# Principal Component Analysis

Aryan Mokhtari, Santiago Paternain, and Alejandro Ribeiro
Dept. of Electrical and Systems Engineering
University of Pennsylvania
aribeiro@seas.upenn.edu
`http://www.seas.upenn.edu/users/~aribeiro/`

March 29, 2022

The Discrete Fourier Transform with Unitary Matrices

Stochastic signals

Principal Component Analysis (PCA) transform

Dimensionality reduction

Principal Components

Face recognition

▶ We are ready to write the DFT with better tools $\Rightarrow$ Which will lead to abstract generalizations

▶ Write the signal $x$ and the complex exponential $e_{kN}$ as vectors in $\mathbb{R}^N$ $\Rightarrow$ Call them x and $e_{kN}$

$$x = \begin{bmatrix} x(0) \\ x(1) \\ \vdots \\ x(N-1) \end{bmatrix} \quad \Rightarrow \quad x^H = (x^*)^T = \begin{bmatrix} x(0), & x(1), & \ldots, & x(N-1) \end{bmatrix}$$

$$e_{kN} = \frac{1}{\sqrt{N}} \begin{bmatrix} e^{j2\pi k0/N} \\ e^{j2\pi k1/N} \\ \vdots \\ e^{j2\pi k(N-1)/N} \end{bmatrix} \quad \Rightarrow \quad e_{kN}^H = (e_{kN}^*)^T = \frac{1}{\sqrt{N}} \begin{bmatrix} e^{-j2\pi k0/N}, & e^{-j2\pi k1/N}, & \ldots, & e^{-j2\pi k(N-1)/N} \end{bmatrix}$$

▶ We can now rewrite the DFT as $\Rightarrow y(k) = \langle x, e_{kN} \rangle = e_{kN}^H x = \frac{1}{\sqrt{N}} \sum_{n=0}^{N-1} x(n) e^{-j2\pi kn/N}$

▶ The $k$th DFT component $\tilde{x}(k)$ is the product of exponential $e_{kN}^H$ with signal $x$ $\Rightarrow y(k) = e_{kN}^H x$

▶ Define the DFT vector $y$ as a stack of all $N$ DFT components and stack individual definitions

$$
y = \begin{bmatrix} y(0) \\ y(1) \\ \vdots \\ y(N-1) \end{bmatrix} = \begin{bmatrix} e_{0N}^H x \\ e_{1N}^H x \\ \vdots \\ e_{(N-1)N}^H x \end{bmatrix} = \begin{bmatrix} e_{0N}^H \\ e_{1N}^H \\ \vdots \\ e_{(N-1)N}^H \end{bmatrix} x = F^H x
$$

▶ Where in the last equality we defined the Fourier matrix Hermitian $F^H$ to write $y = F^H x$

▶ Each row of the DFT Hermitian is the Hermitian of a complex exponential of a different frequency

$$
F^H = \begin{bmatrix} e_{0N}^H \\ e_{1N}^H \\ \vdots \\ e_{kN}^H \\ \vdots \\ e_{(N-1)N}^H \end{bmatrix} = \frac{1}{\sqrt{N}} \begin{bmatrix} 1 & 1 & \cdots & 1 & & \\ 1 & e^{-j2\pi(1)(1)/N} & \cdots & e^{-j2\pi(1)(n)/N} & \cdots & e^{-j2\pi(1)(N-1)/N} \\ \vdots & \vdots & & \vdots & & \\ 1 & e^{-j2\pi(k)(1)/N} & \cdots & e^{-j2\pi(k)(n)/N} & \cdots & e^{-j2\pi(k)(N-1)/N} \\ \vdots & \vdots & & \vdots & & \\ 1 & e^{-j2\pi(N-1)(1)/N} & \cdots & e^{-j2\pi(N-1)(n)/N} & \cdots & e^{-j2\pi(N-1)(N-1)/N} \end{bmatrix}
$$

▶ Each row of DFT Hermitian corresponds to a given frequency. We sweep different time indexes

▶ Each column of DFT Hermitian corresponds to a given time index. We sweep different frequencies

▶ Each row of a matrix-vector product entails sweeping the corresponding row of the matrix



$$e^{-j2\pi(k)(0)/N}x(0)$$

$$e^{-j2\pi(k)(n)/N}x(n)$$

$$e^{-j2\pi(k)(N-1)/N}x(N-1)$$

$$\begin{bmatrix} x(0) \\ \cdot \\ x(n) \\ \cdot \\ x(N-1) \end{bmatrix} = \mathsf{x}$$

$$\mathsf{F}^H = \begin{bmatrix} e^{-j2\pi(0)(0)/N} & \cdot & e^{-j2\pi(0)(n)/N} & \cdot & e^{-j2\pi(0)(N-1)/N} \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ e^{-j2\pi(k)(0)/N} & \cdot & e^{-j2\pi(k)(n)/N} & \cdot & e^{-j2\pi(k)(N-1)/N} \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ e^{-j2\pi(N-1)(0)/N} & \cdot & e^{-j2\pi(N-1)(n)/N} & \cdot & e^{-j2\pi(N-1)(N-1)/N} \end{bmatrix}$$

$$\begin{bmatrix} y(0) \\ \cdot \\ y(k) \\ \cdot \\ y(N-1) \end{bmatrix} = \mathsf{y} = \mathsf{F}^H\mathsf{x}$$

▶ The Hermitian of the DFT Hermitian matrix is the DFT matrix $\Rightarrow \mathsf{F} = (\mathsf{F}^H)^H$

$$\mathsf{F}^H = \begin{bmatrix} \mathsf{e}_{0N}^H \\ \mathsf{e}_{1N}^H \\ \vdots \\ \mathsf{e}_{kN}^H \\ \vdots \\ \mathsf{e}_{(N-1)N}^H \end{bmatrix} \quad \Rightarrow \quad \mathsf{F} = \begin{bmatrix} \mathsf{e}_{0N} & \mathsf{e}_{1N} & \cdots & \mathsf{e}_{kN} & \cdots & \mathsf{e}_{(N-1)N} \end{bmatrix}$$

▶ The conjugate of the $k$th row of the Hermitian $\mathsf{F}^H$ becomes the $k$th column of the DFT matrix F

- Each column of the DFT matrix is a complex exponential of a different frequency

$$F = \begin{bmatrix} e_{0N} & e_{1N} & \cdots & e_{kN} & \cdots & e_{(N-1)N} \end{bmatrix}$$

$$F = \frac{1}{\sqrt{N}} \begin{bmatrix} 1 & 1 & \cdots & 1 & \cdots & 1 \\ 1 & e^{-j2\pi(1)(1)/N} & \cdots & e^{-j2\pi(k)(1)/N} & \cdots & e^{-j2\pi(N-1)(1)/N} \\ \vdots & \vdots & & \vdots & & \vdots \\ 1 & e^{-j2\pi(1)(n)/N} & \cdots & e^{-j2\pi(k)(n)/N} & \cdots & e^{-j2\pi(N-1)(n)/N} \\ \vdots & \vdots & & \vdots & & \vdots \\ 1 & e^{-j2\pi(1)(N-1)/N} & \cdots & e^{-j2\pi(k)(N-1)/N} & \cdots & e^{-j2\pi(N-1)(N-1)/N} \end{bmatrix}$$

- Each column corresponds to a given frequency. Each row corresponds to a given time index.

▶ The product of the Hermitian Fourier matrix $F^H$ and the Fourier matrix F is given by

$$
\begin{bmatrix} e_{0N}^H \\ e_{1N}^H \\ \vdots \\ e_{kN}^H \\ \vdots \\ e_{(N-1)N}^H \end{bmatrix}
\begin{bmatrix} e_{0N} & e_{1N} & \cdots & e_{kN} & \cdots & e_{(N-1)N} \end{bmatrix}
$$

$$
= \begin{bmatrix}
e_{0N}^H e_{0N} & e_{0N}^H e_{1N} & \cdots & e_{0N}^H e_{kN} & \cdots & e_{0N}^H e_{(N-1)N} \\
e_{1N}^H e_{0N} & e_{1N}^H e_{1N} & \cdots & e_{1N}^H e_{kN} & \cdots & e_{1N}^H e_{(N-1)N} \\
\vdots & \vdots & & \vdots & & \vdots \\
e_{kN}^H e_{0N} & e_{kN}^H e_{1N} & \cdots & e_{kN}^H e_{kN} & \cdots & e_{kN}^H e_{(N-1)N} \\
\vdots & \vdots & & \vdots & & \vdots \\
e_{(N-1)N}^H e_{0N} & e_{(N-1)N}^H e_{1N} & \cdots & e_{(N-1)N}^H e_{kN} & \cdots & e_{(N-1)N}^H e_{(N-1)N}
\end{bmatrix} = F^H F
$$

▶ The inner products in the diagonal are one and the inner products outside the diagonal are zero

▶ The product of the Hermitian Fourier matrix $F^H$ and the Fourier matrix $F$ is given by

$$
\begin{bmatrix} e_{0N}^H \\ e_{1N}^H \\ \vdots \\ e_{kN}^H \\ \vdots \\ e_{(N-1)N}^H \end{bmatrix}
\begin{bmatrix}
e_{0N} & e_{1N} & \cdots & e_{kN} & \cdots & e_{(N-1)N}
\end{bmatrix}
=
\begin{bmatrix}
1 & 0 & \cdots & 0 & \cdots & 0 \\
0 & 1 & \cdots & 0 & \cdots & 0 \\
\vdots & \vdots & & \vdots & & \vdots \\
0 & 0 & \cdots & 1 & \cdots & 0 \\
\vdots & \vdots & & \vdots & & \vdots \\
0 & 0 & \cdots & 0 & \cdots & 1
\end{bmatrix}
= F^H F
$$

▶ The inner products in the diagonal are one and the inner products outside the diagonal are zero

▶ We have therefore proved the following fundamental theorem

**Theorem**
*The DFT matrix* F *and its Hermitian are inverses of each other* $\Rightarrow$ *We say they are unitary*

$$F^H F \; = \; I \; = \; FF^H$$

▶ This follows from, and is equivalent to, the orthonormality of complex exponentials

$\Rightarrow$ Orthonormality, the professional way.

▶ We can also write the iDFT of y as a matrix product ⇒ In this case the product is x̃ = Fy



$$F = \begin{bmatrix} e^{j2\pi(0)(0)/N} & . & e^{j2\pi(k)(0)/N} & . & e^{j2\pi(N-1)(0)/N} \\ . & & . & & . \\ e^{j2\pi(0)(n)/N} & . & e^{j2\pi(k)(n)/N} & . & e^{j2\pi(N-1)(n)/N} \\ . & & . & & . \\ e^{j2\pi(0)(N-1)/N} & . & e^{j2\pi(k)(N-1)/N} & . & e^{j2\pi(N-1)(N-1)/N} \end{bmatrix} \begin{bmatrix} \tilde{x}(0) \\ . \\ \tilde{x}(n) \\ . \\ \tilde{x}(N-1) \end{bmatrix} = \tilde{x} = Fy$$

▶ When we proved theorems we had monkey steps and one smart step

$\Rightarrow$ That was orthonormality $\Rightarrow$ matrix F is unitary $\Rightarrow \mathsf{F}^H \mathsf{F} = \mathsf{I}$

Theorem
*The iDFT is, indeed, the inverse of the DFT*

Proof.

▶ Write $\tilde{\mathsf{x}} = \mathsf{F}\mathsf{X}$ and $\mathsf{X} = \mathsf{F}^H \mathsf{x}$ and exploit fact that F is unitary

$$\tilde{\mathsf{x}} \;=\; \mathsf{F}\mathsf{X} \;=\; \mathsf{F}\mathsf{F}^H \mathsf{x} \;=\; \mathsf{I}\mathsf{x} \;=\; \mathsf{x} \qquad\qquad \square$$

▶ Actually, this theorem would be true for any transform pair

$$\mathsf{X} = \mathsf{T}^H \mathsf{x} \qquad \Longleftrightarrow \qquad \tilde{\mathsf{x}} = \mathsf{T}\mathsf{X}$$

▶ As long as the transform matrix T is unitary $\Rightarrow \mathsf{T}^H \mathsf{T} = \mathsf{I}$

Theorem

*The DFT preserves energy* $\Rightarrow \|x\|^2 = x^H x = X^H X = \|X\|^2$

Proof.

- Use iDFT to write $x = FX$ and exploit fact that $F$ is unitary

$$\|x\|^2 = x^H x = (FX)^H FX = X^H F^H FX = X^H X = \|X\|^2$$

- This theorem would also be true for any transform pair

$$X = T^H x \iff \tilde{x} = TX$$

- As long as the transform matrix T is unitary $\Rightarrow T^H T = I$

- Are there other useful transforms defined by unitary matrices T?

  $\Rightarrow$ Many. One we have already found is the DCT

- Define the inverse DCT matrix C to write the iDCT as $\tilde{x} = CX$

$$
C = \frac{1}{\sqrt{N}} \begin{bmatrix} 1 & 1 & \cdots & 1 \\ 1 & \sqrt{2}\cos\left[\frac{2\pi(1)((1)+1/2)}{N}\right] & \cdots & \sqrt{2}\cos\left[\frac{2\pi(N-1)((1)+1/2)}{N}\right] \\ \vdots & \vdots & \ddots & \vdots \\ 1 & \sqrt{2}\cos\left[\frac{2\pi(1)((N-1)+1/2)}{N}\right] & \cdots & \sqrt{2}\cos\left[\frac{2\pi(N-1)((N-1)+1/2)}{N}\right] \end{bmatrix}
$$

- It is ready to verify that C is unitary (the cosines are orthonormal)

- From where the inverse and energy conservation theorems follow

  $\Rightarrow$ Proofs hold for all unitary matrices, C in particular

- A basic information processing theory can be built for any T

- Then, why do we specifically choose the DFT? Or the DCT?

  $\Rightarrow$ Oscillations represent different rates of change

  $\Rightarrow$ Different rates of change represent different aspects of a signal

- Not a panacea, though. E.g., $F^H$ is independent of the signal

- If we know something about signal, we should use it to build better T

- A way of "knowing something" is a stochastic model of the signal

- PCA: Principal component analysis

  $\Rightarrow$ Use the eigenvectors of the covariance matrix to build T

The Discrete Fourier Transform with Unitary Matrices

Stochastic signals

Principal Component Analysis (PCA) transform

Dimensionality reduction

Principal Components

Face recognition

- A random variable $X$ models a random phenomena

  $\Rightarrow$ One in which many different outcomes are possible

  $\Rightarrow$ And one in which some outcomes may be more likely than others

- Thus, a random variable represents two things

  $\Rightarrow$ All possible outcomes and their respective likelihoods



$p_X(x), \; p_Y(y), \; p_Z(y)$

$-\sigma_x \quad \sigma_x \quad \mu_Y - \sigma_Y \quad \mu_Y \quad \mu_Y + \sigma_Y \quad \mu_Z - \sigma_Z \quad \mu_Z \quad \mu_Z + \sigma_Z \qquad x, y, z$

- Random variable $X$ takes values around 0 and $Y$ values around $\mu_Y$

- $Z$ takes values around $\mu_Z$ and the values are more concentrated

▶ Probabilities measure the likelihood of observing different outcomes

$\Rightarrow$ Larger probability means an outcome that is more likely

$\Rightarrow$ Or, observed more often when seeing many realizations

▶ Random variables represented by uppercase $\Rightarrow$ E.g., $X$

▶ Values that it can take represented by lowercase $\Rightarrow$ E.g., $x$

▶ The probability that $X$ takes values between $x$ and $x'$ is written as

$$\mathrm{P}\left(x < X \leq x'\right)$$

▶ Here, we describe probabilities with density functions (pdf) $\Rightarrow$ $p_X(x)$

$$\mathrm{P}\left(x < X \leq x'\right) = \int_x^{x'} p_X(u)\, du$$

▶ $p_X(x) \approx$ How likely random variable $X$ is to take a value around $x$

- A random variable $X$ is Gaussian (or Normal) if its pdf is of the form

$$p_X(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-(x-\mu)^2/\sigma^2}$$

- The mean $\mu$ determines center. The variance $\sigma^2$ determines width



- Means satisfy $0 = \mu_X < \mu_Y < \mu_Z$. Variances are $\sigma_X^2 = \sigma_Y^2 > \sigma_Z^2$

▶ Expectation of random variable is an average weighted by likelihoods

$$\mathbb{E}\left[X\right] = \int_{-\infty}^{\infty} x p_X(x)\, dx$$

▶ Regular average $\Rightarrow$ Sum all values and divide by number of values

▶ Expectation $\Rightarrow$ Weight values $x$ by their relative likelihoods $p_X(x)$

▶ For a Gaussian random variable $X$ the expectation is the mean $\mu$

$$\mathbb{E}\left[X\right] = \int_{-\infty}^{\infty} x \frac{1}{\sqrt{2\pi}\sigma} e^{-(x-\mu)^2/\sigma^2}\, dx = \mu$$

▶ Not difficult to evaluate integral, but besides the point to do so here

▶ Measure of variability around the mean weighted by likelihoods

$$\text{var}\,[X] = \mathbb{E}\left[\left(X - \mathbb{E}\,[X]\right)^2\right] = \int_{-\infty}^{\infty}\left(x - \mathbb{E}\,[X]\right)^2 p_X(x)\,dx$$

▶ Large variance ≡ likely values are spread out around the mean

▶ Small variance ≡ likely values are concentrated around the mean

▶ For a Gaussian random variable $X$ the variance is the variance $\sigma^2$

$$\text{var}\,[X] = \int_{-\infty}^{\infty}\left(x - \mathbb{E}\,[X]\right)^2 \frac{1}{\sqrt{2\pi}\sigma} e^{-(x-\mu)^2/\sigma^2}\,dx = \sigma^2$$

▶ Not difficult to evaluate either. But also besides the point here

▶ A random signal X is a collection of random variables (length $N$)

$$\mathsf{X} = [X(0),\ X(1),\ \ldots,\ X(N-1)]^T$$

▶ Each of the random variables has its own pdf $\Rightarrow p_{X(n)}(x)$

▶ This pdf describes the likelihood of $X(n)$ taking a value around $x$

▶ This is not a sufficient description. Joint outcomes also important

▶ Joint pdf $p_{\mathsf{X}}(\mathsf{x})$ says how likely signal X is to be found around x

$$\mathsf{P}(\mathsf{x} \in \mathcal{X}) = \iint_{\mathcal{X}} p_{\mathsf{X}}(\mathsf{x})\, d\mathsf{x}$$

▶ The individual pdfs $p_{X(n)}(x)$ are said to be marginal pdfs

▶ Random signal X ⇒ All possible images of human faces
▶ More manageable ⇒ X is a collection of 400 face images
   ⇒ The random variable represents all the images
   ⇒ The likelihood of each of them being chosen. E.g., 1/400 each



▶ Random variable specified by all outcomes and respective probabilities

▶ Do observe that the dataset consists of images ≡ matrices

▶ Each image is stored in a matrix of size $112 \times 92$

$$M_i = \begin{bmatrix} m_{1,1} & m_{1,2} & \cdots & m_{1,92} \\ m_{2,1} & m_{2,2} & \cdots & m_{2,92} \\ \vdots & \vdots & \ddots & \vdots \\ m_{112,1} & m_{112,2} & \cdots & m_{112,92} \end{bmatrix}$$

▶ Stack columns of image $M_i$ into the vector $x_i$ with length $10,304$

$$x_i = \begin{bmatrix} m_{1,1}, & m_{21}, & \ldots, & m_{112,1}, & m_{1,2}, & m_{2,2}, & \ldots, & m_{112,2}, & \vdots, & m_{1,92}, & m_{2,92}, & \ldots, & m_{112,92} \end{bmatrix}^T$$

▶ Images are matrices $M_i \in \mathbb{R}^{112 \times 92}$. Signals are vectors $x_i \in \mathbb{R}^{10,304}$

▶ Realization x is an individual face pulled from set of possible outcomes

▶ Three possible realizations shown



▶ Realizations are just regular signals. Nothing random about them

▶ Signal's expectation is the concatenation of individual expectations

$$\mathbb{E}\left[\mathsf{X}\right] = \left[\mathbb{E}\left[X(0)\right], \ \mathbb{E}\left[X(1)\right], \ \ldots \ \mathbb{E}\left[X(N-1)\right]\right]^{T} = \iint \mathsf{x} p_{\mathsf{X}}(\mathsf{x}) \, d\mathsf{x}$$

▶ Variance of $n$th element $\Rightarrow \Sigma_{nn} = \mathrm{var}\left[X(n)\right] = \mathbb{E}\left[\left(X(n) - \mathbb{E}\left[X(n)\right]\right)^2\right]$

▶ Measures variability of $n$th component

▶ Covariance between the signal components $X(n)$ and $X(m)$

$$\Sigma_{nm} = \mathbb{E}\left[\left(X(n) - \mathbb{E}\left[X(n)\right]\right)\left(X(m) - \mathbb{E}\left[X(m)\right]\right)\right] = \Sigma_{mn}$$

▶ Measures how much $X(n)$ predicts $X(m)$. Love, hate, and indifference
   $\Rightarrow \Sigma_{nm} = 0$, components are unrelated. They are orthogonal
   $\Rightarrow \Sigma_{nm} > 0$ ($\Sigma_{nm} < 0$), move in same (opposite) direction

▶ Assume that $\mathbb{E}[X] = 0$ so that covariances are $\Sigma_{nm} = \mathbb{E}[X(n)X(m)]$

▶ Consider the expectation $\mathbb{E}[xx^T]$ of the (outer) product $xx^T$

▶ We can write the outer product $xx^T$ as

$$xx^T = \begin{bmatrix} x(0)x(0) & \cdots & x(0)x(n) & \cdots & x(0)x(N-1) \\ \vdots & \ddots & \vdots & & \vdots \\ x(n)x(0) & \cdots & x(n)x(n) & \cdots & x(n)x(N-1) \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ x(N-1)x(0) & \cdots & x(N-1)x(n) & \cdots & x(N-1)x(N-1) \end{bmatrix}$$

▶

- Assume that $\mathbb{E}[X] = 0$ so that covariances are $\Sigma_{nm} = \mathbb{E}[X(n)X(m)]$

- Consider the expectation $\mathbb{E}\left[xx^T\right]$ of the (outer) product $xx^T$

- Expectation $\mathbb{E}\left[xx^T\right]$ implies expectation of each individual element

$$
\mathbb{E}\left[xx^T\right] = \begin{bmatrix}
\mathbb{E}[x(0)x(0)] & \cdots & \mathbb{E}[x(0)x(n)] & \cdots & \mathbb{E}[x(0)x(N-1)] \\
\vdots & \ddots & \vdots & \ddots & \vdots \\
\mathbb{E}[x(n)x(0)] & \cdots & \mathbb{E}[x(n)x(n)] & \cdots & \mathbb{E}[x(n)x(N-1)] \\
\vdots & \ddots & \vdots & \ddots & \vdots \\
\mathbb{E}[x(N-1)x(0)] & \cdots & \mathbb{E}[x(N-1)x(n)] & \cdots & \mathbb{E}[x(N-1)x(N-1)]
\end{bmatrix}
$$

-

▶ Assume that $\mathbb{E}[X] = 0$ so that covariances are $\Sigma_{nm} = \mathbb{E}[X(n)X(m)]$

▶ Consider the expectation $\mathbb{E}\left[xx^T\right]$ of the (outer) product $xx^T$

▶ The $(n, m)$ element of the matrix $\mathbb{E}\left[xx^T\right]$ is the covariance $\Sigma_{nm}$

$$\mathbb{E}\left[xx^T\right] = \begin{bmatrix} \Sigma_{00} & \cdots & \Sigma_{0n} & \cdots & \Sigma_{0(N-1)} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ \Sigma_{n0} & \cdots & \Sigma_{nn} & \cdots & \Sigma_{n(N-1)} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ \Sigma_{(N-1)0} & \cdots & \Sigma_{(N-1)n} & \cdots & \Sigma_{(N-1)(N-1)} \end{bmatrix}$$

▶ Define the covariance matrix of random signal X as $\Sigma := \mathbb{E}\left[xx^T\right]$

▶ When the mean is not null define the covariance matrix of X as

$$\Sigma := \mathbb{E}\left[\left(\mathsf{x} - \mathbb{E}\left[\mathsf{x}\right]\right)\left(\mathsf{x} - \mathbb{E}\left[\mathsf{x}\right]\right)^T\right]$$

▶ As before, the $(n, m)$ element of $\Sigma$ is the covariance $\Sigma_{nm}$

$$((\Sigma))_{nm} = \mathbb{E}\left[\left(X(n) - \mathbb{E}\left[X(n)\right]\right)\left(X(m) - \mathbb{E}\left[X(m)\right]\right)\right] = \Sigma_{nm}$$

▶ The covariance matrix $\Sigma$ is an arrangement of the covariances $\Sigma_{nm}$

▶ The diagonal of $\Sigma$ contains the (auto)variances $\Sigma_{nn} = \mathrm{var}\left[X(n)\right]$

▶ Covariance matrix is symmetric $\Rightarrow ((\Sigma))_{nm} = \Sigma_{nm} = \Sigma_{mn} = ((\Sigma))_{mn}$

► All images are equally likely $\Rightarrow$ probability $1/400$ for each image

► The mean face is the regular average $\rightarrow \mathbb{E}[x] = \dfrac{1}{400} \sum^{400} x_i$



► Average image looks something, sort of, like an average face

▶ Covariance matrix $\Rightarrow \Sigma = \dfrac{1}{400} \displaystyle\sum_{i=1}^{400} \left( x_i - \mathbb{E}\left[ x \right] \right) \left( x_i - \mathbb{E}\left[ x \right] \right)^T$



▶ Heat map of covariance matrix $\Sigma$ shown on left

▶ Large correlation values around diagonal

▶ Large correlation values every 112 elements (jump a row on matrix)

The Discrete Fourier Transform with Unitary Matrices

Stochastic signals

Principal Component Analysis (PCA) transform

Dimensionality reduction

Principal Components

Face recognition

▶ Consider a vector with $N$ elements $\Rightarrow v = [v(0), v(1), \ldots, v(N-1)]$

▶ We say that v is an eigenvector of $\Sigma$ if for some scalar $\lambda \in \mathbb{R}$

$$\Sigma v = \lambda v$$

▶ We say that $\lambda$ is the eigenvalue associated to v



$\Sigma w$     $\Sigma v_1 = \lambda_1 v_1$     $\Sigma v_2 = \lambda_2 v_2$

w     $v_1$     $v_2$

▶ In general, non-eigenvectors w and $\Sigma w$ point in different directions

▶ But for eigenvectors v, the product vector $\Sigma v$ is collinear with v

▶ If v is an eigenvector, $\alpha$v is also an eigenvector for any scalar $\alpha \in \mathbb{R}$,

$$\Sigma(\alpha v) = \alpha(\Sigma v) = \alpha \lambda v = \lambda(\alpha v)$$

▶ Eigenvectors are defined up to a constant

▶ We use normalized eigenvectors with unit energy $\Rightarrow \|v\|^2 = 1$

▶ If we compute v with $\|v\|^2 \neq 1$ replace v with $v/\|v\|$

▶ There are $N$ eigenvalues and distinct associated eigenvectors

$\Rightarrow$ Some technical qualifications are needed in this statement

### Theorem

*The eigenvalues of $\Sigma$ are real and nonnegative* $\Rightarrow \lambda \in \mathbb{R}$ *and* $\lambda \geq 0$

### Proof.

- Begin by observing that we can write $\lambda = v^H \Sigma v / \|v\|^2$. Indeed

$$v^H \Sigma v = v^H (\Sigma v) = v^H (\lambda v) = \lambda v^H v = \lambda \|v\|^2$$

- Complete by showing that $v^T \Sigma v$ is nonnegative. Indeed (assume $\mathbb{E}[x] = 0$)

$$v^H \Sigma v = v^H \mathbb{E}\left[ x x^H \right] v = \mathbb{E}\left[ v^H x x^H v \right] = \mathbb{E}\left[ \left( v^H x \right) \left( x^H v \right) \right] = \mathbb{E}\left[ \left( v^H x \right)^2 \right] \geq 0$$

□

- Order eigenvalues from largest to smallest $\Rightarrow \lambda_0 \geq \lambda_1 \geq \ldots \geq \lambda_{N-1}$
- Eigenvectors inherit order $\Rightarrow v_0, v_1, \ldots, v_{N-1}$
- The $n$th eigenvector of $\Sigma$ is associated with its $n$th largest eigenvalue

**Theorem**
*Eigenvectors of $\Sigma$ associated with different eigenvalues are orthogonal*

Proof.

▶ Normalized eigenvectors v and u associated with eigenvalues $\lambda \neq \mu$

$$\Sigma v = \lambda v, \qquad \Sigma u = \mu u$$

▶ Since the matrix $\Sigma$ is symmetric we have $\Sigma^H = \Sigma$, and it follows

$$u^H \Sigma v = \left(u^H \Sigma v\right)^H = v^H \Sigma^H u = v^H \Sigma u$$

▶ Make $\Sigma v = \lambda v$ on the leftmost side and $\Sigma u = \mu u$ on the rightmost

$$u^H \lambda v = \lambda u^H v = \mu v^H u = v^H \mu u$$

▶ Eigenvalues are different $\Rightarrow$ Relationship can only be true if $v^H u = 0$

□

▶ One dimensional representation of first four eigenvectors $v_0, v_1, v_2, v_3$

▶ Two dimensional representation of first four eigenvectors $v_0, v_1, v_2, v_3$

▶ Define the matrix $T$ whose $k$th column is the $k$th eigenvector of $\Sigma$

$$T = [v_0, v_1, \ldots, v_{N-1}]$$

▶ Since the eigenvectors $v_k$ are orthonormal, the product $T^H T$ is

$$T^H T = \begin{bmatrix} v_0^H \\ \vdots \\ v_k^H \\ \vdots \\ v_{N-1}^H \end{bmatrix} \begin{bmatrix} v_0 & \cdots & v_k & \cdots & v_{N-1} \end{bmatrix}$$

$$= \begin{bmatrix} v_0^H v_0 & \cdots & v_1^H v_k & \cdots & v_0^H v_{N-1} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ v_k^H v_0 & \cdots & v_k^H v_k & \cdots & v_k^H v_{N-1} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ v_{N-1}^H v_{N-1} & \cdots & v_{N-1}^H v_k & \cdots & v_{N-1}^H v_{N-1} \end{bmatrix} = \begin{bmatrix} 1 & \cdots & 0 & \cdots & 0 \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ 0 & \cdots & 1 & \cdots & 0 \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ 0 & \cdots & 0 & \cdots & 1 \end{bmatrix}$$

▶ The eigenvector matrix T is unitary $\Rightarrow T^H T = I$

▶ Any unitary T can be used to define an info processing transform

▶ Define principal component analysis (PCA) transform $\Rightarrow y = T^H x$

▶ And the inverse (i)PCA transform $\Rightarrow \tilde{x} = Ty$

▶ Since T is unitary, iPCA is, indeed, the inverse of the PCA

$$\tilde{x} \;=\; Ty \;=\; T\left(T^H x\right) \;=\; TT^H x \;=\; Ix \;=\; x$$

▶ Thus y is an equivalent representation of x $\Rightarrow$ Back and forth

▶ And, also because T is unitary, Parseval's theorem holds

$$\|x\|^2 \;=\; x^H x \;=\; (Ty)^H \, Ty \;=\; y^H T^H Ty \;=\; y^H y \;=\; \|y\|^2$$

▶ Modifying elements $y_k$ means altering energy composition of signal

- The PCA transform is defined for any signal (vector) x

  $\Rightarrow$ But we expect to work well only when x is a realization of X

- Write the iPCA in expanded form and compare with the iDFT

$$x(n) = \sum_{k=0}^{N-1} y(k)v_k(n) \quad \Leftrightarrow \quad x(n) = \sum_{k=0}^{N-1} X(k)e_{kN}(n)$$

- The same except that they use different bases for the expansion

- Still, like developing a new sense.

- But not one that is generic. Rather, adapted to the random signal X

- PCA transform coefficients for given face image with 10,304 pixels

- Substantial energy in the first 15 PCA coefficients $y(k)$ with $k \leq 15$
- Almost all energy in the first 50 PCA coefficients $y(k)$ with $k \leq 50$
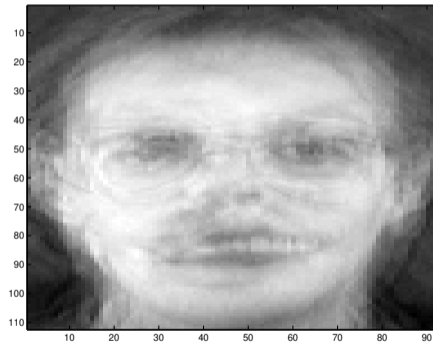  - $\Rightarrow$ This is a compression factor of more than 200



Coefficients for the first 50 eigenvectors

▶ Reconstructed image for increasing number of PCA coefficients
  ⇒ Increasing number of coefficients increases accuracy.
  ⇒ Using 50 coefficients suffices



Figure: image



Figure: No. P.C.s = 1

▶ Reconstructed image for increasing number of PCA coefficients
  ⇒ Increasing number of coefficients increases accuracy.
  ⇒ Using 50 coefficients suffices



Figure: image



Figure: No. P.C.s = 5

► Reconstructed image for increasing number of PCA coefficients
  ⇒ Increasing number of coefficients increases accuracy.
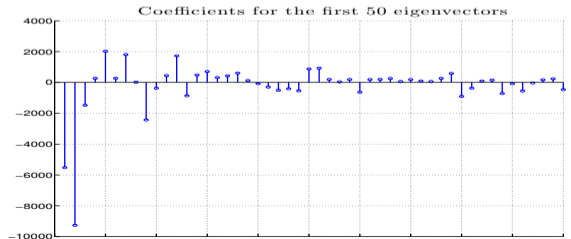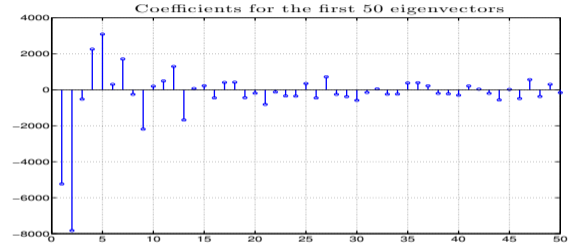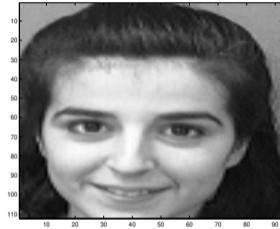  ⇒ Using 50 coefficients suffices



Figure: image



Figure: No. P.C.s = 10

► Reconstructed image for increasing number of PCA coefficients
  ⇒ Increasing number of coefficients increases accuracy.
  ⇒ Using 50 coefficients suffices



Figure: image



Figure: No. P.C.s = 20

- Reconstructed image for increasing number of PCA coefficients
    - ⇒ Increasing number of coefficients increases accuracy.
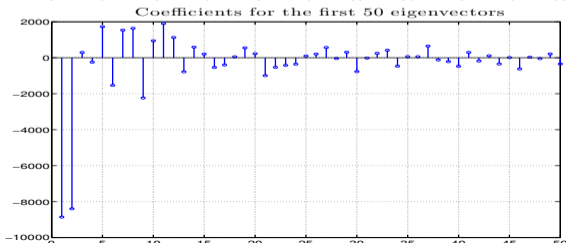    - ⇒ Using 50 coefficients suffices



Figure: image



Figure: No. P.C.s = 30

- Reconstructed image for increasing number of PCA coefficients
  - $\Rightarrow$ Increasing number of coefficients increases accuracy.
  - $\Rightarrow$ Using 50 coefficients suffices



Figure: image



Figure: No. P.C.s = 40

► Reconstructed image for increasing number of PCA coefficients
  ⇒ Increasing number of coefficients increases accuracy.
  ⇒ Using 50 coefficients suffices



Figure: image



Figure: No. P.C.s = 50

- PCA transform $y$ for two different pictures of the same person
- Coefficients are similar, even if pose and attitude are different
  - $\Rightarrow$ E.g., first two coefficients almost identical



Coefficients for the first 50 eigenvectors

Coefficients for the first 50 eigenvectors

- PCA transform $y$ for pictures of different persons
- Similar pose and attitude, but PCA coefficients are still different
  - $\Rightarrow$ Can be used to perform face recognition. More later



Coefficients for the first 50 eigenvectors

Coefficients for the first 50 eigenvectors

The Discrete Fourier Transform with Unitary Matrices

Stochastic signals

Principal Component Analysis (PCA) transform

Dimensionality reduction

Principal Components

Face recognition

▶ Transform signal x into frequency domain with DFT $X = F^H x$

▶ Recover x from X through iDFT matrix multiplication $x = FX$

▶ We compress by retaining $K < N$ DFT coefficients to write

$$\tilde{x}(n) = \sum_{k=0}^{K-1} X(k) e^{j2\pi kn/N}$$

▶ Equivalently, we define the compressed DFT as

$$\tilde{X}(k) = X(k) \quad \text{for} \quad k < K, \qquad \tilde{X}(k) = 0 \text{ otherwise}$$

▶ Reconstructed signal is obtained with iDFT $\Rightarrow \tilde{x} = F\tilde{X}$

▶ Transform signal x into eigenvector domain with PCA $y = T^H x$

▶ Recover x from y through iPCA matrix multiplication $x = Ty$

▶ We compress by retaining $K < N$ PCA coefficients to write

$$\tilde{x}(n) = \sum_{k=0}^{K-1} y(k) v_k(n)$$

▶ Equivalently, we define the compressed PCA as

$$\tilde{y}(k) = y(k) \quad \text{for} \quad k < K, \qquad \tilde{y}(k) = 0 \text{ otherwise}$$

▶ Reconstructed signal is obtained with iPCA $\Rightarrow \tilde{x} = T\tilde{y}$

▶ Why do we keep the first $K$ DFT coefficients?

⇒ Because faster oscillations tend to represent faster variation

⇒ Also, not always, sometimes we keep the largest coefficients

▶ Why do we keep the first $K$ PCA coefficients?

⇒ Eigenvectors with lower ordinality have larger eigenvalues

⇒ Larger eigenvalues entail more variability

⇒ And more variability signifies more dominant features

▶ Eigenvectors with large ordinality represent finer signal features

⇒ And can often be omitted

▶ PCA compression is (more accurately) called dimensionality reduction
  ⇒ Do not compress signal. Reduce number of dimensions

$$\Sigma = \begin{bmatrix} 3/2 & 1/2 \\ 1/2 & 3/2 \end{bmatrix}$$

▶ Covariance eigenvectors mix coordinates

$$v_0 = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad v_1 = \begin{bmatrix} 1 \\ -1 \end{bmatrix}$$

▶ Eigenvalues are $\lambda_0 = 2$ and $\lambda_1 = 1$



▶ Signal varies more in $v_0 = [1, 1]^T$ direction than in $v_1 = [1, -1]^T$
  ⇒ Study one dimensional signal $\tilde{x} = y(0)v_0$
  ⇒ instead of the original two dimensional signal x

► PCA dimensionality reduction minimizes the expected error energy

► To see that this is true, define the error signal as $\Rightarrow \mathrm{e} := \mathrm{x} - \tilde{\mathrm{x}}$

► The energy of the error signal is $\Rightarrow \|\mathrm{e}\|^2 = \|\mathrm{x} - \tilde{\mathrm{x}}\|^2$

► The expected value of the energy of the error signal is

$$\mathbb{E}\left[\|\mathrm{e}\|^2\right] = \mathbb{E}\left[\|\mathrm{x} - \tilde{\mathrm{x}}\|^2\right]$$

► Keeping the first $K$ PCA coefficients minimizes $\mathbb{E}\left[\|\mathrm{e}\|^2\right]$

$\Rightarrow$ Among all reconstructions that use, at most, $K$ coefficients

### Theorem

*The expectation of the reconstruction error is the sum of the eigenvalues corresponding to the eigenvectors of the coefficients that are discarded*

$$\mathbb{E}\left[\|\mathbf{e}\|^2\right] = \sum_{k=K}^{N-1} \lambda_k$$

▶ It follows that keeping the first $K$ PCA coefficients is optimal

$\Rightarrow$ In the sense that it minimizes the Expected error energy

▶ Good on average. Across realizations of the stochastic signal X

▶ Need not be good for given realization (but we expect it to be good)

Proof.

- Error signal signal is $e := x - \tilde{x}$. Define error PCA transform as $f = T^H x$
- Using Parseval's (energy conservation) we can write the energy of e as

$$\|e\|^2 = \|f\|^2 = \sum_{k=K}^{N-1} y^2(k)$$

- In the last equality we used that $f = y - \tilde{y} = [0, \ldots, 0, y(K), \ldots, y(N-1)]$
- Here, we are interested in the expected value of the error's energy

- Take expectation on both sides of equality $\Rightarrow \mathbb{E}\left[\|e\|^2\right] = \sum_{k=K}^{N-1} \mathbb{E}\left[y^2(k)\right]$

- Used the fact that expectations are linear operators

Proof.

▶ Compute expected value $\mathbb{E}\left[y^2(k)\right]$ of the squared PCA coefficient $y(k)$

▶ As per PCA transform definition $y(k) = v^H x$, which implies

$$\mathbb{E}\left[y^2(k)\right] \;=\; \mathbb{E}\left[(v_k^H x)^2\right] \;=\; \mathbb{E}\left[v_k^H x x^T v_k\right] \;=\; v_k^H \mathbb{E}\left[x x^T\right] v_k$$

▶ Covariance matrix: $\Sigma := \mathbb{E}\left[x x^T\right]$. Eigenvector definition $\Sigma v_k = \lambda_k$. Thus

$$\mathbb{E}\left[y^2(k)\right] \;=\; v_k^H \Sigma v_k \;=\; v_k^H \lambda_k v_k \;=\; \lambda_k \|v_k\|^2$$

▶ Substitute into expression for $\mathbb{E}\left[\|e\|^2\right]$ to write $\Rightarrow \mathbb{E}\left[\|e\|^2\right] = \displaystyle\sum_{k=K}^{N-1} \lambda_k$ $\qquad\square$

▶ Covariance matrix eigenvalues for faces dataset.
▶ Expected approximation error ⇒ Tail sum of eigenvalue distribution
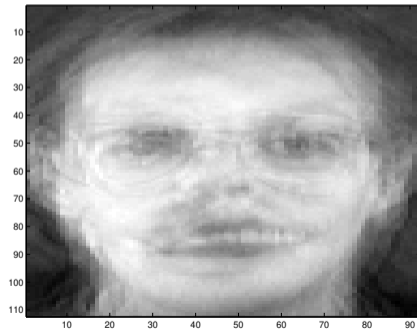  ⇒ Average across all realizations. Not the same as actual error



▶ First 10 coefficients have 98% of energy.
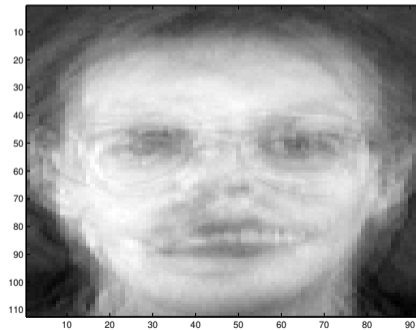▶ Eigenvectors with index $k > 50$ have $10^{-3}$% of energy on average

- Increasing number of coefficients reduces reconstruction error
- Average and actual reconstruction not the same (although "close")

- Keep 1 coefficient $\Rightarrow$ Reconstruction error $\Rightarrow$ 0.06
  $\Rightarrow$ Sum of removed eigenvalues $\Rightarrow$ 0.52

- Increasing number of coefficients reduces reconstruction error
- <span style="color:red">Average and actual reconstruction not the same</span> (although "close")

- Keep 5 coefficients $\Rightarrow$ Reconstruction error $\Rightarrow$ 0.03
  $\Rightarrow$ Sum of removed eigenvalues $\Rightarrow$ 0.11

▶ Increasing number of coefficients reduces reconstruction error

▶ Average and actual reconstruction not the same (although "close")

▶ Keep 10 coefficients $\Rightarrow$ Reconstruction error $\Rightarrow$ 0.02
$\qquad\qquad\qquad \Rightarrow$ Sum of removed eigenvalues $\Rightarrow$ 0.04

▶ Increasing number of coefficients reduces reconstruction error

▶ Average and actual reconstruction not the same (although "close")

▶ Keep 20 coefficients $\Rightarrow$ Reconstruction error $\Rightarrow$ 0.01
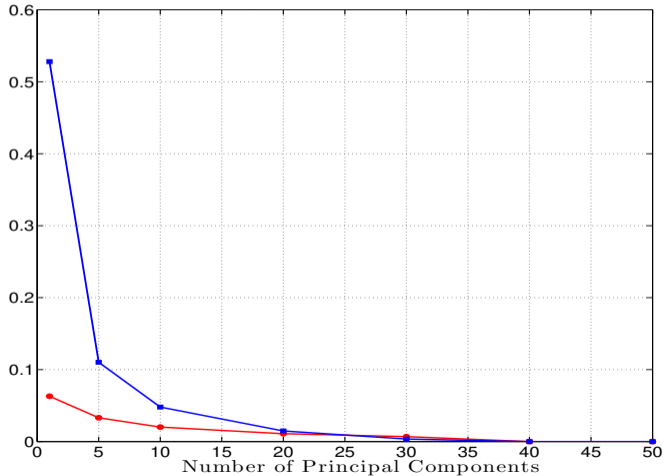
$\Rightarrow$ Sum of removed eigenvalues $\Rightarrow$ 0.01

▶ Increasing number of coefficients reduces reconstruction error

▶ Average and actual reconstruction not the same (although "close")

▶ Keep 30 coefficients ⇒ Reconstruction error ⇒ 0.006

⇒ Sum of removed eigenvalues ⇒ 0.003

- Increasing number of coefficients reduces reconstruction error
- Average and actual reconstruction not the same (although "close")

- Keep 40 coefficients $\Rightarrow$ Reconstruction error $\Rightarrow 0$
  $\Rightarrow$ Sum of removed eigenvalues $\Rightarrow 0$

▶ Increasing number of coefficients reduces reconstruction error

▶ Average and actual reconstruction not the same (although "close")

▶ Keep 50 coefficients $\Rightarrow$ Reconstruction error $\Rightarrow$ 0

$\Rightarrow$ Sum of removed eigenvalues $\Rightarrow$ 0

- ▶ Error for reconstruction process
- ▶ one realization (red), energy of removed eigenvalues (blue)

The Discrete Fourier Transform with Unitary Matrices

Stochastic signals

Principal Component Analysis (PCA) transform

Dimensionality reduction

Principal Components

Face recognition

▶ A random signal $X$ with uncorrelated components is one with

$$\Sigma_{nm} = \mathbb{E}\left[\left(X(n) - \mathbb{E}\left[X(n)\right]\right)\left(X(m) - \mathbb{E}\left[X(m)\right]\right)\right] = 0$$

▶ Different components are unrelated to each other.
▶ They represent different (orthogonal) aspects of signal

▶ Components uncorrelated $\Rightarrow$ The covariance matrix is diagonal

$$\Sigma = \mathbb{E}\left[\left(x - \mathbb{E}\left[x\right]\right)\left(x - \mathbb{E}\left[x\right]\right)^{T}\right] = \begin{bmatrix} \Sigma_{00} & \cdots & \Sigma_{0n} & \cdots & \Sigma_{0(N-1)} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ \Sigma_{n0} & \cdots & \Sigma_{nn} & \cdots & \Sigma_{n(N-1)} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ \Sigma_{(N-1)0} & \cdots & \Sigma_{(N-1)n} & \cdots & \Sigma_{(N-1)(N-1)} \end{bmatrix}$$
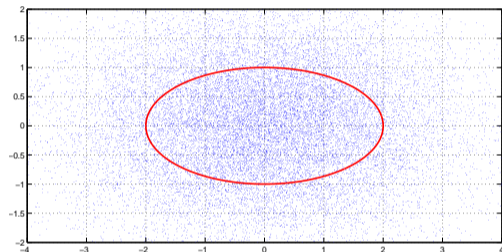
▶ How do eigenvectors (principal components) of uncorrelated signals look?

▶ Signal $X = [X(0), X(1)]^T$ with 2 components and diagonal covariance

$$\Sigma = \left[ \begin{array}{cc} 2 & 0 \\ 0 & 1 \end{array} \right]$$

▶ Covariance eigenvectors are

$$v_0 = \left[ \begin{array}{c} 1 \\ 0 \end{array} \right] \quad v_1 = \left[ \begin{array}{c} 0 \\ 1 \end{array} \right]$$



▶ The respective associated eigenvalues are $\lambda_0 = 2$ and $\lambda_1 = 1$

▶ Eigenvectors are orthogonal, as they should.
  $\Rightarrow$ Represent directions of separate signal variability
  $\Rightarrow$ Rate of variability given by associated eigenvalue

▶ Signal $X = [X(0), X(1)]^T$ with 2 components and diagonal covariance

$$\Sigma = \left[ \begin{array}{cc} 1 & 0 \\ 0 & 2 \end{array} \right]$$

▶ Covariance eigenvectors reverse order

$$v_0 = \left[ \begin{array}{c} 0 \\ 1 \end{array} \right] \quad v_1 = \left[ \begin{array}{c} 1 \\ 0 \end{array} \right]$$

▶ Associated eigenvalues are $\lambda_0 = 2$ and $\lambda_1 = 1$

▶ Eigenvectors still orthogonal, as they should.
⇒ Directions of separate signal variability
⇒ Rate given by associated eigenvalue

▶ Signal $X = [X(0), X(1)]^T$ with 2 components and diagonal covariance

$$\Sigma = \begin{bmatrix} 3/2 & 1/2 \\ 1/2 & 3/2 \end{bmatrix}$$

▶ Covariance eigenvectors mix coordinates

$$v_0 = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad v_1 = \begin{bmatrix} 1 \\ -1 \end{bmatrix}$$

▶ Eigenvalues are $\lambda_0 = 2$ and $\lambda_1 = 1$



▶ The eigenvalues are orthogonal. This is true for any covariance matrix
  ⇒ Mix coordinates but still represent directions of separate variability
  ⇒ Rate of change also given by associated eigenvalue

- Uncorrelated components means diagonal covariance matrix

$$\Sigma = \begin{bmatrix} \mathbf{\Sigma_{00}} & \cdots & \Sigma_{0n} & \cdots & \Sigma_{0(N-1)} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ \Sigma_{n0} & \cdots & \mathbf{\Sigma_{nn}} & \cdots & \Sigma_{n(N-1)} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ \Sigma_{(N-1)0} & \cdots & \Sigma_{(N-1)n} & \cdots & \mathbf{\Sigma_{(N-1)(N-1)}} \end{bmatrix}$$

- If variances are ordered, $k$th eigenvector is $k$-shifted delta $\delta(n - k)$

- The corresponding variance $\Sigma_{kk}$ is the associated eigenvalue

- Eigenvectors represent directions of orthogonal variability

- Rate of variability given by associated eigenvalue

- Correlated components means a full covariance matrix

$$\Sigma = \begin{bmatrix} \Sigma_{00} & \cdots & \Sigma_{0n} & \cdots & \Sigma_{0(N-1)} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ \Sigma_{n0} & \cdots & \Sigma_{nn} & \cdots & \Sigma_{n(N-1)} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ \Sigma_{(N-1)0} & \cdots & \Sigma_{(N-1)n} & \cdots & \Sigma_{(N-1)(N-1)} \end{bmatrix}$$

- The eigenvectors $v_k$ now mix different components

    $\Rightarrow$ But they still represent directions of orthogonal variability

    $\Rightarrow$ With the rate of variability given by associated eigenvalue

- PCA transform represents a signal as a sum of orthonormal vectors

    $\Rightarrow$ Each of which represents independent variability

- Principal components (eigenvectors) with larger eigenvalues represent directions in which the signal has more variability

The Discrete Fourier Transform with Unitary Matrices

Stochastic signals

Principal Component Analysis (PCA) transform

Dimensionality reduction

Principal Components

Face recognition

- ▶ Observe faces of known people ⇒ Use them to train classifier

- ▶ Observe a face of unknown character ⇒ Compare and classify

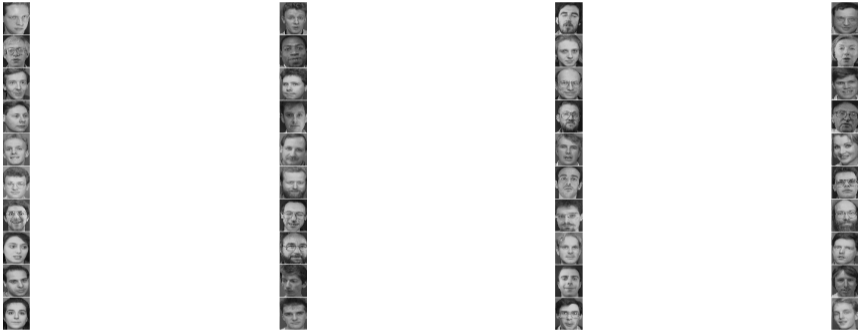- ▶ The dataset we've used contains 10 different images of 40 people

▶ Separate the first 9 of each person to construct training set



▶ Interpret these images as know, and use them to train classifier

▶ Utilize the last image of each person to construct a test set



▶ Interpret these images as unknown, and use them to test classifier

▶ Training set contains (signal, label) pairs $\Rightarrow \mathcal{T} = \{(x_i, z_i)\}_{i=1}^{N}$

▶ Signal $x$ is the face image. Label $z$ is the person's "name"

▶ Given (unknown) signals $x$, we want to assign a label

▶ Nearest neighbor classification rule

$\Rightarrow$ Find nearest neighbor signal in the training set

$$x_{NN} := \underset{x_i \in \mathcal{T}}{\operatorname{argmin}} \|x_i - x\|^2$$

$\Rightarrow$ Assign the label associated with the nearest neighbor

$$x_{NN} \quad \Rightarrow \quad (x_i, z_i) \quad \Rightarrow \quad z = z_i$$

▶ Reasonable enough. It should work. But it doesn't

- Image has a part that is inherent to the person ⇒ The actual signal
- But it also contains variability ⇒ Which we model as noise

$$x_i = \tilde{x}_i + w$$

- Problem is, there is more variability (noise) than signal
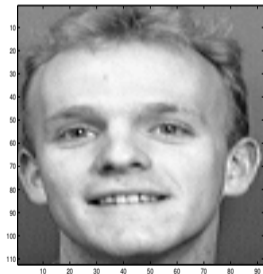


Figure: Test image



Figure: Nearest neighbor

▶ Compute PCA for all elements of training set $\Rightarrow$ $y_i = T^H x_i$

▶ Redefine training set as one with PCA transforms $\Rightarrow$ $\mathcal{T} = \{(y_i, z_i)\}_{i=1}^N$

▶ Compute PCA transform of (unknown) signal $x$ $\Rightarrow$ $y = T^H x$

▶ PCA nearest neighbor classification rule

$\Rightarrow$ Find nearest neighbor signal in training set with PCA transforms

$$y_{NN} := \underset{y_i \in \mathcal{T}}{\operatorname{argmin}} \|y_i - y\|^2$$

$\Rightarrow$ Assign the label associated with the nearest neighbor

$$y_{NN} \quad \Rightarrow \quad (y_i, z_i) \quad \Rightarrow \quad z = z_i$$

▶ Reasonable enough. It should work. And it does

- Recall: image = a part that belongs to the person + noise

$$x_i = \tilde{x}_i + w$$

- PCA transformation $T = [v_0^T; \ldots; v_{N-1}^T]$ leads to

$$y_i = Tx_i = T\tilde{x}_i + Tw$$

- PCA concentrates energy of $\tilde{x}_i$ on a few components

- But it keeps the energy of the noise on all components

- Keeping principal components improves the accuracy of classification

$\Rightarrow$ Because it increases the signal to noise ratio

▶ The training set $D = \{x_1, \ldots, x_{360}\}$ where $x_i \in \mathbb{R}^{10304}$ is given

▶ Compute the mean vector and the covariance matrix as

$$\bar{x} = \frac{1}{n} \sum_{i=1}^{n} x_i \quad \text{and} \quad \Sigma := \frac{1}{n} \sum_{i=1}^{n} (x_i - \bar{x}_i)(x_i - \bar{x}_i)^T.$$
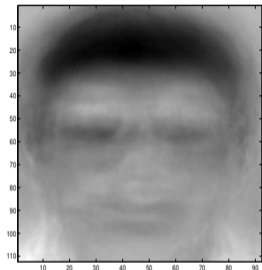
▶ Find the $k$ largest eigenvalues of $\Sigma$

▶ Store their corresponding eigenvalues $v_0, \ldots, v_{k-1} \in \mathbb{R}^{10304}$ as P.C.

$\Rightarrow$ The Principal Components $v_0, \ldots, v_{k-1}$ are called eigenfaces

▶ Create the PCA transform matrix as $T = [v_0^T; \ldots; v_{k-1}^T]$

▶ Project the training set into the space of P.C.s $y_i = Tx_i$

▶ $\Sigma$ depends training set, but is also a good description of the test set
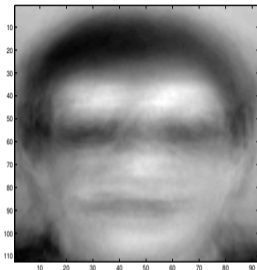
▶ The average face of the training set
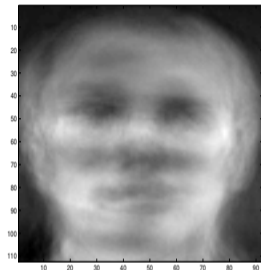
▶ The top 6 eigenfaces of the training set
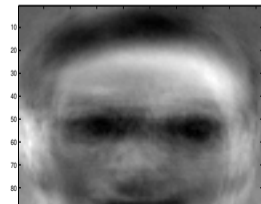


(1)



(2)



(3)

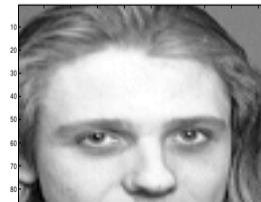Num. of P.C.          test point          N.N. in the training set



$k = 1$



$k = 5$

Classification method      test point      result of classification

Naive N.N.



PCA-ed($k = 5$) N.N.