# FSAN/ELEG815: Statistical Learning

Gonzalo R. Arce

**Department of Electrical and Computer Engineering**
**University of Delaware**

2. Eigen Analysis, SVD, PCA, and Matrix Completion

# Outline

Eigen Analysis

Eigen Properties

SVD

PCA

# Eigen Analysis

Objective: Utilize tools from linear algebra to characterize and analyze matrices, especially the correlation matrix

- ▶ The correlation matrix plays a large role in statistical characterization and processing.
- ▶ Previously result: $\mathbf{R}$ is Hermitian.
- ▶ Further insight into the correlation matrix is achieved through eigen analysis

Objective: For a Hermitian matrix $\mathbf{R}$, find a vector $\mathbf{q}$ satisfying

$$\mathbf{R}\mathbf{q} = \lambda\mathbf{q}$$

▶ Interpretation: Linear transformation by $\mathbf{R}$ changes the scale, but not the direction of $\mathbf{q}$

▶ Fact: A $M \times M$ matrix $\mathbf{R}$ has $M$ eigenvectors and eigenvalues

$$\mathbf{R}\mathbf{q}_i = \lambda_i\mathbf{q}_i \quad i = 1, 2, 3, \cdots, M$$

To see this, note

$$(\mathbf{R} - \lambda\mathbf{I})\mathbf{q} = \mathbf{0}$$

For this to be true, the row/columns of $(\mathbf{R} - \lambda\mathbf{I})$ must be linearly dependent,

$$\Rightarrow \det(\mathbf{R} - \lambda\mathbf{I}) = 0$$

Note: $\det(\mathbf{R} - \lambda \mathbf{I})$ is a $M$th order polynomial in $\lambda$

▶ The roots of the polynomial are the eigenvalues $\lambda_1, \lambda_2, \cdots, \lambda_M$

$$\mathbf{R}\mathbf{q}_i = \lambda_i \mathbf{q}_i$$

▶ Each eigenvector $\mathbf{q}_i$ is associated with one eigenvalue $\lambda_i$

▶ The eigenvectors are not unique

$$\begin{aligned} \mathbf{R}\mathbf{q}_i &= \lambda_i \mathbf{q}_i \\ \Rightarrow \mathbf{R}(a\mathbf{q}_i) &= \lambda_i(a\mathbf{q}_i) \end{aligned}$$

Consequence: eigenvectors are generally normalized, e.g., $|\mathbf{q}_i| = 1$ for $i = 1, 2, \ldots, M$

## Example (General two dimensional case)

Let $M = 2$ and

$$\mathbf{R} = \left[ \begin{array}{cc} R_{1,1} & R_{1,2} \\ R_{2,1} & R_{2,2} \end{array} \right]$$

Determine the eigenvalues and eigenvectors.

Thus

$$
\begin{aligned}
\det(\mathbf{R} - \lambda \mathbf{I}) &= 0 \\
\Rightarrow \left| \begin{array}{cc} R_{1,1} - \lambda & R_{1,2} \\ R_{2,1} & R_{2,2} - \lambda \end{array} \right| &= 0 \\
\Rightarrow \lambda^2 - \lambda(R_{1,1} + R_{2,2}) + (R_{1,1}R_{2,2} - R_{1,2}R_{2,1}) &= 0 \\
\Rightarrow \lambda_{1,2} = \frac{1}{2} \left[ (R_{1,1} + R_{2,2}) \pm \sqrt{4R_{1,2}R_{2,1} + (R_{1,1} - R_{2,2})} \right]
\end{aligned}
$$

Back substitution yields the eigenvectors:

$$\left[ \begin{array}{cc} R_{1,1} - \lambda & R_{1,2} \\ R_{2,1} & R_{2,2} - \lambda \end{array} \right] \left[ \begin{array}{c} q_1 \\ q_2 \end{array} \right] = \left[ \begin{array}{c} 0 \\ 0 \end{array} \right]$$

In general, this yields a set of linear equations. In the $M = 2$ case:

$$(R_{1,1} - \lambda)q_1 + R_{1,2}q_2 = 0$$
$$R_{2,1}q_1 + (R_{2,2} - \lambda)q_2 = 0$$

▶ Solving the set of linear equations for a specific eigenvalue $\lambda_i$ yields the corresponding eigenvector, $\mathbf{q}_i$

## Example (Two–dimensional white noise)

Let $\mathbf{R}$ be the correlation matrix of a two–sample vector of zero mean white noise

$$\mathbf{R} = \left[ \begin{array}{cc} \sigma^2 & 0 \\ 0 & \sigma^2 \end{array} \right]$$

Determine the eigenvalues and eigenvectors.

Carrying out the analysis yields eigenvalues

$$\begin{aligned} \lambda_{1,2} &= \frac{1}{2} \left[ (R_{1,1} + R_{2,2}) \pm \sqrt{4 R_{1,2} R_{2,1} + (R_{1,1} - R_{2,2})} \right] \\ &= \frac{1}{2} \left[ (\sigma^2 + \sigma^2) \pm \sqrt{0 + (\sigma^2 - \sigma^2)} \right] = \sigma^2 \end{aligned}$$

and eigenvectors

$$\mathbf{q}_1 = \left[ \begin{array}{c} 1 \\ 0 \end{array} \right] \quad \text{and} \quad \mathbf{q}_2 = \left[ \begin{array}{c} 0 \\ 1 \end{array} \right]$$

Note: The eigenvectors are unit length (and orthogonal)

## Property (Correlation matrix eigenvalues are real & nonnegative)

The eigenvalues of $\mathbf{R}$ are real and nonnegative.

Proof:

$$
\begin{aligned}
\mathbf{R}\mathbf{q}_i &= \lambda_i \mathbf{q}_i \\
\Rightarrow \mathbf{q}_i^H \mathbf{R}\mathbf{q}_i &= \lambda_i \mathbf{q}_i^H \mathbf{q}_i \qquad [\text{pre--multiply by } \mathbf{q}_i^H] \\
\Rightarrow \lambda_i &= \frac{\mathbf{q}_i^H \mathbf{R}\mathbf{q}_i}{\mathbf{q}_i^H \mathbf{q}_i} \geq 0
\end{aligned}
$$

Follows from the facts: $\mathbf{R}$ is positive semi-definite and $\mathbf{q}_i^H \mathbf{q}_i = |\mathbf{q_i}|^2 > 0$

Note: In most cases, $\mathbf{R}$ is positive definite and

$$
\lambda_i > 0, \qquad i = 1, 2, \cdots, M
$$

## Property (Unique eigenvalues $\Rightarrow$ orthogonal eigenvectors)

If $\lambda_1, \lambda_2, \cdots, \lambda_M$ are unique eigenvalues of $\mathbf{R}$, then the corresponding eigenvectors, $\mathbf{q}_1, \mathbf{q}_2, \cdots, \mathbf{q}_M$, are orthogonal.

Proof:

$$
\begin{aligned}
\mathbf{R}\mathbf{q}_i &= \lambda_i \mathbf{q}_i \\
\Rightarrow \mathbf{q}_j^H \mathbf{R} \mathbf{q}_i &= \lambda_i \mathbf{q}_j^H \mathbf{q}_i \qquad (*)
\end{aligned}
$$

Also, since $\lambda_j$ is real and $\mathbf{R}$ is Hermitian

$$
\begin{aligned}
\mathbf{R}\mathbf{q}_j &= \lambda_j \mathbf{q}_j \\
\Rightarrow \mathbf{q}_j^H \mathbf{R} &= \lambda_j \mathbf{q}_j^H \\
\Rightarrow \mathbf{q}_j^H \mathbf{R} \mathbf{q}_i &= \lambda_j \mathbf{q}_j^H \mathbf{q}_i
\end{aligned}
$$

Substituting the LHS from $(*)$

$$
\Rightarrow \lambda_i \mathbf{q}_j^H \mathbf{q}_i = \lambda_j \mathbf{q}_j^H \mathbf{q}_i
$$

Thus

$$\lambda_i \mathbf{q}_j^H \mathbf{q}_i = \lambda_j \mathbf{q}_j^H \mathbf{q}_i$$
$$\Rightarrow (\lambda_i - \lambda_j) \mathbf{q}_j^H \mathbf{q}_i = 0$$

Since $\lambda_1, \lambda_2, \cdots, \lambda_M$ are unique

$$\mathbf{q}_j^H \mathbf{q}_i = 0 \qquad i \neq j$$

$\Rightarrow \mathbf{q}_1, \mathbf{q}_2, \cdots, \mathbf{q}_M$ are orthogonal.

QED

# Diagonalization of $\mathbf{R}$

Objective: Find a transformation that transforms the correlation matrix into a diagonal matrix.

Let $\lambda_1, \lambda_2, \cdots, \lambda_M$ be unique eigenvectors of $\mathbf{R}$ and take $\mathbf{q}_1, \mathbf{q}_2, \cdots, \mathbf{q}_M$ to be the $M$ orthonormal eigenvectors

$$\mathbf{q}_i^H \mathbf{q}_j = \left\{ \begin{array}{ll} 1 & i = j \\ 0 & i \neq j \end{array} \right.$$

Define $\mathbf{Q} = [\mathbf{q}_1, \mathbf{q}_2, \cdots, \mathbf{q}_M]$ and $\mathbf{\Omega} = \mathrm{diag}(\lambda_1, \lambda_2, \cdots, \lambda_M)$. Then consider

$$\mathbf{Q}^H \mathbf{R} \mathbf{Q} \;\; = \;\; \left[ \begin{array}{c} \mathbf{q}_1^H \\ \mathbf{q}_2^H \\ \vdots \\ \mathbf{q}_M^H \end{array} \right] \mathbf{R}[\mathbf{q}_1, \mathbf{q}_2, \cdots, \mathbf{q}_M]$$

$$\mathbf{Q}^H \mathbf{R} \mathbf{Q} = \begin{bmatrix} \mathbf{q}_1^H \\ \mathbf{q}_2^H \\ \vdots \\ \mathbf{q}_M^H \end{bmatrix} \mathbf{R}[\mathbf{q}_1, \mathbf{q}_2, \cdots, \mathbf{q}_M]$$

$$= \begin{bmatrix} \mathbf{q}_1^H \\ \mathbf{q}_2^H \\ \vdots \\ \mathbf{q}_M^H \end{bmatrix} [\lambda_1 \mathbf{q}_1, \lambda_2 \mathbf{q}_2, \cdots, \lambda_N \mathbf{q}_M]$$

$$= \begin{bmatrix} \lambda_1 & 0 & \cdots & 0 \\ 0 & \lambda_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \lambda_M \end{bmatrix}$$

$$\Rightarrow \mathbf{Q}^H \mathbf{R} \mathbf{Q} = \mathbf{\Omega} \quad \text{(eigenvector diagonalization of } \mathbf{R})$$

## Property ($\mathbf{Q}$ is unitary)

$\mathbf{Q}$ is unitary, i.e., $\mathbf{Q}^{-1} = \mathbf{Q}^H$

Proof: Since the $\mathbf{q}_i$ eigenvectors are orthonormal

$$\mathbf{Q}^H\mathbf{Q} = \begin{bmatrix} \mathbf{q}_1^H \\ \mathbf{q}_2^H \\ \vdots \\ \mathbf{q}_M^H \end{bmatrix} [\mathbf{q}_1, \mathbf{q}_2, \cdots, \mathbf{q}_M] = \mathbf{I}$$

$$\Rightarrow \mathbf{Q}^{-1} = \mathbf{Q}^H$$

## Property (Eigen decomposition of $\mathbf{R}$)

The correlation matrix can be expressed as

$$\mathbf{R} = \sum_{i=1}^{M} \lambda_i \mathbf{q}_i \mathbf{q}_i^H$$

Proof: The correlation diagonalization result states

$$\mathbf{Q}^H \mathbf{R} \mathbf{Q} = \mathbf{\Omega}$$

Isolating $\mathbf{R}$ and expanding,

$$
\begin{aligned}
\mathbf{R} &= \mathbf{Q}\mathbf{\Omega}\mathbf{Q}^H = [\mathbf{q}_1, \mathbf{q}_2, \cdots, \mathbf{q}_M] \mathbf{\Omega}
\begin{bmatrix}
\mathbf{q}_1^H \\
\mathbf{q}_2^H \\
\vdots \\
\mathbf{q}_M^H
\end{bmatrix} \\
&= [\mathbf{q}_1, \mathbf{q}_2, \cdots, \mathbf{q}_M]
\begin{bmatrix}
\lambda_1 \mathbf{q}_1^H \\
\lambda_2 \mathbf{q}_2^H \\
\vdots \\
\lambda_M \mathbf{q}_M^H
\end{bmatrix}
= \sum_{i=1}^{M} \lambda_i \mathbf{q}_i \mathbf{q}_i^H
\end{aligned}
$$

Note: This also gives

$$\mathbf{R}^{-1} = (\mathbf{Q}^H)^{-1} \mathbf{\Omega}^{-1} \mathbf{Q}^{-1} = \mathbf{Q}\mathbf{\Omega}^{-1}\mathbf{Q}^H$$

where $\mathbf{\Omega}^{-1} = \text{diag}(1/\lambda_1, 1/\lambda_2, \cdots, 1/\lambda_M)$

## Aside (trace & determinant for matrix products)

Note $\mathsf{trace}(\boldsymbol{A}) \stackrel{\triangle}{=} \sum_i A_{i,i}$. Also,

$$\mathsf{trace}(\boldsymbol{AB}) = \mathsf{trace}(\boldsymbol{BA}) \qquad \text{similarly} \qquad \det(\boldsymbol{AB}) = \det(\boldsymbol{A})\det(\boldsymbol{B})$$

## Property (Determinant–Eigenvalue Relation)

The determinant of the correlation matrix is related to the eigenvalues as follows:

$$\det(\mathbf{R}) = \prod_{i=1}^{M} \lambda_i$$

Proof: Using $\mathbf{R} = \boldsymbol{Q}\boldsymbol{\Omega}\boldsymbol{Q}^H$ and the above,

$$
\begin{aligned}
\det(\mathbf{R}) &= \det(\boldsymbol{Q}\boldsymbol{\Omega}\boldsymbol{Q}^H) \\
&= \det(\mathbf{Q})\det(\mathbf{Q}^H)\det(\boldsymbol{\Omega}) = \det(\boldsymbol{\Omega}) = \prod_{i=1}^{M} \lambda_i
\end{aligned}
$$

## Property (Trace–Eigenvalue Relation)

The trace of the correlation matrix is related to the eigenvalues as follows:

$$\text{trace}(\mathbf{R}) = \sum_{i=1}^{M} \lambda_i$$

Proof: Note

$$
\begin{aligned}
\text{trace}(\mathbf{R}) &= \text{trace}(Q\Omega Q^H) \\
&= \text{trace}(\mathbf{Q}^H Q\Omega) \\
&= \text{trace}(\mathbf{\Omega}) \\
&= \sum_{i=1}^{M} \lambda_i
\end{aligned}
$$

QED

## Definition (Normal Matrix)

A complex square matrix $\mathbf{A}$ is a normal matrix if

$$\mathbf{A}^H \mathbf{A} = \mathbf{A} \mathbf{A}^H$$

That is, a matrix is normal if it commutes with its conjugate transpose.

Note

▶ All Hermitian symmetric matrices are normal

▶ Every matrix that can be diagonalized by the unitary transform is normal

## Definition (Condition Number)

The condition number reflects how numerically well–conditioned a problem is, i.e, a low condition number $\Rightarrow$ well–conditioned; a high condition number $\Rightarrow$ ill–conditioned.

## Definition (Condition Number for Linear Systems)

For a linear system

$$\mathbf{Ax} = \mathbf{b}$$

defined by a normal matrix $\mathbf{A}$, the condition number is

$$\chi(\mathbf{A}) = \frac{\lambda_{\max}}{\lambda_{\min}}$$

where $\lambda_{\max}$ and $\lambda_{\min}$ are the maximum/minimum eigenvalues of $\mathbf{A}$

Observations:

▶ Large eigenvalue spread $\Rightarrow$ ill–conditioned

▶ Small eigenvalue spread $\Rightarrow$ well–conditioned

# Outline

Eigen Analysis

Eigen Properties

SVD

PCA

# Matrix-Vector Multiplication

Example in 2D:

$$\mathbf{A} = \begin{bmatrix} 2 & 1 \\ -1 & 1 \end{bmatrix} \qquad \mathbf{x} = \begin{bmatrix} 1 \\ 3 \end{bmatrix}$$

and,

$$\mathbf{y} = \mathbf{A}\mathbf{x} = \begin{bmatrix} 2 & 1 \\ -1 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 3 \end{bmatrix} = \begin{bmatrix} 5 \\ 2 \end{bmatrix}$$

What is the geometrical meaning of the matrix-vector multiplication?

# Matrix-Vector Multiplication

$$\mathbf{y} = \mathbf{A}\mathbf{x} = \begin{bmatrix} 2 & 1 \\ -1 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 3 \end{bmatrix} = \begin{bmatrix} 5 \\ 2 \end{bmatrix}$$



▶ Rotates the vector $\angle\theta$
▶ Stretches the vector

# Matrix-Vector Multiplication

To rotate **x** by an angle $\theta$, we pre-multiply by

$$\mathbf{A} = \begin{bmatrix} \cos\theta & -\sin\theta \\ \sin\theta & \cos\theta \end{bmatrix}$$

Stretch **x** by factor $\alpha$, pre-multiply by

$$\mathbf{A} = \begin{bmatrix} \alpha & 0 \\ 0 & \alpha \end{bmatrix}$$

# Matrix-Vector Multiplication

Consider the vectors $\mathbf{v}_1$ and $\mathbf{v}_2$ depicting a circle. What happens to the circle under matrix multiplication?



| 2-D Circle | 3-D Sphere | n-D Hypersphere |
|---|---|---|
| $\mathbf{A}[\mathbf{v}_1 \quad \mathbf{v}_2]$ | $\mathbf{A}[\mathbf{v}_1 \quad \mathbf{v}_2 \quad \mathbf{v}_3]$ | $\mathbf{A}[\mathbf{v}_1 \quad \cdots \quad \mathbf{v}_n]$ |
| $\mathbf{v}_i \in \mathbb{C}^2$ | $\mathbf{v}_i \in \mathbb{C}^3$ | $\mathbf{v}_i \in \mathbb{C}^n$ |

# Matrix-Vector Multiplication

What happens to the 2D circle under matrix multiplication?

$$[\mathbf{A}][\mathbf{v_1} \quad \mathbf{v_2}] = [\hat{\mathbf{u}}_1 \quad \hat{\mathbf{u}}_2] \begin{bmatrix} \sigma_1 & 0 \\ 0 & \sigma_2 \end{bmatrix}$$

$\mathbf{v_1}, \mathbf{v_2}$   $\mathbf{A} \in \mathbb{C}^{2 \times 2}$

$\mathbf{v_2}$   $\mathbf{v_1}$

$\sigma_2 \hat{\mathbf{u}}_2$  Principal axis

$\hat{\mathbf{u}}_2$
$\hat{\mathbf{u}}_1$

Singular values

$\sigma_1 \hat{\mathbf{u}}_1$

$\hat{\mathbf{u}}_1, \hat{\mathbf{u}}_2$   Unitary orthonormal vectors
$\sigma_1, \sigma_2$   "Stretching" constant

**Note:** Ortogonality holds since they are all rotated by the same angle.

# Matrix-Vector Multiplication

What happens to the n-D hyper-sphere under matrix multiplication?

$$\mathbf{A}\mathbf{v}_j = \sigma_j \widehat{\mathbf{u}}_j$$

$$\mathbf{v}_1, \mathbf{v}_2, \ldots, \mathbf{v}_n \quad \mathbf{A} \in \mathbb{C}^{m \times n} \quad [\mathbf{A}][\mathbf{v}_1 \quad \cdots \quad \mathbf{v}_n] = [\widehat{\mathbf{u}}_j \quad \cdots \quad \widehat{\mathbf{u}}_j]\begin{bmatrix} \sigma_1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \sigma_n \end{bmatrix}$$



Unitary orthonormal vectors
$\widehat{\mathbf{u}}_1, \widehat{\mathbf{u}}_2, \ldots, \widehat{\mathbf{u}}_n$

"Stretching" constant
$\sigma_1, \sigma_2, \ldots, \sigma_n$

n-dim hyper-ellipse

# n-dim Hyper-Sphere Mapping to n-dim Hyper-Ellipsoid

The mapping can be written as

$$\mathbf{A}\mathbf{v}_1 = \sigma_1 \hat{\mathbf{u}}_1$$
$$\vdots \qquad \vdots$$
$$\mathbf{A}\mathbf{v}_n = \sigma_j \hat{\mathbf{u}}_n$$

Expressed in matrix form as

$$\underbrace{\begin{bmatrix} \mathbf{A} \end{bmatrix}}_{\mathbf{A} \in \mathbb{C}^{m \times n}} \underbrace{[\mathbf{v}_1 \ \mathbf{v}_2 \ \dots \mathbf{v}_n]}_{\mathbf{V} \ \mathbb{C}^{n \times n}} = \underbrace{[\hat{\mathbf{u}}_1 \ \hat{\mathbf{u}}_2 \ \dots \hat{\mathbf{u}}_n]}_{\hat{\mathbf{U}} \ \mathbb{C}^{m \times n}} \underbrace{\begin{bmatrix} \sigma_1 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \sigma_n \end{bmatrix}}_{\hat{\mathbf{\Sigma}} \ \mathbb{C}^{n \times n}}$$

$$\mathbf{A}\mathbf{V} = \hat{\mathbf{U}}\hat{\mathbf{\Sigma}}$$

# n-dim Hyper-Sphere Mapping to n-dim Hyper-Ellipsoid

Let $\mathbf{v}_1, \ldots, \mathbf{v}_n$ be unitary orthonormal vectors, then $\mathbf{V} = [\mathbf{v}_1 \ \mathbf{v}_2 \ \ldots \ \mathbf{v}_n]$ is a unitary transformation matrix, that is

$$\mathbf{V}^{-1} = \mathbf{V}^H.$$

Let $\hat{\mathbf{u}}_1, \ldots, \hat{\mathbf{u}}_n$ be unitary orthonormal vectors, then $\hat{\mathbf{U}} = [\hat{\mathbf{u}}_1 \ \hat{\mathbf{u}}_2 \ \ldots \ \hat{\mathbf{u}}_n]$ is a unitary transformation matrix, that is

$$\mathbf{U}^{-1} = \hat{\mathbf{U}}^H.$$

# Reduced Singular Value Decomposition

The mapping is thus given by,

$$\mathbf{AV} = \hat{\mathbf{U}}\hat{\mathbf{\Sigma}}$$

Multiply both sides by $\mathbf{V}^{-1}$ we obtain:

$$\begin{aligned}
\mathbf{AVV}^{-1} &= \hat{\mathbf{U}}\hat{\mathbf{\Sigma}}\mathbf{V}^{-1} \\
\mathbf{AVV}^{H} &= \hat{\mathbf{U}}\hat{\mathbf{\Sigma}}\mathbf{V}^{H} \\
\mathbf{AI} &= \hat{\mathbf{U}}\hat{\mathbf{\Sigma}}\mathbf{V}^{H} \\
\mathbf{A} &= \hat{\mathbf{U}}\hat{\mathbf{\Sigma}}\mathbf{V}^{H}
\end{aligned}$$

where $\mathbf{\Sigma} = \text{diag}\left([\sigma_1, \sigma_2, \ldots, \sigma_n]\right)$, such that $\sigma_1 \geq \sigma_2 \geq \ldots \sigma_p \geq 0$.

# Singular Value Decomposition

▶ Reduced SVD



▶ SVD

# Theorem 1

Every matrix $\mathbf{A} \in \mathbb{C}^{m \times n}$ has a singular value decomposition (SVD).

▶ Singular values $\sigma_j$ are uniquely determined.

▶ If $\mathbf{A}$ is square $\sigma_j$ are distinct.

▶ $\mathbf{u}_j$ and $\mathbf{v}_j$ are also unique up to a complex sign. (unique if the complex sign is ignored)

## SVD calculation

Start with $\mathbf{A}^\mathsf{T}\mathbf{A}$:

$$
\begin{aligned}
\mathbf{A}^\mathsf{H}\mathbf{A} &= \left(\mathbf{U}\boldsymbol{\Sigma}\mathbf{V}^H\right)^\mathsf{H}\left(\mathbf{U}\boldsymbol{\Sigma}\mathbf{V}^\mathsf{H}\right) \\
&= \mathbf{V}\boldsymbol{\Sigma}\mathbf{U}^\mathsf{H}\mathbf{U}\boldsymbol{\Sigma}\mathbf{V}^H \\
\mathbf{A}^\mathsf{H}\mathbf{A}\mathbf{V} &= \mathbf{V}\boldsymbol{\Sigma}^2\mathbf{V}^\mathsf{H}\mathbf{V} \\
\mathbf{A}^\mathsf{H}\mathbf{A}\mathbf{V} &= \mathbf{V}\boldsymbol{\Sigma}^2
\end{aligned}
$$

Reduces to an eigenvalue decomposition problem of the form:

$$
\underbrace{\mathbf{A}^\mathsf{T}\mathbf{A}}_{\mathbf{B}}\mathbf{V} = \mathbf{V}\underbrace{\boldsymbol{\Sigma}^2}_{\boldsymbol{\Lambda}},
$$

where $\boldsymbol{\Lambda}$ is a diagonal matrix with the eigenvalues of $\mathbf{B}$ and $\mathbf{V}$ corresponds to the eigenvectors of $\mathbf{B}$.

## SVD calculation

How do we calculate $\mathbf{U}$:

$$
\begin{aligned}
\mathbf{A}\mathbf{A}^{\mathsf{H}} &= \left(\mathbf{U}\mathbf{\Sigma}\mathbf{V}^{\mathsf{H}}\right)\left(\mathbf{U}\mathbf{\Sigma}\mathbf{V}^{\mathsf{H}}\right)^{\mathsf{H}} \\
&= \mathbf{U}\mathbf{\Sigma}\mathbf{V}^{\mathsf{H}}\mathbf{V}\mathbf{\Sigma}\mathbf{U}^{\mathsf{H}} \\
\mathbf{A}\mathbf{A}^{\mathsf{H}}\mathbf{U} &= \mathbf{U}\mathbf{\Sigma}^{2}\mathbf{U}^{\mathsf{H}}\mathbf{U} \\
\underbrace{\mathbf{A}\mathbf{A}^{\mathsf{H}}}_{\mathbf{B}}\mathbf{U} &= \mathbf{U}\underbrace{\mathbf{\Sigma}^{2}}_{\mathbf{\Lambda}}
\end{aligned}
$$

Eigenvalue problem where $\mathbf{\Lambda}$ is a diagonal matrix with the eigenvalues of $\mathbf{B}$ and $\mathbf{U}$ corresponds to the eigenvectors of $\mathbf{B}$.

# Netflix Movie Challenge

- Dataset: $n = 17,770$ movies (columns) and $m = 480,189$ customers (rows).
- Customers rated movies on a scale from 1 to 5. Matrix is very sparse with "only" 100 million of the ratings present in the training set.
- Goal: Predict the ratings for unrated movies.



- (2006) "Cinematch" algorithm used by Netflix RMSE=0.9525 over a large test set.
- Competition started in 2006, winner should improve this RMSE by at least 10%.
- 2009 "Bellkor's Pragmatic Chaos," uses a combination of many statistical techniques to win.

# Movie Rating - A Solution

- Describe a movie as an array of factors, e.g. comedy, action…

- Describe each viewer using same factors, e.g. likes comedy, likes action, etc

- Rating based on match/mismatch

- More factors → better prediction

# Singular Value Decomposition Solution

Viewers rated movies on a scale from 1 to 5. 0 for movies that were not rated by the user.

▶ Each column $j$ is a different movie

▶ Each row $i$ is a different viewer

▶ Each element $a_{i,j}$ represents the rating of movie $j$ by viewer $i$

|  | Movie 1 | Movie 2 | Movie 3 | Movie 4 | Movie 5 |
|---|---|---|---|---|---|
| Viewer 1 | 0 | 1 | 0 | 0 | 5 |
| Viewer 2 | 4 | 2 | 0 | 0 | 0 |
| Viewer 3 | 0 | 0 | 3 | 3 | 0 |
| Viewer 4 | 4 | 2 | 0 | 0 | 0 |
| Viewer 5 | 0 | 0 | 0 | 0 | 5 |
| Viewer 6 | 0 | 0 | 3 | 3 | 0 |
| Viewer 7 | 1 | 0 | 0 | 0 | 4 |
| Viewer 8 | 2 | 1 | 0 | 0 | 4 |
| Viewer 9 | 1 | 0 | 0 | 0 | 4 |

$$\mathbf{A} = \begin{bmatrix} a_{1,1} & \cdots & \cdots & a_{1,n} \\ \vdots & \ddots & \ddots & \vdots \\ a_{m,1} & \cdots & \cdots & a_{m,n} \end{bmatrix}$$

**Goal:** Use SVD to predict unobserved data or the rating of a movie that hasn't come out yet.

# Singular Value Decomposition Solution

We want to classify Movies and Viewers:

$$Movies = \begin{cases} \text{Category 1} \\ \text{Category 2} \\ \text{Category 3} \\ \quad\vdots \end{cases}$$

Intuitively, if $Movie_1 \approx Movie_2$, these movies are similar (same category).

|          | Movie 1 | Movie 2 | Movie 3 | Movie 4 | Movie 5 |
|----------|---------|---------|---------|---------|---------|
| Viewer 1 | 0       | 1       | 0       | 0       | 5       |
| Viewer 2 | 4       | 2       | 0       | 0       | 0       |
| Viewer 3 | 0       | 0       | 3       | 3       | 0       |
| Viewer 4 | 4       | 2       | 0       | 0       | 0       |
| Viewer 5 | 0       | 0       | 0       | 0       | 5       |
| Viewer 6 | 0       | 0       | 3       | 3       | 0       |
| Viewer 7 | 1       | 0       | 0       | 0       | 4       |
| Viewer 8 | 2       | 1       | 0       | 0       | 4       |
| Viewer 9 | 1       | 0       | 0       | 0       | 4       |

Categories are determined by matrix $\mathbf{A}$ and SVD algorithm.

# Singular Value Decomposition Solution

Now, consider that each movie belongs to more than one category e.g. half comedy and half action. This can be written as:

$$Movie_j = v_1\mathsf{Cat1} + v_2\mathsf{Cat2} + \cdots + v_m\mathsf{Catn}$$
$$\text{s.t.} ||\mathbf{v}||_2 = 1$$

where the set of categories $\{\mathsf{Cat}j \in \mathbb{R}^{n \times 1}\}$ forms an orthonormal basis.

$$\mathsf{Cat} = \begin{bmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \ddots & 0 \\ 0 & 0 & \cdots & 1 \end{bmatrix}_{n \times n}$$

|  | Movie 1 | Movie 2 | Movie 3 | Movie 4 | Movie 5 |
|---|---|---|---|---|---|
| Viewer 1 | 0 | 1 | 0 | 0 | 5 |
| Viewer 2 | 4 | 2 | 0 | 0 | 0 |
| Viewer 3 | 0 | 0 | 3 | 3 | 0 |
| Viewer 4 | 4 | 2 | 0 | 0 | 0 |
| Viewer 5 | 0 | 0 | 0 | 0 | 5 |
| Viewer 6 | 0 | 0 | 3 | 3 | 0 |
| Viewer 7 | 1 | 0 | 0 | 0 | 4 |
| Viewer 8 | 2 | 1 | 0 | 0 | 4 |
| Viewer 9 | 1 | 0 | 0 | 0 | 4 |

# Singular Value Decomposition Solution

In the case of $Viewers$, we use the same $Movies$' categories:

$$Movies = \left\{ \begin{array}{c} \text{Category 1} \\ \text{Category 2} \\ \text{Category 3} \\ \vdots \end{array} \right\} = Viewers.$$

| | Movie 1 | Movie 2 | Movie 3 | Movie 4 | Movie 5 |
|---|---|---|---|---|---|
| Viewer 1 | 0 | 1 | 0 | 0 | 5 |
| Viewer 2 | 4 | 2 | 0 | 0 | 0 |
| Viewer 3 | 0 | 0 | 3 | 3 | 0 |
| Viewer 4 | 4 | 2 | 0 | 0 | 0 |
| Viewer 5 | 0 | 0 | 0 | 0 | 5 |
| Viewer 6 | 0 | 0 | 3 | 3 | 0 |
| Viewer 7 | 1 | 0 | 0 | 0 | 4 |
| Viewer 8 | 2 | 1 | 0 | 0 | 4 |
| Viewer 9 | 1 | 0 | 0 | 0 | 4 |

E.g. a viewer that loves comedy is represented with the same unit vector of the comedy category movies ($\text{Cat}i \in \mathbb{R}^{1 \times n}$) .

Each $Viewer$ is represented as:

$$Viewer_i = \quad u_1\text{Cat}1 + u_2\text{Cat}2 + \cdots + u_n\text{Cat}n$$
$$\text{s.t.} ||\mathbf{u}||_2 = 1$$

# Singular Value Decomposition Solution

If $m > n$ i.e # of Viewers $>$ # of Movies, each $Viewer$ is represented as:

$$Viewer_i = u_1\text{Cat1} + u_2\text{Cat2} + \cdots + u_n\text{Catn} + \cdots + u_m\text{Catm}$$
$$\text{s.t.}||\mathbf{u}||_2 = 1$$

where $\text{Cat}i \in \mathbb{R}^{1\times m}$. Thus, useless categories vectors with zero rating value are added.

# Singular Value Decomposition Solution

From Theorem 1:
There exist a unique decomposition into categories. Every matrix $\mathbf{A} \in \mathbb{C}^{m \times n}$ can be factorized as $\mathbf{A} = \hat{\mathbf{U}} \mathbf{\Sigma} \mathbf{V}^H$ where:

# Singular Value Decomposition Solution

We have more viewers than movies:



New categories are created. The new vectors are still unit vectors orthonormal to all the basis vectors but the ratings of these useless categories are zero.

Note: consider reduced SVD i.e. consider only useful categories.

# Singular Value Decomposition Solution



$$\mathbf{A} = \hat{\mathbf{U}} \; \boldsymbol{\Sigma} \; \mathbf{V}^{\mathrm{H}}$$

$\mathbb{C}^{m \times n}$   $\mathbb{C}^{m \times n}$   $\mathbb{C}^{n \times n}$   $\mathbb{C}^{n \times n}$

▶ Each row vector ($\mathbf{u}_i$) in $\hat{\mathbf{U}}$ represents the taste of a $Viewer_i$ on the corresponding categories.

$$\hat{\mathbf{U}} = \begin{bmatrix} u_{1,1} & \cdots & \cdots & u_{1,n} \\ \vdots & \ddots & \ddots & \vdots \\ u_{m,1} & \cdots & \cdots & u_{m,n} \end{bmatrix} = \begin{bmatrix} \mathbf{u}_1 \\ \mathbf{u}_2 \\ \vdots \\ \mathbf{u}_m \end{bmatrix}$$

# Singular Value Decomposition Solution



▶ Each column ($\mathbf{v}_j$) in $\mathbf{V}^H$ represents the content of a $Movie_j$ on the corresponding categories.

$$\mathbf{V}^H \;=\; \begin{bmatrix} v_{1,1} & \cdots & \cdots & v_{1,n} \\ \vdots & \ddots & \ddots & \vdots \\ v_{n,1} & \cdots & \cdots & v_{n,n} \end{bmatrix} = \begin{bmatrix} \mathbf{v}_1 & \mathbf{v}_2 & \ldots & \mathbf{v}_n \end{bmatrix}$$

# Singular Value Decomposition Solution



▶ Each singular value $\sigma_{ii}$ in $\Sigma$ computes how a viewer of category $i$ rates a movie of the same category $i$.

$$\Sigma = \begin{bmatrix} \sigma_{1,1} & 0 & \cdots & 0 \\ \vdots & \sigma_{2,2} & \ddots & \vdots \\ 0 & 0 & \ldots & \sigma_{n,n} \end{bmatrix}$$

# Singular Value Decomposition Solution

The representation of each movie can be obtained by

$$
\begin{aligned}
Movie_j &= v_{1,j}\mathsf{Cat1} + v_{2,j}\mathsf{Cat2} + \cdots + v_{n,j}\mathsf{Catn} \qquad \text{s.t.}\,||\mathbf{v}_j||_2 = 1 \\
&= v_{1,j}\begin{bmatrix} \sqrt{\sigma_{1,1}} \\ 0 \\ \vdots \\ 0 \end{bmatrix} + v_{2,j}\begin{bmatrix} 0 \\ \sqrt{\sigma_{2,2}} \\ \vdots \\ 0 \end{bmatrix} + \cdots + v_{n,j}\begin{bmatrix} 0 \\ 0 \\ \vdots \\ \sqrt{\sigma_{n,n}} \end{bmatrix} \\
&= \sqrt{\mathbf{\Sigma}}\mathbf{v}_j \qquad \in \mathbb{C}^{n\times 1}
\end{aligned}
$$

# Singular Value Decomposition Solution

The representation of each viewer can be obtained by

$$
\begin{aligned}
Viewer_i &= u_{i,1}\mathsf{Cat1} + u_{i,2}\mathsf{Cat2} + \cdots + u_{i,n}\mathsf{Catn} + \cdots + u_{i,m}\mathsf{Catm} \\
&\qquad \text{s.t.} \|\mathbf{u}_i\|_2 = 1, \qquad \mathsf{Cat}j = 0 \text{ for } j > n \rightarrow \text{useless categories} \\
&= u_{i,1}\begin{bmatrix} \sqrt{\sigma_{1,1}} \\ 0 \\ \vdots \\ 0 \end{bmatrix}^H + u_{i,2}\begin{bmatrix} 0 \\ \sqrt{\sigma_{2,2}} \\ \vdots \\ 0 \end{bmatrix}^H + \cdots + u_{i,n}\begin{bmatrix} 0 \\ 0 \\ \vdots \\ \sqrt{\sigma_{n,n}} \end{bmatrix}^H \\
&= \mathbf{u}_i\sqrt{\mathbf{\Sigma}}^H \qquad \in \mathbb{C}^{1 \times n}
\end{aligned}
$$

# Singular Value Decomposition Solution

Given the decomposition of a movie and a viewer, the rating is estimated by:

$$
\begin{aligned}
Viewer_i Movie_j &= u_{i,1}v_{1,j}\sigma_{1,1} + u_{i,2}v_{2,j}\sigma_{2,2} + \cdots + u_{i,n}v_{n,j}\sigma_{n,n} \\
&= (\mathbf{u}_i\sqrt{\mathbf{\Sigma}}^H)(\sqrt{\mathbf{\Sigma}}\mathbf{v}_j) \\
&= \mathbf{u}_i\mathbf{\Sigma}\mathbf{v}_i
\end{aligned}
$$

# Singular Value Decomposition - Example

Considering the rating from 60 viewers to 16 movies of 4 different genres(action, romance, sci-fi, comedy), we generate $\mathbf{A} \in \mathbb{R}^{60 \times 16}$

▶ Viewers rated movies on a scale from 1 to 5, 0 for movies that were not rated by the user.

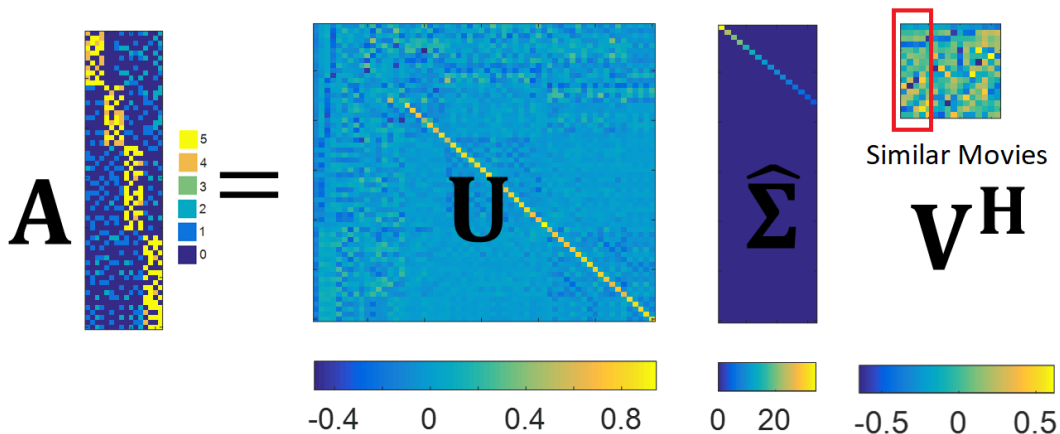▶ Observe the same 4 categories of viewers.

$\mathbf{A} =$

# Singular Value Decomposition - Example

# Singular Value Decomposition - Example

# Singular Value Decomposition - Example



Similar Movies

# Singular Value Decomposition - Example

To estimate not rated movies (zero entries in **A**), we use additional information: **A** is known to be low-rank or approximately low-rank.

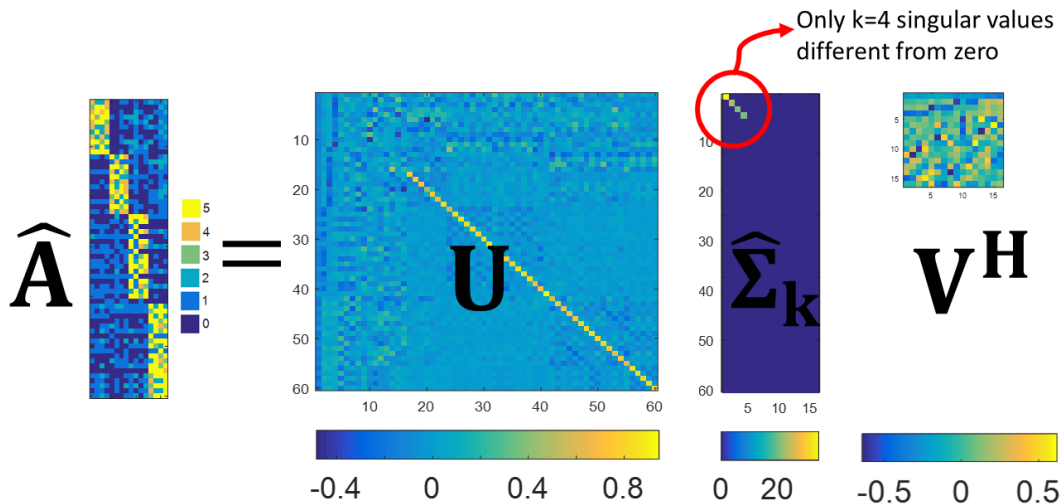Thus, we are going to use the k-rank approximation of the matrix **A** that is:

$$\hat{\mathbf{A}} = \mathbf{U}\hat{\Sigma}_k\mathbf{V}^H$$

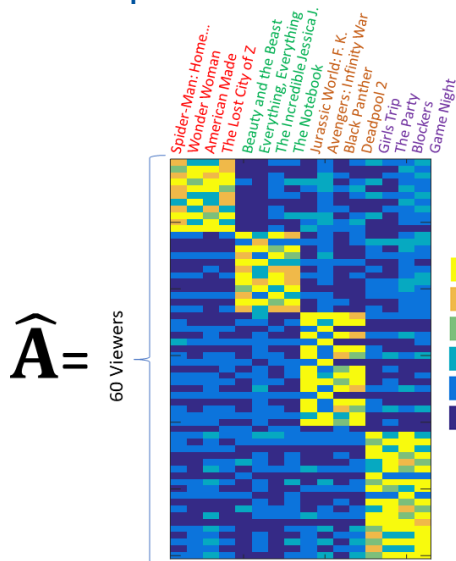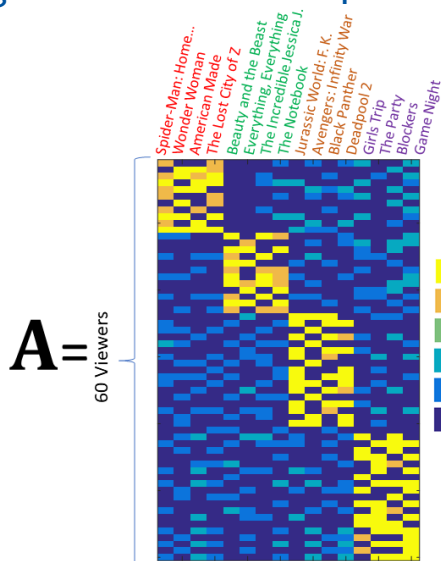where $\hat{\Sigma}_k$ has all but the first $k$ singular values $\sigma_{ii}$ set to zero.

The ratings different from zero in **A** are set to its original value.

Note: The ratings matrix **A** is expected to be low-rank since user preferences can be described by a few categories ($k$), such as the movie genres.

# Singular Value Decomposition - Example



Only k=4 singular values different from zero

$$\widehat{\mathbf{A}} = \mathbf{U} \; \widehat{\mathbf{\Sigma}}_{\mathbf{k}} \; \mathbf{V}^{\mathbf{H}}$$
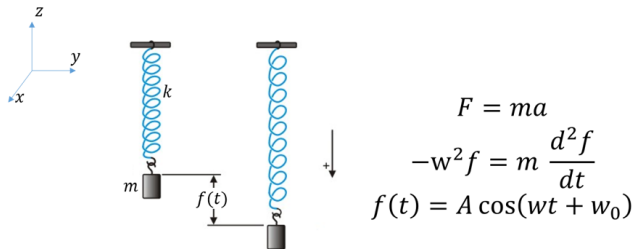
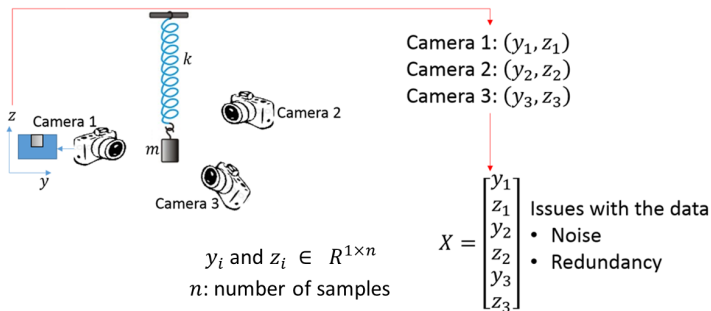# Singular Value Decomposition - Example

# Principal Component Analysis (PCA)

▶ Simple, method for extracting relevant information from confusing data sets.

▶ How to reduce a complex data set to a lower dimension?

▶ Consider a mass attached to a spring which oscillates as shown below.



$$F = ma$$

$$-w^2 f = m \, \frac{d^2 f}{dt}$$
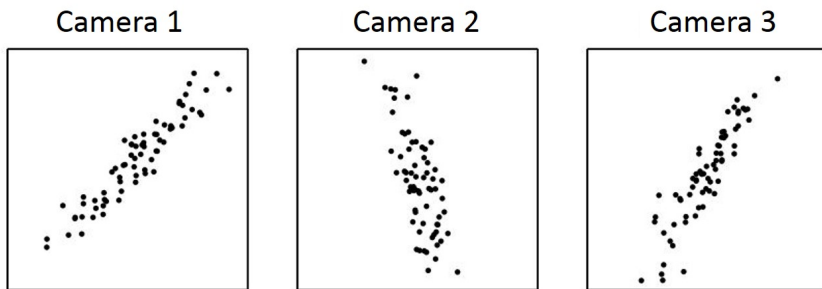
$$f(t) = A\cos(wt + w_0)$$

What if we did not know that $F = ma$?

# PCA - Motivation: Toy example

▶ Since we live in a 3D world → use three cameras to capture data from the system.

▶ No information about the real x,y, and z axes → camera positions are chosen arbitrarily.

▶ How do we get from this data set to a simple equation of $z$ ?



Camera 1: $(y_1, z_1)$
Camera 2: $(y_2, z_2)$
Camera 3: $(y_3, z_3)$

$$X = \begin{bmatrix} y_1 \\ z_1 \\ y_2 \\ z_2 \\ y_3 \\ z_3 \end{bmatrix}$$

Issues with the data
• Noise
• Redundancy

$y_i$ and $z_i \in R^{1 \times n}$
$n$: number of samples

# PCA - Motivation: Toy example

▶ Three cameras give redundant information.

▶ Only one camera at a specific angle necessary to describe the system behavior.

▶ PCA is used to avoid redundancy.



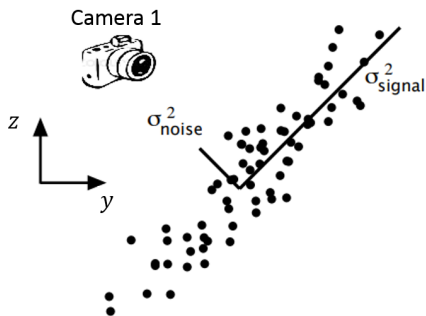Camera 1      Camera 2      Camera 3

# Change of Basis

▶ PCA: Is there another basis, which is a linear combination of the original basis, that best respresents the data set?

▶ Let $\mathbf{X}$ be the original data set, where each column is a single measurements set.

▶ Let $\mathbf{Y}$ be a linear transformation by $\mathbf{P}$, i.e. $\mathbf{Y} = \mathbf{PX}$, where $\mathbf{X} = [\mathbf{x}_1 | \ldots | \mathbf{x}_n]$ and $\mathbf{x}_i \in \mathbb{R}^{m \times 1}$ represents a sampled vector.

**Implications:**

▶ Geometrically $\mathbf{P}$ is a rotation and a stretch which transforms $\mathbf{X}$ into $\mathbf{Y}$.

▶ The rows of $\mathbf{P}$, $\{\mathbf{p}_1, \ldots, \mathbf{p}_m\}$ are a set of new basis vectors for expressing the columns of $\mathbf{X}$.
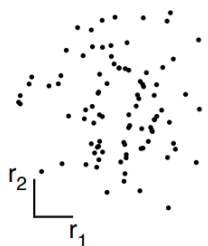
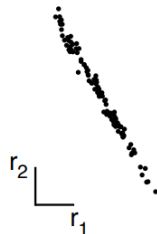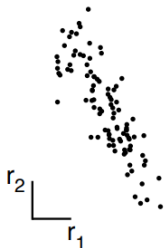What is the best way to re-express $\mathbf{X}$?, what is a good choice for $\mathbf{P}$?

## Noise



- ▶ Signal and noise variances are depicted as $\sigma^2_{\text{signal}}$ and $\sigma^2_{\text{noise}}$.
- ▶ The largest direction of variance is not along the natural basis but along the best-fit line.
- ▶ The directions with largest variances contain the dynamics of interest.
- ▶ Intuition: Find the direction indicated by $\sigma_{\text{signal}}$.

# Redundancy



low redundancy                                   high redundancy

- ▶ Figures depict possible plots between two arbitrary measurement types $r_1$ and $r_2$.
- ▶ Low redundancy $\rightarrow$ uncorrelated recordings
- ▶ High redundancy $\rightarrow$ correlated recordings, e.g. the sensors are too close or the measured variables are equivalent.
- ▶ If recordings are highly correlated it is not necessary to measure both of them.

# PCA - Basic concepts

Let $\mathbf{a} = [a_1, a_2, \ldots, a_n]$ and $\mathbf{b} = [b_1, b_2, \ldots, b_n]$ be two sets of measurements. Are they related?
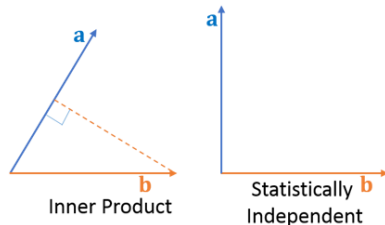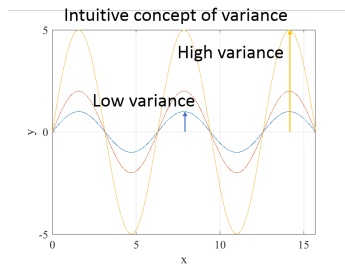
If the mean of $a$ and $b$ is zero, then:

▶ Variance: How large the change is in each vector.

$$\sigma_a^2 = \frac{1}{n}\mathbf{a}\mathbf{a}^T = \frac{1}{n}\sum_i a_i^2$$

$$\sigma_b^2 = \frac{1}{n}\mathbf{b}\mathbf{b}^T = \frac{1}{n}\sum_i b_i^2$$

▶ Covariance: Statistical relationship between data in $\mathbf{a}$ and $\mathbf{b}$.

$$\sigma_{ab}^2 = \frac{1}{n}\mathbf{a}\mathbf{b}^T = \frac{1}{n}\sum_i a_i b_i$$



Intuitive concept of variance

High variance

Low variance



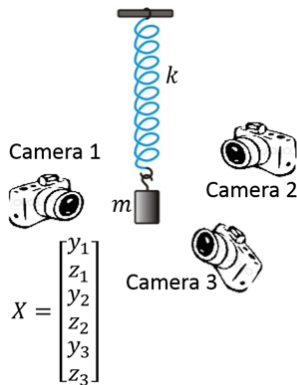Inner Product

Statistically Independent

## Variance and Covariance

Let $\mathbf{X}$ be defined as $\mathbf{X} = [\mathbf{x}_1^T | \ldots | \mathbf{x}_m^T]$, where $\mathbf{x}_i \in \mathbb{R}^{n \times 1}$ is a column vector that corresponds to all measurements of a particular type. Then the covariance matrix is defined as:

$$\mathbf{C_X} = \frac{1}{n} \mathbf{X} \mathbf{X}^T$$

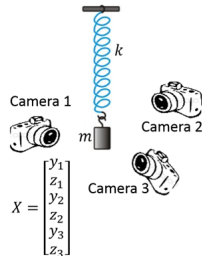The covariance values reflect the noise and redundacy in the measurements.



$$X = \begin{bmatrix} y_1 \\ z_1 \\ y_2 \\ z_2 \\ y_3 \\ z_3 \end{bmatrix}$$

# Variance and Covariance

Recall $\mathbf{C_X}$ is the covariance matrix of $\mathbf{X}$ defined as $\mathbf{C_X} = \frac{1}{n}\mathbf{X}\mathbf{X}^T$.



$$X = \begin{bmatrix} y_1 \\ z_1 \\ y_2 \\ z_2 \\ y_3 \\ z_3 \end{bmatrix}$$

▶ Covariance matrix in the spring example is $\mathbf{C_X} \in \mathbb{R}^{6\times6}$:

$$\mathbf{C_X} = \begin{bmatrix} \sigma_{y_1y_1}^2 & \sigma_{y_1z_1}^2 & \sigma_{y_1y_2}^2 & \sigma_{y_1z_2}^2 & \sigma_{y_1y_3}^2 & \sigma_{y_1z_3}^2 \\ \sigma_{z_1y_1}^2 & \sigma_{z_1z_1}^2 & \sigma_{z_1y_2}^2 & \sigma_{z_1z_2}^2 & \sigma_{z_1y_3}^2 & \sigma_{z_1z_3}^2 \\ \sigma_{y_2y_1}^2 & \sigma_{y_2z_1}^2 & \sigma_{y_2y_2}^2 & \sigma_{y_2z_2}^2 & \sigma_{y_2y_3}^2 & \sigma_{y_2z_3}^2 \\ \sigma_{z_2y_1}^2 & \sigma_{z_2z_1}^2 & \sigma_{z_2y_2}^2 & \sigma_{z_2z_2}^2 & \sigma_{z_2y_3}^2 & \sigma_{z_2z_3}^2 \\ \sigma_{y_3y_1}^2 & \sigma_{y_3z_1}^2 & \sigma_{y_3y_2}^2 & \sigma_{y_3z_2}^2 & \sigma_{y_3y_3}^2 & \sigma_{y_3z_3}^2 \\ \sigma_{z_3y_1}^2 & \sigma_{z_3z_1}^2 & \sigma_{z_3y_2}^2 & \sigma_{z_3z_2}^2 & \sigma_{z_3y_3}^2 & \sigma_{z_3z_3}^2 \end{bmatrix}$$

▶ Diagonal: Variance measures; Off-diagonal: covariance between all pairs.

▶ $\mathbf{C_X}$ is hermitian and symmetric, i.e. $\mathbf{C_X} = \mathbf{C_X}^{T*} = \mathbf{C_X}^T$.

## Covariance Matrix Interpretation

$$\mathbf{C_X} = \begin{bmatrix} \sigma^2_{y_1 y_1} & \sigma^2_{y_1 z_1} & \sigma^2_{y_1 y_2} & \sigma^2_{y_1 z_2} & \sigma^2_{y_1 y_3} & \sigma^2_{y_1 z_3} \\ \sigma^2_{z_1 y_1} & \sigma^2_{z_1 z_1} & \sigma^2_{z_1 y_2} & \sigma^2_{z_1 z_2} & \sigma^2_{z_1 y_3} & \sigma^2_{z_1 z_3} \\ \sigma^2_{y_2 y_1} & \sigma^2_{y_2 z_1} & \sigma^2_{y_2 y_2} & \sigma^2_{y_2 z_2} & \sigma^2_{y_2 y_3} & \sigma^2_{y_2 z_3} \\ \sigma^2_{z_2 y_1} & \sigma^2_{z_2 z_1} & \sigma^2_{z_2 y_2} & \sigma^2_{z_2 z_2} & \sigma^2_{z_2 y_3} & \sigma^2_{z_2 z_3} \\ \sigma^2_{y_3 y_1} & \sigma^2_{y_3 z_1} & \sigma^2_{y_3 y_2} & \sigma^2_{y_3 z_2} & \sigma^2_{y_3 y_3} & \sigma^2_{y_3 z_3} \\ \sigma^2_{z_3 y_1} & \sigma^2_{z_3 z_1} & \sigma^2_{z_3 y_2} & \sigma^2_{z_3 z_2} & \sigma^2_{z_3 y_3} & \sigma^2_{z_3 z_3} \end{bmatrix}$$

Off-diagonal terms
- ▶ If covariance is large then components are statistically dependent.
- ▶ If covariance is small then components are statistically independent.

Diagonal terms:
- ▶ If variance is large it contains a lot of information about the system.
- ▶ If variance is small it does not provide significant information about the system.

## PCA

Goal: Change basis such that the covariance matrix of the data is diagonal.

▶ If off-diagonal terms $\approx 0$, the redundancies are eliminated.

▶ Diagonal terms represent the variance of each component.

▶ Components with large variance are the most representative.



$$C_X =$$

Looks like the SVD!

# PCA and Eigenvalue Decomposition

How to solve the problem?

- ▶ Data set: $\mathbf{X} \in \mathbb{R}^{m \times n}$, where $m$ is the number of measurement types and $n$ is the number of samples.

- ▶ PCA : Find an orthonormal matrix $\mathbf{P}$ in $\mathbf{Y} = \mathbf{P}\mathbf{X}$ such that $\mathbf{C_Y} = \frac{1}{n}\mathbf{Y}\mathbf{Y}^T$ is a diagonal matrix.

- ▶ The rows of $\mathbf{P}$ are the principal components of $\mathbf{X}$

## PCA and Eigenvalue Decomposition

We begin rewriting $\mathbf{C_Y}$ in terms of the unknown variable.

$$
\begin{aligned}
\mathbf{C_Y} &= \frac{1}{n}\mathbf{YY}^T \\
&= \frac{1}{n}(\mathbf{PX})(\mathbf{PX})^T \\
&= \frac{1}{n}\mathbf{PXX}^T\mathbf{P}^T \\
&= \mathbf{P}\left(\frac{1}{n}\mathbf{XX}^T\right)\mathbf{P}^T \\
&= \mathbf{PC_XP}^T
\end{aligned}
$$

# PCA and Eigenvalue Decomposition

$\mathbf{C_X}$ can be diagonalized by an orthogonal matrix of its eigenvectors since it is a symmetric matrix. Let $\mathbf{P} = \mathbf{Q}^T$, where $\mathbf{Q}$ is a matrix with the eigenvectors of $\frac{1}{n}\mathbf{X}\mathbf{X}^T$, then:

$$
\begin{aligned}
\mathbf{C_Y} &= \mathbf{P}\mathbf{C_X}\mathbf{P}^T \\
&= \mathbf{P}\left(\mathbf{Q}\Omega\mathbf{Q}^T\right)\mathbf{P}^T \\
&= \mathbf{P}\left(\mathbf{P}^T\Omega\mathbf{P}\right)\mathbf{P}^T \\
&= \left(\mathbf{P}\mathbf{P}^{-1}\right)\Omega\left(\mathbf{P}\mathbf{P}^{-1}\right) \\
&= \Omega
\end{aligned}
$$

The transformation $\mathbf{Y} = \mathbf{P}\mathbf{X}$ diagonalizes the system. Covariance of $\mathbf{Y}$ is a diagonal matrix with the eigenvalues of $\frac{1}{n}\mathbf{X}\mathbf{X}^T$.

## PCA and SVD

The SVD of $\mathbf{X}$ is given by $\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$. Let $\mathbf{P} = \mathbf{U}^T$, then:

$$\mathbf{Y} = \mathbf{U}^T\mathbf{X},$$

The covariance matrix of $\mathbf{Y}$ is given by:

$$
\begin{aligned}
\mathbf{C_Y} &= \frac{1}{n}\mathbf{Y}\mathbf{Y}^T \\
&= \frac{1}{n}\mathbf{U}^T\mathbf{X}\mathbf{X}^T\mathbf{U} \\
&= \frac{1}{n}\mathbf{U}^T\mathbf{U}\mathbf{\Sigma}\mathbf{V}^T\mathbf{V}\mathbf{\Sigma}\mathbf{U}^T\mathbf{U} \\
&= \frac{1}{n}\mathbf{\Sigma}^2
\end{aligned}
$$

## PCA

▶ The transformation $\mathbf{Y} = \mathbf{U}^T\mathbf{X}$ diagonalized the system. Covariance of $\mathbf{Y}$ is a diagonal matrix with the squared singular values of $\mathbf{X}$ multiplied by a factor of $\frac{1}{n}$.

▶ It can be concluded that $\mathbf{\Sigma}^2 = \mathbf{\Omega}$, and $\sigma_i^2 = \lambda_i$.

▶ The principal components of the data matrix are given by $\mathbf{U}^T$.
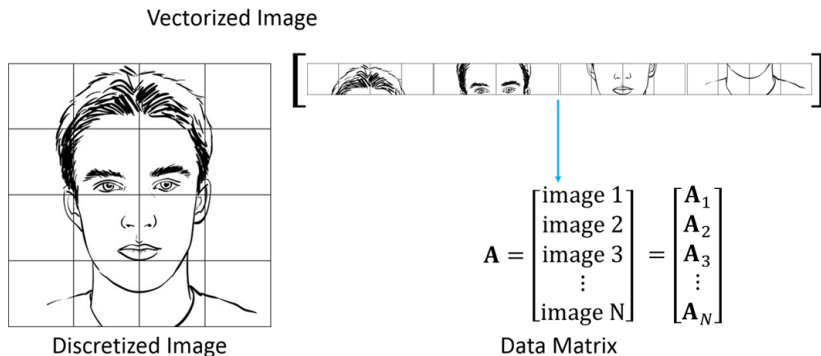
## Application: Face Recognition

▶ PCA in face recognition $\triangleq$ Eigenfaces

▶ Intuition: Figure out the correlation between the rows/ colums of **A** from the SVD.

$$\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T \tag{1}$$

▶ How important each direction is: $\mathbf{\Sigma}$

▶ Principal Directions: **U**

▶ How each individual component (row/column) projects onto the principal components: **V**.

# Data in Face Recognition

The data matrix is constructed by vectorizing the face images as shown below, i.e. $\mathbf{A} = [\mathbf{A}_1^T | \mathbf{A}_2^T | \ldots | \mathbf{A}_N^T]^T$. The matrix will be $N \times M$, where $N$ is the number of images in the data base and $M$ is the number of pixels of each image.



Vectorized Image

$$\mathbf{A} = \begin{bmatrix} \text{image 1} \\ \text{image 2} \\ \text{image 3} \\ \vdots \\ \text{image N} \end{bmatrix} = \begin{bmatrix} \mathbf{A}_1 \\ \mathbf{A}_2 \\ \mathbf{A}_3 \\ \vdots \\ \mathbf{A}_N \end{bmatrix}$$

Discretized Image

Data Matrix

# Example - Celebrity Images

Example, take 5 images of each celebrity: George Clooney, Bruce Willis, Margaret Thatcher and Matt Damon. In the example, $M = 240 * 160 = 38400$ and $N = 20$.



$$\mathbf{A} = \begin{bmatrix} -\!-\!-\!-\!-\!-\!- \text{ Image 1 } -\!-\!-\!-\!-\!-\!- \\ -\!-\!-\!-\!-\!-\!- \text{ Image 2 } -\!-\!-\!-\!-\!-\!- \\ -\!-\!-\!-\!-\!-\!- \text{ Image 3 } -\!-\!-\!-\!-\!-\!- \\ \\ \\ -\!-\!-\!-\!-\!-\!-\text{Image 20 } -\!-\!-\!-\!-\!-\!- \end{bmatrix}_{20 \times 38400}$$

## Average Faces

How do the average of the faces of these celebrities look like?

$$\bar{\mathbf{a}}_i = \frac{1}{5} \sum_{j=1}^{5} \mathbf{A}_j \quad \text{where} \quad \mathbf{A}_j \in \mathbb{R}^{1 \times M}$$
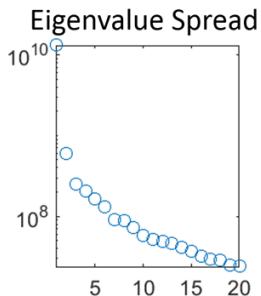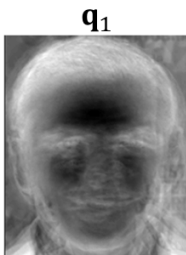
# Average Faces

What defines George Clooney's face?

▶ Data matrix $\mathbf{A} \in \mathbb{R}^{N \times M}$ with the images of the example.

▶ Compute the correlation matrix of the features of the dataset, i.e. the pixels.

▶ The correlation matrix is $\mathbf{C} = \mathbf{A}^T \mathbf{A} \in \mathbb{R}^{M \times M}$, here $M = 38400$.

▶ High correlation values $\rightarrow$ everybody has eyes, a nose and a mouth.

▶ Correlations between images of the same person will be higher.



Average Face

# Eigendecomposition

- ▶ Obtain the eigenvalue decomposition of $\mathbf{C} = \mathbf{A}^T\mathbf{A}$. That is $\mathbf{C} = \mathbf{Q\Omega Q}^{-1}$.

- ▶ First eigenvectors $\mathbf{q}_i \in \mathbb{R}^{M \times 1}$ are called the principal components (eigenfaces).

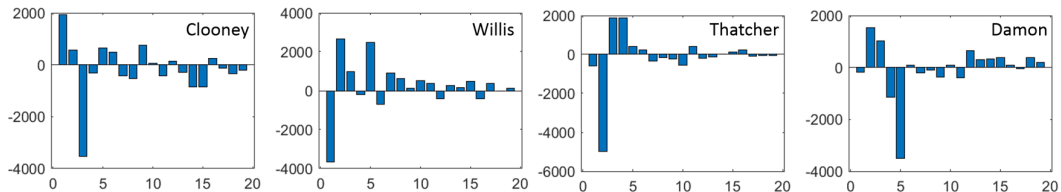- ▶ One can reconstruct each face as a weighted sum of the eigenvectors.



$\mathbf{q}_1$    $\mathbf{q}_2$    $\mathbf{q}_3$

$\mathbf{q}_4$    $\mathbf{q}_5$    Eigenvalue Spread

## Representing Faces onto Basis

Each face $\mathbf{A}_i \in \mathbb{R}^{1\times M}$ in the data set $\mathbf{A} = [\mathbf{A}_1^T | \mathbf{A}_2^T | \ldots | \mathbf{A}_N^T]^T$, can be represented as a linear combination of the best K eigenvectors:

$$\mathbf{A}_i^T = \sum_{j=1}^{K} w_j \mathbf{q}_j, \text{ where } w_j = \mathbf{q}_j^T \mathbf{A}_i^T \tag{2}$$



| $K = 1$ | $K = 5$ | $K = 10$ | $K = 15$ | $K = 20$ |

# Projection of the Average faces into the K=20 largest Eigenvectors

- ▶ **Q** is $M \times M$, from now on let **V** be the matrix formed by the first K=20 eigenvectors, i.e. $\mathbf{V} \in \mathbb{R}^{M \times K}$.
- ▶ Project the average faces $\bar{\mathbf{a}}_i \in \mathbb{R}^{1 \times M}$ onto the reduced eigenvector space, i.e. $\mathbf{p}_{\bar{\mathbf{a}}_i} = \bar{\mathbf{a}}_i \mathbf{V} \in \mathbb{R}^{1 \times K}$
- ▶ Projections for each face are characteristic of each average face and could be used for classification purposes.

# Projection of new images

▶ Test set: New image of Margaret Thatcher, Maryl Streep as Margaret Thatcher in "The Iron Lady", Betty White.

▶ Project test images onto eigenvector space, $\mathbf{p} = \mathbf{x}\mathbf{V} \in \mathbb{R}^{1 \times K}$, where $\mathbf{x} \in \mathbb{R}^{1 \times M}$ is the new vectorized image and $\mathbf{V} \in \mathbb{R}^{M \times K}$ is the matrix with the first 20 eigenvectors of the database.

▶ Reconstruct images as $\hat{\mathbf{x}} = \mathbf{V}\mathbf{p}^{T}$.

▶ Error defined as the difference between the projection of the new image and the projection of the original Margaret Thatcher images $\mathbf{o}_j\mathbf{V}$ where $j = 1, \ldots, 5$, that is
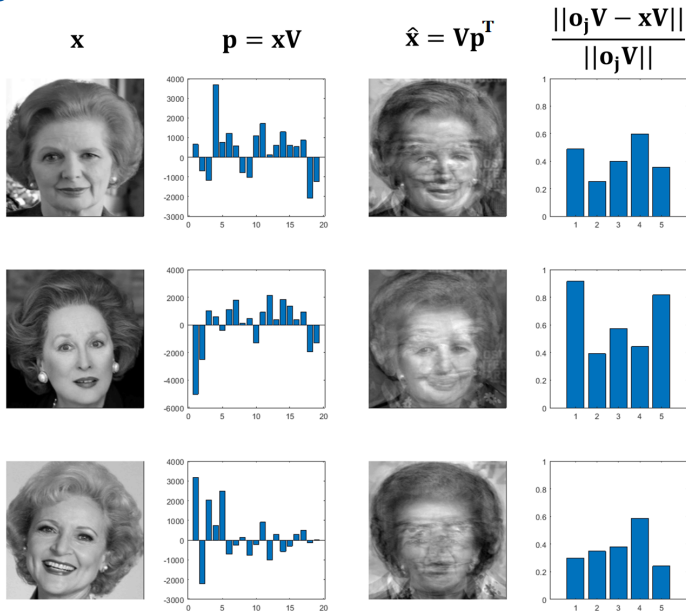
$$E_j = \frac{||\mathbf{o}_j\mathbf{V} - \mathbf{x}\mathbf{V}||}{||\mathbf{o}_j\mathbf{V}||},$$

where $\mathbf{o}_j$ are the original images of the database, in this case the 5 images of Margareth Thatcher.
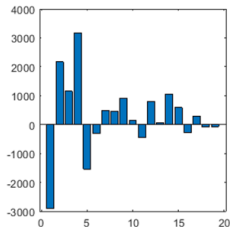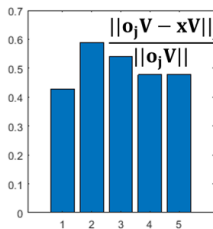
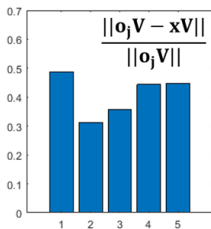# Projection of new images

Image depicts, from left to right

- ▶ Test images.
- ▶ Projection of the test images onto the eigenvector space $\mathbf{p} = \mathbf{xV}$.
- ▶ Reconstructed images using the first 20 eigenvectors of the database $\hat{\mathbf{x}} = \mathbf{Vp}^T$.
- ▶ Error of the projection with respect to each original Margareth Thatcher Image $\mathbf{o}_j$ for $j = 1, ..., 5$.
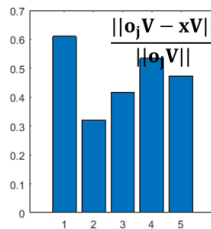
| $\mathbf{x}$ | $\mathbf{p} = \mathbf{xV}$ | $\hat{\mathbf{x}} = \mathbf{Vp^T}$ | $\dfrac{||\mathbf{o_j V} - \mathbf{xV}||}{||\mathbf{o_j V}||}$ |

# Projection of new images



$$\mathbf{x} \qquad \mathbf{p} = \mathbf{xV} \qquad \hat{\mathbf{x}} = \mathbf{Vp^T}$$

Clooney — Willis — Thatcher — Damon

$$\frac{||\mathbf{o_jV} - \mathbf{xV}||}{||\mathbf{o_jV}||}$$