lambda_r)$  + corresponding eigenvectors
- approximate data  $\mathbf{D}' = \mathbf{D}\mathbf{P}_r$



# Example: Dimension reduction with PCA

Data **D**

$$\begin{bmatrix} 2 & 2 & 1 & 2 & 0 & 0 \\ 2 & 3 & 3 & 3 & 0 & 0 \\ 1 & 1 & 1 & 1 & 0 & 0 \\ 2 & 2 & 2 & 3 & 1 & 1 \\ 0 & 0 & 0 & 1 & 1 & 1 \\ 0 & 0 & 0 & 2 & 1 & 2 \end{bmatrix}$$

Mean-centered data <sup>a</sup>

$$\begin{bmatrix} 0.83 & 0.67 & -0.17 & 0.00 & -0.50 & -0.67 \\ 0.83 & 1.67 & 1.83 & 1.00 & -0.50 & -0.67 \\ -0.17 & -0.33 & -0.17 & -1.00 & -0.50 & -0.67 \\ 0.83 & 0.67 & 0.83 & 1.00 & 0.50 & 0.33 \\ -1.17 & -1.33 & -1.17 & -1.00 & 0.50 & 0.33 \\ -1.17 & -1.33 & -1.17 & 0.00 & 0.50 & 1.33 \end{bmatrix}$$

---

<sup>a</sup>rounded for the slide presentation only!

Mean vector  $\approx [1.167, 1.333, 1.167, 2, 0.5, 0.667]$



## ***Example (continued)***

---

Covariance matrix <sup>a</sup>:

$$\mathbf{C} \approx \begin{bmatrix} 0.97 & 1.13 & 0.97 & 0.6 & -0.3 & -0.53 \\ 1.13 & 1.47 & 1.33 & 0.8 & -0.4 & -0.67 \\ 0.97 & 1.33 & 1.37 & 0.80 & -0.3 & -0.53 \\ 0.60 & 0.80 & 0.80 & 0.80 & 0.00 & 0.00 \\ -0.30 & -0.40 & -0.30 & 0.00 & 0.30 & 0.40 \\ -0.53 & -0.67 & -0.53 & 0.00 & 0.40 & 0.67 \end{bmatrix}$$

---

<sup>a</sup>rounded for the slide presentation only!

## Example (continued)

Eigenvalues  $\lambda_i$ : 4.43, 0.89, 0.17, 0.066, 0.014,  $1.2e-17$

Eigenvectors  $\mathbf{v}_i$  (column vectors) <sup>a</sup>:

$$\begin{bmatrix} 0.44 & 0.058 & 0.68 & -0.34 & -0.47 & 1.8e-15 \\ 0.57 & 0.027 & 0.092 & 0.18 & 0.54 & 0.58 \\ 0.53 & -0.14 & -0.71 & -0.26 & -0.35 & 1.5e-15 \\ 0.32 & -0.62 & 0.14 & 0.37 & 0.16 & -0.58 \\ -0.15 & -0.42 & 0.041 & -0.78 & 0.43 & 8.4e-17 \\ -0.26 & -0.65 & 0.052 & 0.19 & -0.38 & 0.58 \end{bmatrix}$$

**What shall we do if we want a 2-dimensional representation of data?**

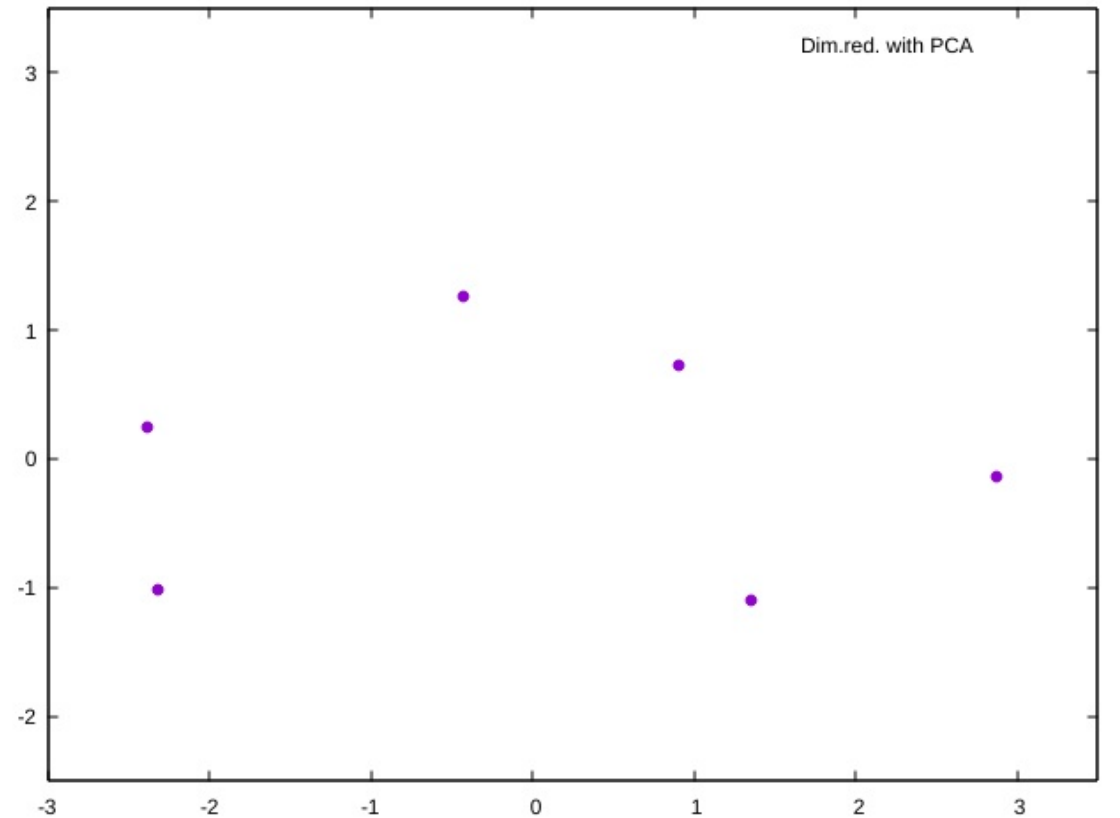
---

<sup>a</sup>All rounded for the presentation

## Example (continued)

Since  $\lambda_1$  and  $\lambda_2$  largest, set  $\mathbf{P}_2 = [\mathbf{v}_1, \mathbf{v}_2]$ .

$$\mathbf{D}' = \mathbf{D}\mathbf{P}_2 \approx \begin{bmatrix} 0.91 & 0.73 \\ 2.87 & -0.14 \\ -0.43 & 1.26 \\ 1.35 & -1.09 \\ -2.38 & 0.25 \\ -2.32 & -1.01 \end{bmatrix}$$



# When PCA can fail?

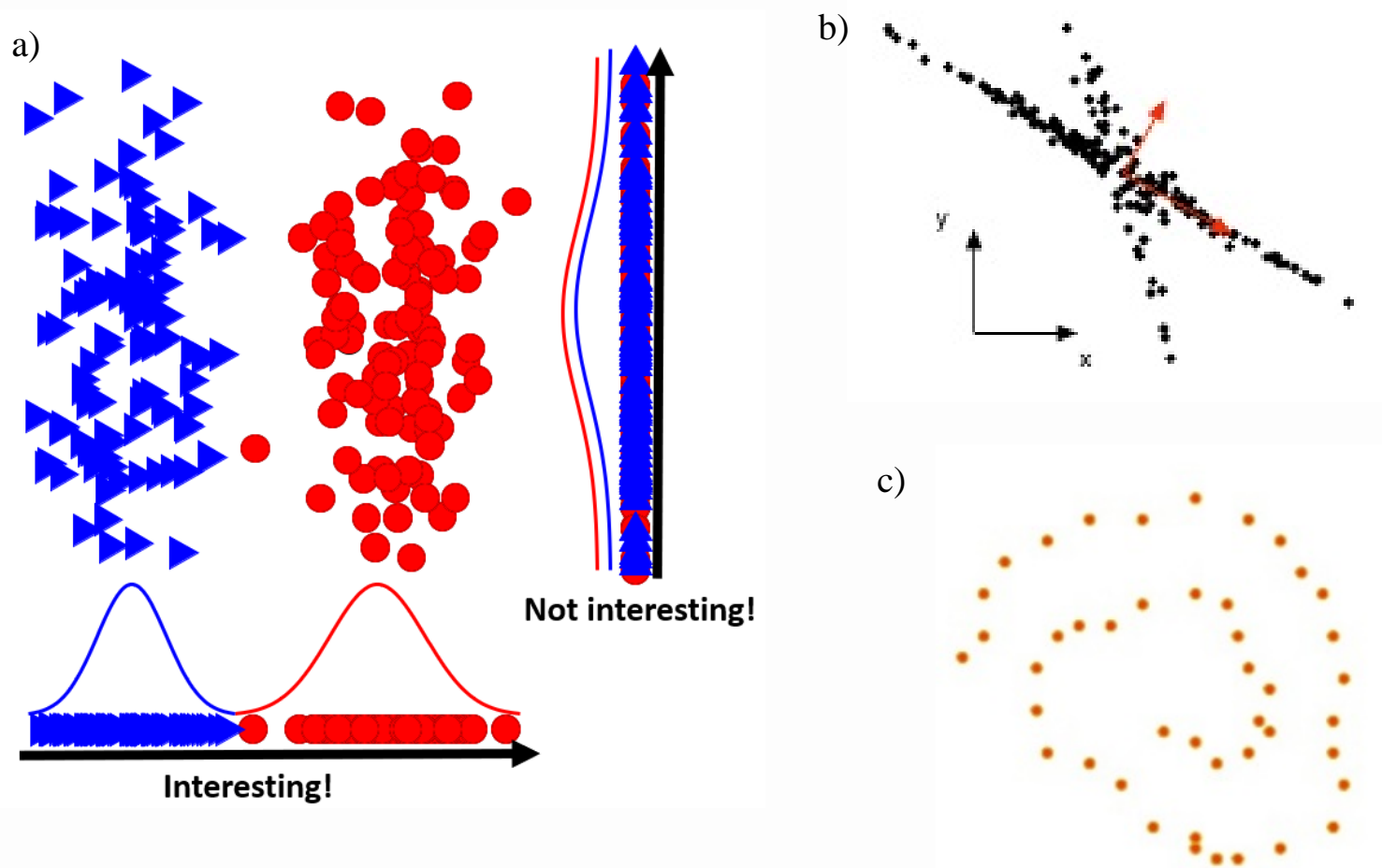


Image sources: <https://towardsdatascience.com/interesting-projections-where-pca-fails-fe64ddca73e6> and Shlen's good tutorial: <https://arxiv.org/abs/1404.1100>

# Singular value decompositions (SVD)

---

Factorize  $\mathbf{D}$  as  $\mathbf{D} = \mathbf{Q}\mathbf{\Sigma}\mathbf{P}^T$  <sup>a</sup>

- $\mathbf{\Sigma}$  diagonal,  $\Sigma_{ii} = \sigma_i$  **singular values**
- $\mathbf{Q}$ 's columns **left singular vectors**, (orthonormal eigenvectors of  $\mathbf{D}\mathbf{D}^T$ )
- $\mathbf{P}$ 's columns **right singular vectors**, (orthonormal eigenvectors of  $\mathbf{D}^T\mathbf{D}$ )
- transformed data  $\mathbf{D}' = \mathbf{D}\mathbf{P}$ 
  - if  $\mathbf{D}$  mean-centered, same basis vectors as PCA <sup>b</sup>
  - mean-centering often skipped, if  $\mathbf{D}$  sparse, non-negative (e.g., document-word matrices)

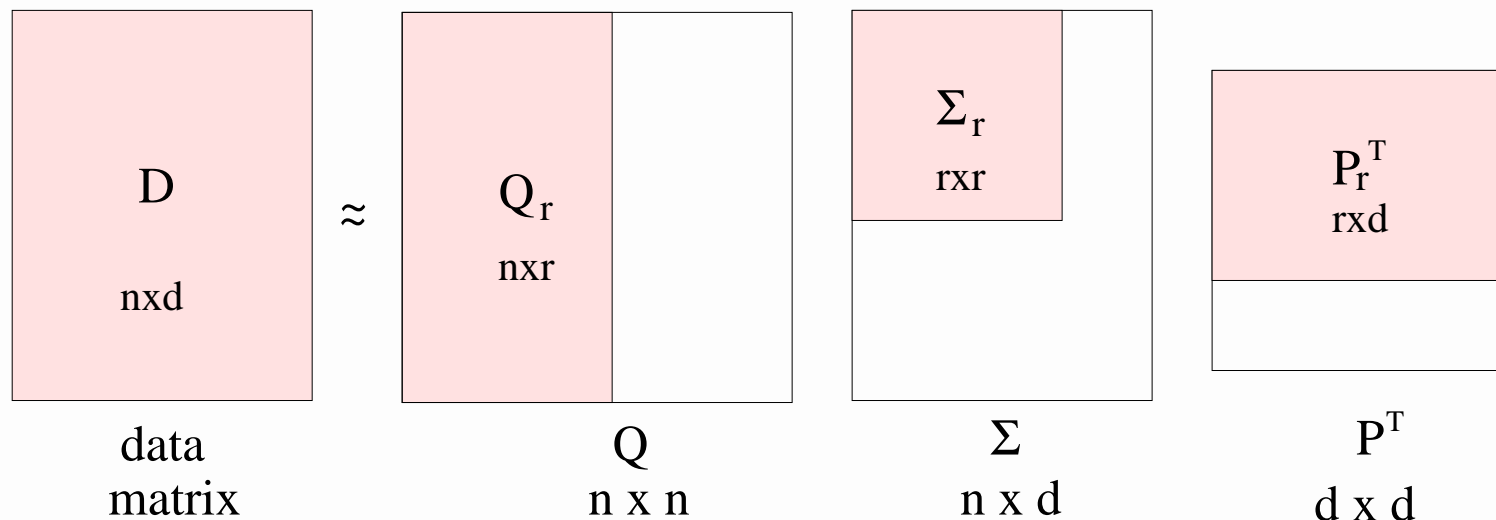
---

<sup>a</sup>always possible

<sup>b</sup>may be opposite direction

# Truncated SVD: Dimension reduction

- Assume  $\Sigma$  ordered into decreasing order by  $\sigma_i = \Sigma_{ii}$ .
- Keep only  $r$  largest  $(\sigma_1, \dots, \sigma_r)$  + corresponding singular vectors of  $\mathbf{P}$
- approximate data  $\mathbf{D}' = \mathbf{D}\mathbf{P}_r$



## Previous example with SVD

Let's first test without mean-centering:

```
Q
[[-4.10936057e-01  1.74569815e-01  8.24528531e-01  2.52257347e-01 -2.39114763e-01  1.29454797e-16]
 [-6.45804322e-01  3.14417109e-01 -5.61570935e-01  3.01160896e-01 -2.79318562e-01  1.25912076e-16]
 [-2.31559546e-01  1.26698028e-01  3.39347806e-02 -9.93766673e-02  5.02626927e-01  8.16496581e-01]
 [-5.62143219e-01 -2.03086484e-01  4.38362617e-02 -6.02554461e-01  3.33302743e-01 -4.08248290e-01]
 [-9.90241264e-02 -4.56482541e-01 -2.40332995e-02 -4.03801126e-01 -6.71951112e-01  4.08248290e-01]
 [-1.86112160e-01 -7.77813886e-01 -3.37638835e-02  5.56475872e-01  2.22626206e-01 -3.01457954e-16]]

Sigma
[8.42523943e+00  3.26119142e+00  9.87979172e-01  5.74286469e-01  2.72145612e-01  2.01264667e-17]

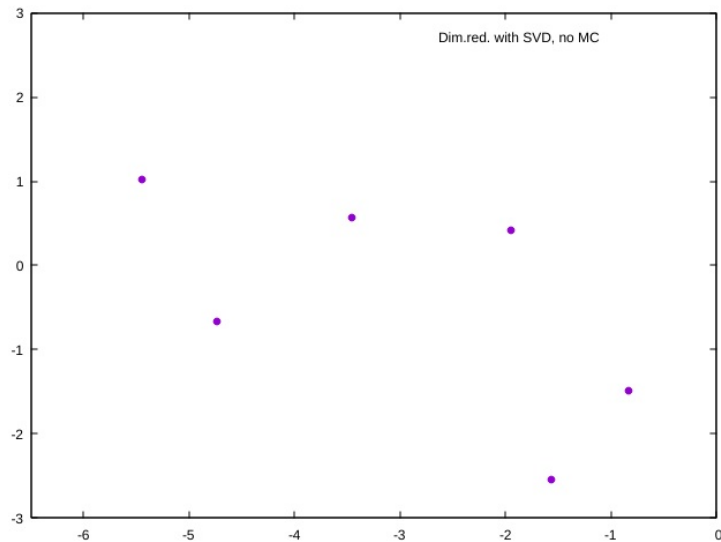
P^T
[[-4.11777822e-01 -4.88428975e-01 -4.39654568e-01 -6.11083254e-01 -1.00564442e-01 -1.22654279e-01]
 [ 2.14185191e-01  3.10596922e-01  2.57067462e-01 -3.68663051e-01 -4.40753922e-01 -6.79259973e-01]
 [ 6.55400957e-01  8.69973391e-02 -7.47563302e-01  3.86918624e-02 -1.41307851e-02 -4.83054767e-02]
 [-3.44164653e-01  1.80244178e-01 -2.59009332e-01  3.65859132e-01 -7.83371609e-01  1.85614954e-01]
 [ 4.86378460e-01 -5.39978569e-01  3.38649460e-01 -1.48261638e-01 -4.26323839e-01  3.91716930e-01]
 [-0.00000000e+00 -5.77350269e-01  1.66533454e-16  5.77350269e-01  9.12464548e-16 -5.77350269e-01]]
```

$\sigma_1$  and  $\sigma_2$  largest  $\Rightarrow \mathbf{P}_2$  = first two columns of  $\mathbf{P}$   
+ calculate  $\mathbf{D}' = \mathbf{D}\mathbf{P}_2$

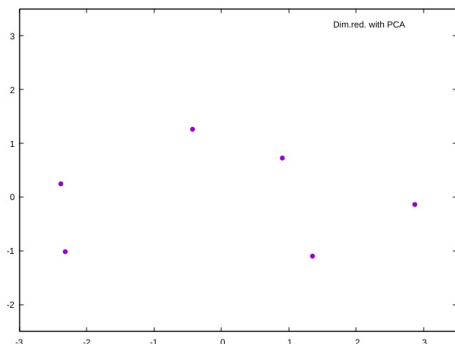
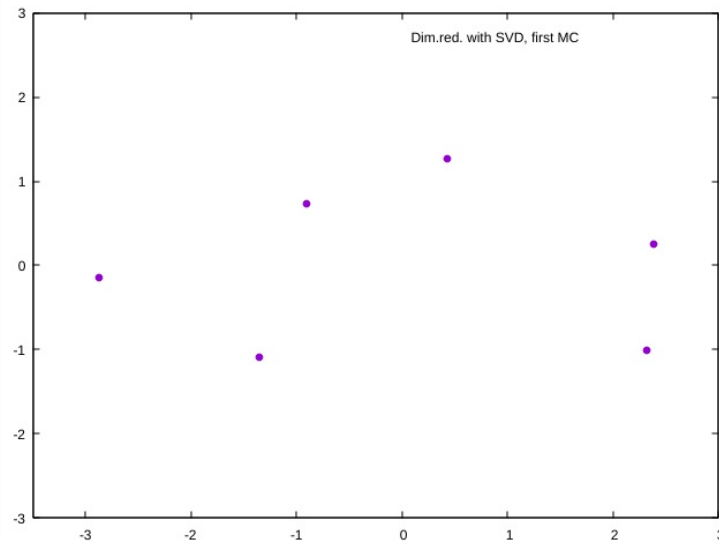
# Example continued

Note: Different  $\mathbf{Q}$ ,  $\Sigma$ , and  $\mathbf{P}$ , if mean-centered data

$\mathbf{D}'$  without mean-centering



$\mathbf{D}'$  with mean-centering





# ***SVD applications***

---

- dimension reduction:  $\mathbf{D}' = \mathbf{D}\mathbf{P}_r$
- sometimes  $\mathbf{Q}$  also useful, e.g., user-item rating matrix
  - $\mathbf{D}^T \mathbf{Q}_r$  describes items by  $r$  latent components
- **Latent semantic analysis (LSA)** applies SVD on document-term matrix (e.g., tf-idf matrix)
  - often drastic dimension reduction!
  - helps with noise due to synonymous words
  - e.g., {(car), (truck), (flower)}  $\rightarrow$  {(1.3452 · car + 0.2828 · truck), (flower)}
- noise reduction: truncated SVD tends to correct inconsistencies

# Example of truncated SVD (Aggarwal p. 45)

$$D = \begin{pmatrix} 2 & 2 & 1 & 2 & 0 & 0 \\ 2 & 3 & 3 & 3 & 0 & 0 \\ 1 & 1 & 1 & 1 & 0 & 0 \\ 2 & 2 & 2 & 3 & 1 & 1 \\ 0 & 0 & 0 & 1 & 1 & 1 \\ 0 & 0 & 0 & 2 & 1 & 2 \end{pmatrix} \approx Q_2 \Sigma_2 P_2^T$$

$$\approx \begin{pmatrix} -0.41 & 0.17 \\ -0.65 & 0.31 \\ -0.23 & 0.13 \\ -0.56 & -0.20 \\ -0.10 & -0.46 \\ -0.19 & -0.78 \end{pmatrix} \begin{pmatrix} 8.4 & 0 \\ 0 & 3.3 \end{pmatrix} \begin{pmatrix} -0.41 & -0.49 & -0.44 & -0.61 & -0.10 & -0.12 \\ 0.21 & 0.31 & 0.26 & -0.37 & -0.44 & -0.68 \end{pmatrix}$$

$$= \begin{pmatrix} 1.55 & 1.87 & \underline{1.67} & 1.91 & 0.10 & 0.04 \\ 2.46 & 2.98 & 2.66 & 2.95 & 0.10 & -0.03 \\ 0.89 & 1.08 & 0.96 & 1.04 & 0.01 & -0.04 \\ 1.81 & 2.11 & 1.91 & 3.14 & 0.77 & 1.03 \\ 0.02 & -0.05 & -0.02 & 1.06 & 0.74 & 1.11 \\ 0.10 & -0.02 & 0.04 & 1.89 & 1.28 & 1.92 \end{pmatrix}$$

# Summary

---

1. Factorize  $\mathbf{C} = \mathbf{P}\mathbf{\Lambda}\mathbf{P}^T$  or  $\mathbf{D} = \mathbf{Q}\mathbf{\Sigma}\mathbf{P}^T$
2. Use  $\mathbf{P}_r$  (best components of  $\mathbf{P}$ ).  
Reduced data  $\mathbf{D}' = \mathbf{D}\mathbf{P}_r$

Only heuristics! Work well only if the underlying assumptions are true.