# Lecture 10: { Karhunen-Loève Transform / Principal Component Analysis

Consider a stochastic process in $\mathbb{C}^n$ s.t.

**R.V.:**
**upper-**
**case**

$\mathbb{X} \in \mathbb{C}^n \sim f_{\mathbb{X}}(x_1, \cdots, x_n)$ pdf.
$\quad$ ↳ random vector, i.e., each coordinate is a r.v.
$\quad$ (or random signal)

Now consider its **covariance**:

$$\Gamma_{\mathbb{X}}[k, \ell] := E\left[(X_k - EX_k)\overline{(X_\ell - EX_\ell)}\right]$$

$\quad\quad\quad\quad\quad\quad$ ↳ expectation $\quad\quad\quad\quad k, \ell = 1:n$

i.e., $\quad \Gamma_{\mathbb{X}} := E(\mathbb{X} - E\mathbb{X})(\mathbb{X} - E\mathbb{X})^*$

$$= E\,\mathbb{X}\,\mathbb{X}^* - (E\mathbb{X})(E\mathbb{X})^* \in \mathbb{C}^{n \times n}$$

**Realizations**
**(obs.'s)**
**⇒ lower-**
**case**

Let $\{\mathbb{x}_1, \cdots, \mathbb{x}_N\}$ be $N$ realizations of $\mathbb{X}$. Then, the **sample estimate** of $\Gamma_{\mathbb{X}}$ is:

$$\hat{\Gamma}_{\mathbb{X}} := \frac{1}{N} \sum_{j=1}^{N} \mathbb{x}_j \, \mathbb{x}_j^* - \bar{\mathbb{x}}\,\bar{\mathbb{x}}^*$$

where $\quad \bar{\mathbb{x}} := \frac{1}{N} \sum_{j=1}^{N} \mathbb{x}_j$

<u>If</u> we define the **data matrix**

$$X := [\mathbb{x}_1 \mid \mathbb{x}_2 \mid \cdots \mid \mathbb{x}_N] \in \mathbb{C}^{n \times N},$$

<u>then</u> $\quad \hat{\Gamma}_{\mathbb{X}} = \frac{1}{N} X X^* - \bar{\mathbb{x}}\,\bar{\mathbb{x}}^*.$

Suppose we want a <u>**data-adaptive ONB**</u> of $\mathbb{C}^n$ (unlike fixed ONBs such as DFT, DCT, DST) s.t. the realizations of this stochastic process as a whole (i.e., <u>on average</u>) can be **best** approximated by $m$ <u>coordinates</u> with $\underline{m \ll n}$ in the mean-squared ($L^2$) error sense.

Let $W \in U(n) :=$ a set of all <span style="color:red">unitary</span> matrices in $\mathbb{C}^n$
and let $Y = W^* X$.

$\Rightarrow$ Viewing $X$ relative to $W$ (or the ONB consisting of col's of $W$)
$X$ is viewed relative to $I_n$.

$$\Rightarrow X = WY = [w_1 | \cdots | w_n] \begin{bmatrix} Y_1 \\ \vdots \\ Y_n \end{bmatrix}$$

$$= [e_1 | \cdots | e_n] \begin{bmatrix} X_1 \\ \vdots \\ X_n \end{bmatrix}$$

$$= Y_1 w_1 + \cdots + Y_n w_n.$$

Suppose we retain $Y_1 \sim Y_m$ and replace $Y_{m+1} \sim Y_n$ by $\alpha_{m+1} \sim \alpha_n$ (predetermined val's).
Then the approx. of $X$ is given by

$$X^{(m)} := \sum_{j=1}^{m} Y_j w_j + \sum_{j=m+1}^{n} \alpha_j w_j$$

The error is $\Delta X := X - X^{(m)} = \sum_{j=m+1}^{n} (Y_j - \alpha_j) w_j$,
and the mean-squared error is

$$\mathcal{E}(m) := E \| \Delta X \|^2 = E (\Delta X)^* (\Delta X)$$

$$= E \left[ \sum_{j=m+1}^{n} \sum_{k=m+1}^{n} \overline{(Y_j - \alpha_j)} (Y_k - \alpha_k) \underbrace{w_j^* w_k}_{= \delta_{jk}} \right]$$

$$= E \left[ \sum_{j=m+1}^{n} |Y_j - \alpha_j|^2 \right]$$

<u>Step 1</u>. $\dfrac{\partial \mathcal{E}^{(m)}}{\partial \alpha_j} = -2\,E\,(Y_j - \alpha_j) = 0,\quad j = m+1 : n.$

$\Rightarrow \quad \alpha_j = E\,Y_j = E\,[\,w_j^* X\,] = w_j^*\,E[X].$

Then, $\mathcal{E}(m) = E\left[\displaystyle\sum_{j=m+1}^{n} (Y_j - \alpha_j)\overline{(Y_j - \alpha_j)}\right]$

$\qquad\qquad = \displaystyle\sum_{j=m+1}^{n} E\,w_j^*(X - EX)\big(w_j^*(X - EX)\big)^*$

$\qquad\qquad = \displaystyle\sum_{j=m+1}^{n} E\,w_j^*(X - EX)(X - EX)^* w_j$

$\qquad\qquad = \displaystyle\sum_{j=m+1}^{n} w_j^*\,E\,(X - EX)(X - EX)^* w_j$

$\qquad\qquad = \displaystyle\sum_{j=m+1}^{n} w_j^*\,\Gamma_X\,w_j$

<u>Step 2</u>.  What kind of $\{w_j\}$ minimizes
the above quantity subject to $w_j^* w_j = 1$.

$\Rightarrow$ Use the <span style="color:red">**Lagrange multiplier**</span>:

$\tilde{\mathcal{E}}(m) := \mathcal{E}(m) - \displaystyle\sum_{j=m+1}^{n} \lambda_j\,(w_j^* w_j - 1)$

$\qquad\quad = \displaystyle\sum_{j=m+1}^{n} \left[\,w_j^*\,\Gamma_X\,w_j - \lambda_j\,(w_j^* w_j - 1)\,\right]$

$\dfrac{\partial \tilde{\mathcal{E}}(m)}{\partial w_j} = 2\,\Gamma_X\,w_j - 2\,\lambda_j\,w_j = 0$

$\qquad\qquad \Rightarrow \quad \Gamma_X\,w_j = \lambda_j\,w_j \qquad$ Eigenvalue problem!

$1 \le m \le n$ was arbitrary.  So, for any $1 \le m \le n$,
we need to solve $\Gamma_X\,w_j = \lambda_j\,w_j$, $j = 1 : m$.

$\Gamma_X$: hermitian $\to$ unitarily diagonalizable & $\lambda_j \in \mathbb{R}$.

## Remarks

(1) In practice, $\hat{\Gamma}_X$ is used for $\Gamma_X$.
The quality of $\hat{\Gamma}_X$ depends on $n > N$ or $n < N$.

$\begin{cases} \text{classical setting}: \ n \ll N \ (\text{e.g., census}) \Rightarrow \hat{\Gamma}_X : \text{good} \\ \textcolor{red}{\text{neo}}\text{classical setting}: \ n \gg N \ (\text{e.g., images}) \Rightarrow \hat{\Gamma}_X : \text{poor} \end{cases}$

(2) "Only" optimal in terms of the mean-squared (or the entropy) criterion.

Satosi Watanabe (1965) : The Entropy Minimization Criterion.

Define $\delta_X := \text{diag}(\Gamma_X) / \|\text{diag}(\Gamma_X)\|_1, \in \mathbb{R}^n$

Consider all possible $W \in U(n)$ and the coordinate transf's $W^* X$.

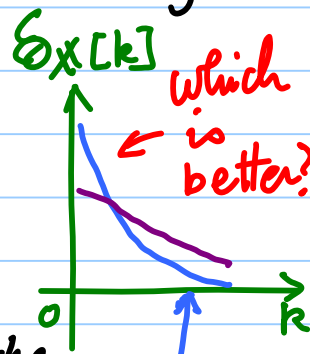Define $H(\mathbb{P}) := -\sum_{i=1}^{n} p_i \log p_i, \quad p_i \geq 0, \ \sum_{1}^{n} p_i = 1$.

<span style="color:red">the Shannon entropy</span>

$\delta_X$ above qualifies as $\mathbb{P}$ thanks to the normalization.
Then Watanabe showed that

<span style="color:green">$\delta_{X[k]}$</span> <span style="color:red">which ← is better?</span>

$$H(\delta_{W_{KL}^* X}) = \min_{W \in U(n)} H(\delta_{W^* X})$$

i.e., $W_{KL} = \arg\min_{W \in U(n)} H(\delta_{W^* X})$.



That is, the KLT (or PCA) provides the <span style="color:red">minimum entropy</span> coordinates (<span style="color:blue">sharper distribution !</span>)
$\Rightarrow$ Packing more energy (or variance) into the first few coordinates !

(3) $\mathbb{Y} = W_{KL}^* \mathbb{X}$, $W_{KL} = [w_1 | \cdots | w_n]$
eigenvectors of $\Gamma_{\mathbb{X}}$.

Then, what about $\Gamma_{\mathbb{Y}}$?

$$\Gamma_{\mathbb{Y}} = E(\mathbb{Y} - E\mathbb{Y})(\mathbb{Y} - E\mathbb{Y})^*$$
$$= E W_{KL}^* (\mathbb{X} - E\mathbb{X})(\mathbb{X} - E\mathbb{X})^* W_{KL}$$
$$= W_{KL}^* E(\mathbb{X} - E\mathbb{X})(\mathbb{X} - E\mathbb{X})^* W_{KL}$$
$$= W_{KL}^* \Gamma_{\mathbb{X}} W_{KL} = \text{diag}(\lambda_1, \cdots, \lambda_n)$$
$$\Gamma_{\mathbb{X}} w_j = \lambda_j w_j, \quad \lambda_j > 0.$$

$\Rightarrow$ The components of $\mathbb{Y}$ are <span style="color:red">decorrelated</span>!
$$E(Y_i - EY_i)(Y_j - EY_j) = \lambda_i \delta_{ij}$$

## ☆ Relationship between KLT/PCA and SVD

For simplicity, let's consider the centered data matrix $\tilde{X}$ of $X \in \mathbb{C}^{n \times N}$
$$\tilde{X} := X - \frac{1}{N} X \mathbb{1}\mathbb{1}^T = X \underbrace{(I - \frac{1}{N}\mathbb{1}\mathbb{1}^T)}_{\text{centering matrix}}$$
where $\mathbb{1} := (1, 1, \cdots, 1)^T \in \mathbb{R}^N$.
Note $\overline{X} = \frac{1}{N} X \mathbb{1}$, so multiplying the centering matrix from right subtracts $\overline{X}$ from each col. vector $X_j$ of $X$, $j = 1, \cdots, N$.

$\Rightarrow$ The mean of the col. vectors of $\tilde{X} = 0$.

Let the <span style="color:red">Singular Value Decomposition (SVD)</span> of $\tilde{X}$ be $\tilde{X} = U\Sigma V^*$ (full SVD)

where $U \in U(n) \subset \mathbb{C}^{n \times n}$
$V \in U(N) \subset \mathbb{C}^{N \times N}$
$\Sigma \in \mathbb{R}^{n \times N}$ : diagonal $\sigma_1, \cdots, \sigma_{min(n,N)}$

$$\Sigma = n\begin{bmatrix} \sigma_1 & & & & \\ & \ddots & & & O \\ & & \sigma_r & & \\ & & & 0 \cdots & \\ & & & & 0 \end{bmatrix} \quad \text{or} \quad N\left\{\begin{bmatrix} \sigma_1 & & & \\ & \ddots & & \\ & & \sigma_r & 0 \cdots \\ & & & 0 \\ \hline & & O & \end{bmatrix}\right.$$

$\underbrace{\quad}_{n} \quad \underbrace{\quad}_{N-n}$ $\quad$ $n-N\left\{\right.$ $\underbrace{\qquad}_{N}$

<span style="color:green">$n < N$</span>
<span style="color:green">classical</span>
$\qquad$ <span style="color:green">$n > N$</span>
<span style="color:green">neoclassical</span>

$r = rank(\tilde{X}) \leq min(n, N{\color{red}-1})$.

Note that if $\{x_1, \cdots, x_N\}$ are <u>linearly independent</u>, then $rank(X) \leq min(n, N)$. However, in the case of its centered version $\tilde{X} = [\tilde{x}_1 | \cdots | \tilde{x}_N]$, $\tilde{x}_j : x_j - \bar{X}$, $j = 1:N$, its column vectors are <span style="color:red">not</span> linearly indep. because $\tilde{x}_1 + \cdots + \tilde{x}_N = 0$.

$\Rightarrow \quad rank(\tilde{X}) \leq min(n, N-1)$.

(1) **If** $n < N$, **then**

$$\hat{\Gamma}_{\tilde{X}} = \frac{1}{N} \tilde{X} \tilde{X}^* = \frac{1}{N} U \Sigma V^* V \Sigma^* U^*$$

$$= \frac{1}{N} U \Sigma \Sigma^* U^*$$

$$= U \begin{bmatrix} \sigma_1^2/N & & 0 \\ & \ddots & \\ 0 & & \sigma_n^2/N \end{bmatrix} U^* \Rightarrow \begin{array}{l} \text{The eigenvalue} \\ \text{decomposition} \\ \text{of } \hat{\Gamma}_{\tilde{X}} \text{ !!} \end{array}$$

Hence, in this case

$$U = W_{KL} \text{ !}$$

(2) **If** $n > N$, **then** we should not compute all $n$ KLB vectors because $rank(\tilde{X}) \leq N-1$, so computing more than $N-1$ KLB vectors is useless. Moreover $n$ could be huge.

$\Rightarrow$ How to compute the top $N-1$ KLB vec's?

<span style="color:red">The first $N-1$ column vectors of $\tilde{X} V \in \mathbb{C}^{n \times N}$
= the top $N-1$ KLB vectors!</span>

<u>why?</u>  $\hat{\Gamma}_{\tilde{X}} = \frac{1}{N} \tilde{X} \tilde{X}^*$

So, $\hat{\Gamma}_{\tilde{X}} \tilde{X} V = \frac{1}{N} \tilde{X} \tilde{X}^* \tilde{X} V$

$$= \frac{1}{N} \tilde{X} (V \Sigma^* U^* U \Sigma V^*) V$$

$$= \frac{1}{N} \tilde{X} V \underbrace{\Sigma^* \Sigma}_{N \times N \text{ diagonal!}} = \tilde{X} V \begin{bmatrix} \sigma^2_{1/N} & & & O \\ & \ddots & & \\ & & \sigma^2_{N-1/N} & \\ O & & & 0 \end{bmatrix}$$

$\Rightarrow$ Each column vector of $\tilde{X} V$ is a linear combination of the column vectors of $\tilde{X}$, i.e., belongs to span$\{ \tilde{X}_1, \cdots, \tilde{X}_N \}$ !!
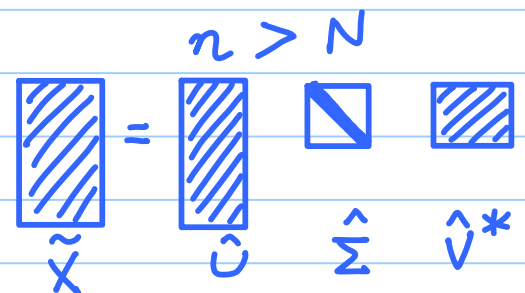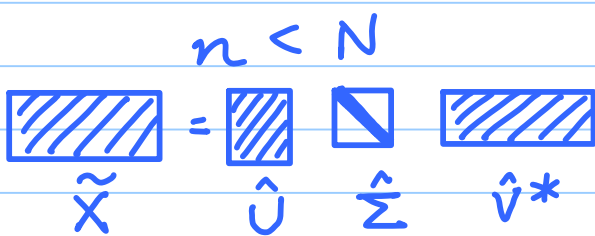
<u>Remark</u>: In either case ($n < N$ or $n > N$), you do not need to compute $\Gamma_{\tilde{X}}$. Moreover, you should use the <span style="color:red">reduced SVD</span> instead of the full SVD for the purpose of KLB/PCA computation.

The reduced SVD of $\tilde{X} \in \mathbb{C}^{n \times N}$

$$\tilde{X} = \hat{U} \hat{\Sigma} \hat{V}^*$$

Let $p := \min(n, N-1) \geq r = \text{rank}(\tilde{X})$.
Then $\hat{U} \in \mathbb{C}^{n \times p}, \quad \hat{\Sigma} \in \mathbb{R}^{p \times p}, \quad \hat{V} \in \mathbb{C}^{N \times p}$



In MATLAB, this is done by
$$\gg [\hat{U}, \hat{\Sigma}, \hat{V}] = \text{svd}(\tilde{X}, \text{`econ'});$$

Note also $\tilde{X}\hat{V} = \hat{U}\hat{\Sigma}\hat{V}^*\hat{V} = \hat{U}\hat{\Sigma}$

$$= [\sigma_1 u_1, \cdots \sigma_p u_p]$$

$$= [\tilde{X}v_1, \cdots, \tilde{X}v_p]$$

So, $u_j = \frac{1}{\sigma_j}\tilde{X}v_j$ , $j = 1, \cdots, p = \min(n, N-1)$.
$$= N-1 \text{ if } n \gg N.$$

In other words, each principal axis $u_j$ is just a linear combination of the (centered) input vectors $\tilde{X}_1, \cdots, \tilde{X}_N$ !

Example: The Rogues' Gallery Dataset
- Through the courtesy of Prof. Larry Sirovich
- A set of digitized photos of 143 faces each of which has $128 \times 128$ ixels i.e., $n = 128^2 = 16384$.
- These faces were of a specific group of people, i.e., Caucasian male students (and some faculty) at Brown Univ., without glasses, mustache, beard.
- Horizontal dilation has been applied s.t. the pupils are placed on two fixed common points.
- 143 faces were randomly divided into 72 and 71 faces for the training and test datasets, i.e., $N = 72$.
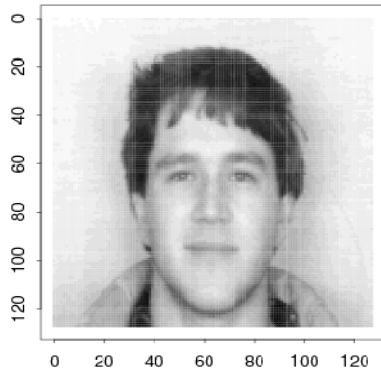- KLB/PCA were computed on the training dataset.
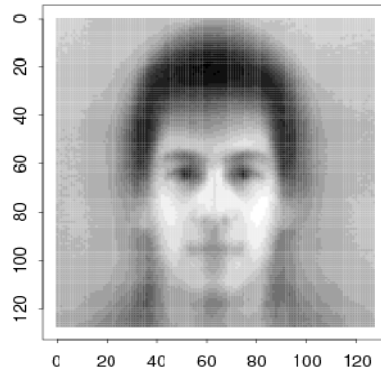
# 72 Faces !

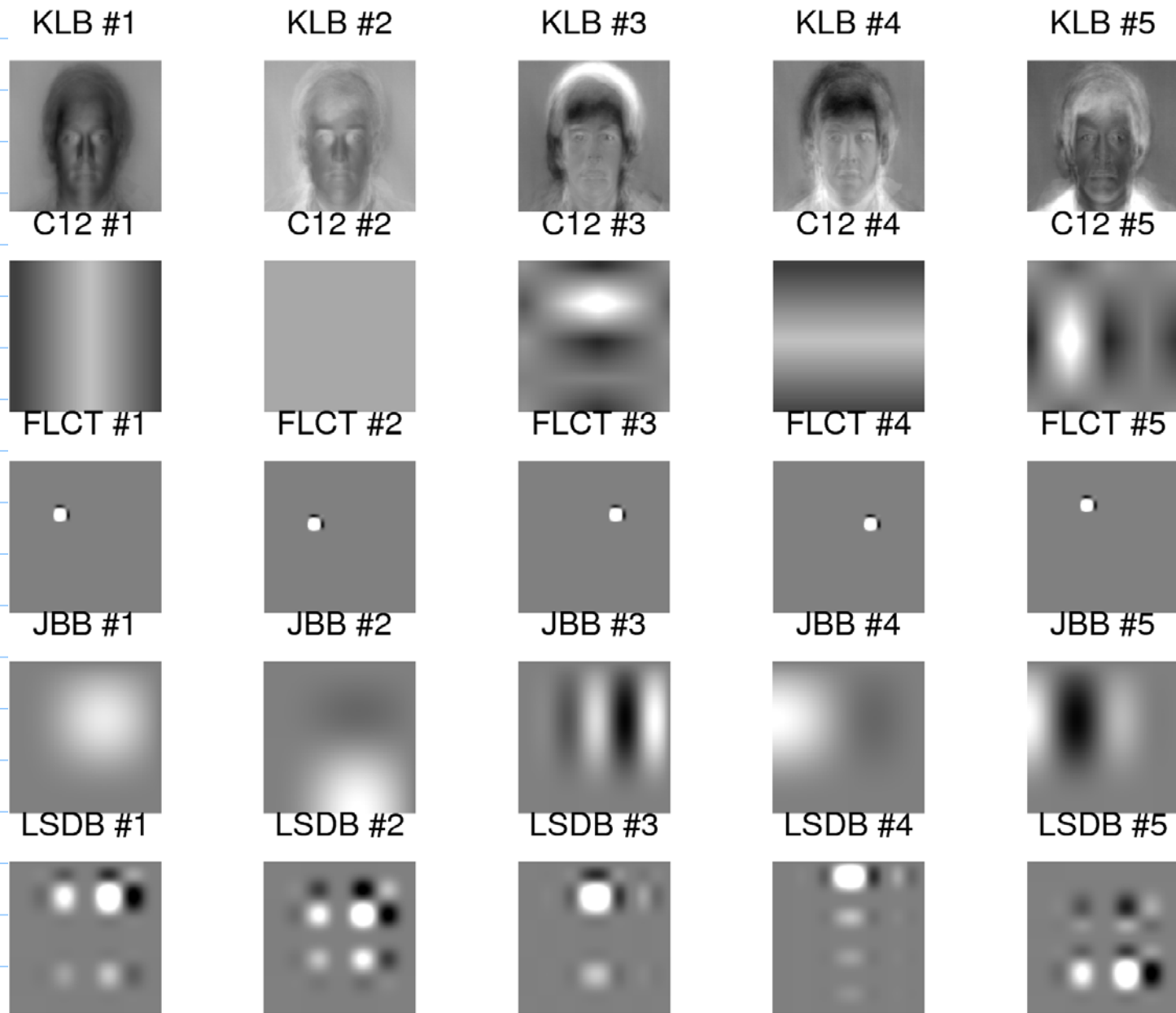$x_1$  $x_2$  $\cdots$



$x_{71}$  $x_{72}$

Original



one of the
test image

Average face



$\overline{X}$

# Comparison of Bases

| KLB #1 | KLB #2 | KLB #3 | KLB #4 | KLB #5 |
|--------|--------|--------|--------|--------|



| C12 #1 | C12 #2 | C12 #3 | C12 #4 | C12 #5 |
|--------|--------|--------|--------|--------|



| FLCT #1 | FLCT #2 | FLCT #3 | FLCT #4 | FLCT #5 |
|---------|---------|---------|---------|---------|



| JBB #1 | JBB #2 | JBB #3 | JBB #4 | JBB #5 |
|--------|--------|--------|--------|--------|



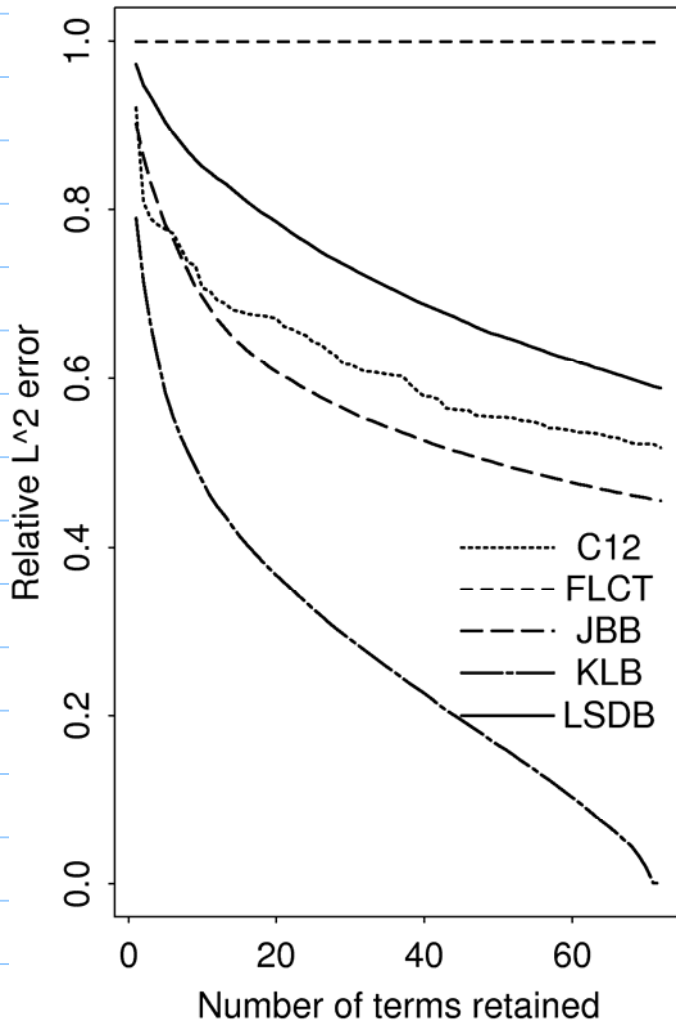| LSDB #1 | LSDB #2 | LSDB #3 | LSDB #4 | LSDB #5 |
|---------|---------|---------|---------|---------|



## Remarks:

(1) I am only displaying the top 5 KLB vectors. There are totally 71 KLB vectors plus the mean vector $\bar{X}$.
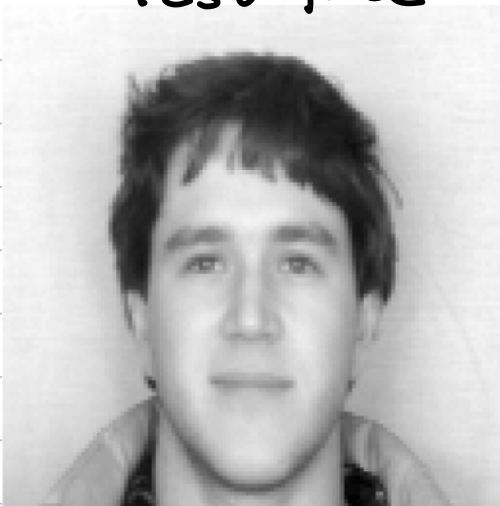
(2) As you can see these KLB vectors = lin. combi's of the 72 faces.
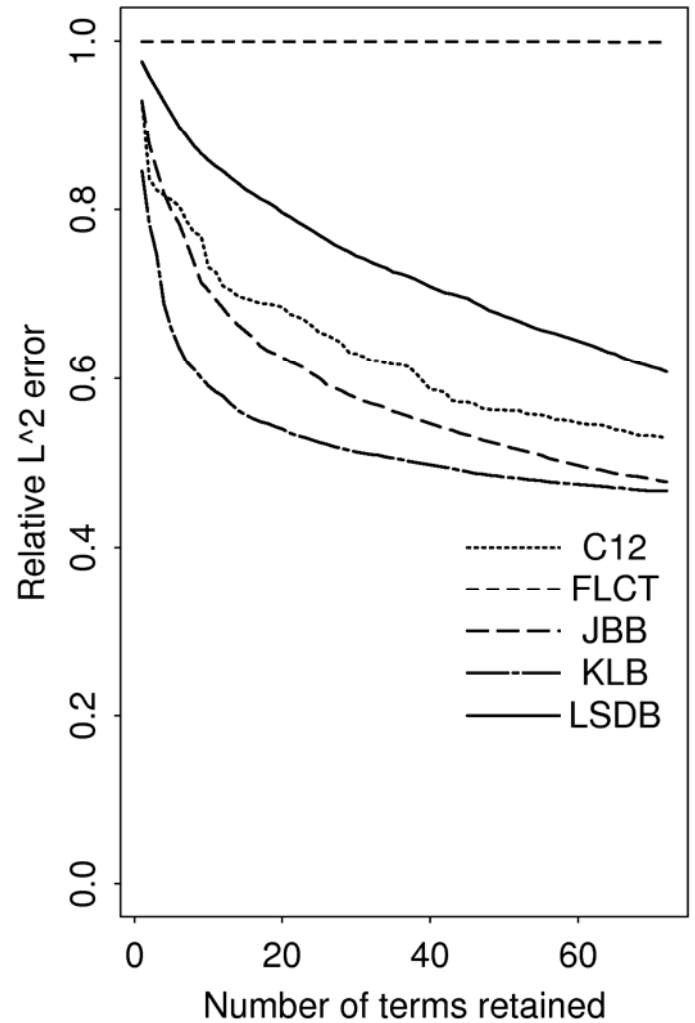
# Approximation Errors

## (a) Training Dataset Errors



## (b) Test Dataset Errors



test face



71 term approx. + ave.

## Remarks:

Note that for each training face, one can reconstruct it within machine precision using the full 71 KLB vectors and the average face $\bar{X}$.

However, for the test faces that were **not** used to generate the KLB, the reconstructions are not perfect.

Why? Basically, viewing these faces as the realizations of a stochastic process using only up to the 2nd order statistics (i.e., mean + covariance) and only using 72 realizations lead to serious limitation of the KLB applicability.

For simplicity, let's assume $E X_j = 0$, $j = 1:n$.
Suppose the underlying stochastic process
is stationary in wide sense, i.e.,
$$\Gamma_X[k, \ell] = \Gamma_X[k+m, \ell+m], \quad \forall m \in \mathbb{Z}$$
(appropriate mod N needs here)

Furthermore, assume
$$\Gamma_X[k, \ell] = \rho^{|k-\ell|} \qquad 0 < \rho < 1.$$
This is called the 1st order Markov model,
i.e., $X_{k+1}$ depends only on $X_k$ (plus
independent noise). This can be
explained as follows.

Suppose $\quad X_{k+1} = \rho X_k + Z_k$
where $\rho \in (0,1)$, $Z_k$ : iid. $E Z_k = 0$
$$\operatorname{var} Z_k = 1.$$
Then
$$E X_k X_{k+1}^* = \rho E |X_k|^2 + E X_k Z_k$$
$$\underset{= E X_k E Z_k}{\underbrace{\phantom{E X_k Z_k}}} = E X_k E Z_k$$
$$= 0 \qquad \overset{"}{0}$$
$$E X_k X_{k+2}^* = E X_k (\rho X_{k+1} + Z_{k+1})$$
$$= \rho E X_k X_{k+1} + E X_k Z_{k+1}$$
$$= \rho^2 E |X_k|^2$$
$$\vdots$$
$$E X_k X_{k+\ell}^* = \rho^\ell E |X_k|^2$$

Since we assumed the wide sense stationarity, $E|X_k|^2$ does not depend on $k$, i.e., we can assume

$$E|X_k|^2 = \sigma^2 = \text{const.}, \quad \text{and}$$

$$\Gamma_X = \sigma^2 \begin{bmatrix} 1 & \rho & \cdots & \rho^{n-1} \\ \rho & 1 & \ddots & \vdots \\ \vdots & \ddots & \ddots & \rho \\ \rho^{n-1} & \cdots & \rho & 1 \end{bmatrix} \Rightarrow \text{Toeplitz matrix!}$$

For further simplicity, let's assume $\sigma^2 = 1$. One can show that The eigenvalues & eigenvectors of this Toeplitz matrix are:

$$\begin{cases} \lambda_k = \dfrac{1 - \rho^2}{1 - 2\rho\cos\omega_k + \rho^2}, & k = 0, 1, \cdots, n-1. \\[3mm] \omega_k[\ell] = \sqrt{\dfrac{2}{n + \lambda_k}} \sin\left[\omega_k\left(\ell - \left(\tfrac{n-1}{2}\right)\right) + (k+1)\tfrac{\pi}{2}\right] \end{cases}$$

$$\ell = 0, 1, \cdots, n-1.$$

where $\omega_k$ is a solution of the following secular equation:

$$\tan(n\omega_k) = -\frac{(1 - \rho^2)\sin\omega_k}{(1 + \rho^2)\cos\omega_k - 2\rho}$$

Now consider the case when $\rho \uparrow 1$. Then $\tan n\omega_k = 0 \iff \omega_k = \frac{k\pi}{n}, \; k \neq 0$

For $k = 0$, we use the small angle perturbation: $\tan\alpha \approx \alpha$, $\sin\alpha \approx \alpha$, $\cos\alpha \approx 1 - \frac{\alpha^2}{2}$ using these, we get

$$n \omega_0' = - \frac{(1-\rho^2) \omega_0}{(1+\rho^2)(1 - \omega_0^2/2) - 2\rho}$$

$$\Leftrightarrow \omega_0^2 = \frac{2}{1+\rho^2} \cdot \frac{1-\rho^2}{n} + \frac{2(1-\rho)^2 \rho}{1+\rho^2}$$

$$\approx \frac{1}{n}(1-\rho^2) \quad \text{as } \rho \uparrow 1.$$

So, one can say $\omega_0 \to 0$ as $\rho \uparrow 1$.
after all, $\omega_k = k\pi/n$, $k = 0, 1, \cdots; n-1$.

Now, $\widetilde{w_k}[l] = \sqrt{\frac{2}{n+\lambda_k}} \sin\left[\frac{k\pi}{n}\left(l - \left(\frac{n-1}{2}\right)\right) + (k+1)\frac{\pi}{2}\right]$

$$= \sqrt{\frac{2}{n+\lambda_k}} \sin\left[\frac{\pi k}{n}\left(l + \frac{1}{2}\right) + \frac{\pi}{2}\right]$$

$$= \sqrt{\frac{2}{n+\lambda_k}} \cos\left[\frac{\pi k}{n}\left(l + \frac{1}{2}\right)\right]$$

How about $\lambda_k$ as $\rho \uparrow 1$?
For $k \neq 0$, $\lambda_k = \frac{1-\rho^2}{1 - 2\rho \cos\omega_k + \rho^2} \to 0$ as $\rho \uparrow 1$

For $k = 0$, $\lambda_0 \approx \frac{1-\rho^2}{1 - 2\rho(1 - \omega_0^2/2) + \rho^2}$

$$= \frac{1+\rho}{1 - \rho + \rho(1+\rho)/n}$$

$$\to n \quad \text{as } \rho \uparrow 1.$$

So, $\begin{cases} \widetilde{w_0}[l] = \frac{1}{\sqrt{n}} \\ \widetilde{w_k}[l] = \sqrt{\frac{2}{n}} \cos\left[\frac{\pi k(l + \frac{1}{2})}{n}\right] \end{cases}$ DCT-II !!