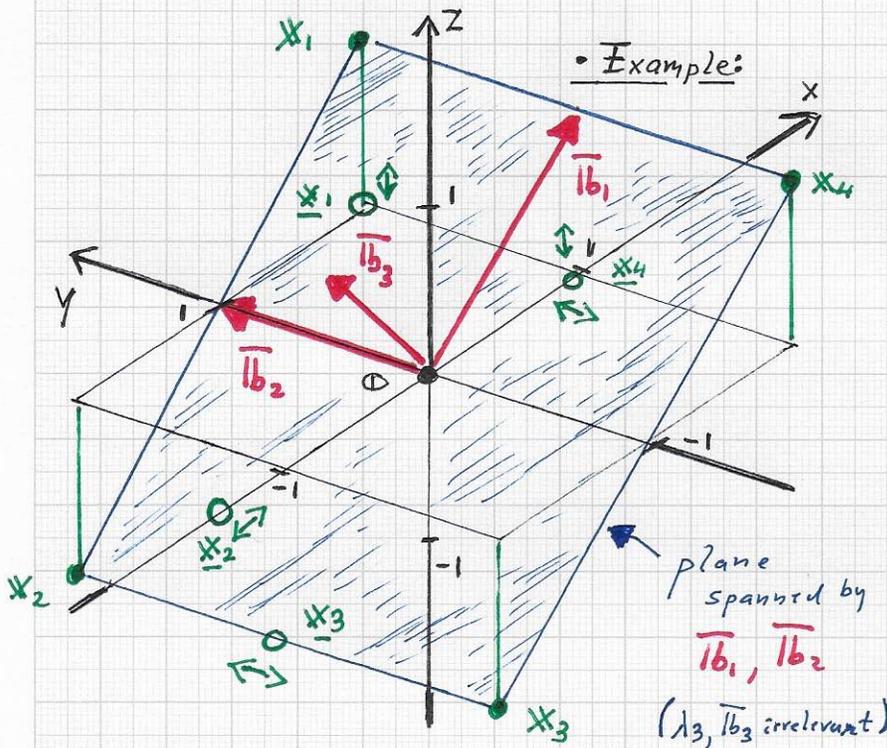


Stratovan

■ DIMENSION REDUCTION - Cont'd.

• Missing data II: A given sample data set $\{x_i\}$ includes several data with missing coordinate values.

An approximation of a PCA basis $\{\bar{b}_j\}$ can be constructed and improved iteratively.



→ complete data set:

$$\{x_i\} = \left\{ \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}, \begin{pmatrix} -1 \\ 1 \\ 1 \end{pmatrix}, \begin{pmatrix} -1 \\ -1 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 \\ -1 \\ 1 \end{pmatrix} \right\}$$

→ introduce missing values, '?':

$$\{x_i\} = \left\{ \begin{pmatrix} 1 \\ ? \\ 1 \end{pmatrix}, \begin{pmatrix} ? \\ 1 \\ 1 \end{pmatrix}, \begin{pmatrix} -1 \\ ? \\ 1 \end{pmatrix}, \begin{pmatrix} 1 \\ ? \\ 1 \end{pmatrix} \right\}$$

→ PCA for complete data set:

$$C = \begin{pmatrix} 4 & 0 & 4 \\ 0 & 4 & 0 \\ 4 & 0 & 4 \end{pmatrix};$$

$$\lambda_1 = 8, \lambda_2 = 4, (\lambda_3 = 0);$$

$$\bar{b}_1 = \begin{pmatrix} \sqrt{2}/2 \\ 0 \\ \sqrt{2}/2 \end{pmatrix}, \bar{b}_2 = \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix},$$

$$\bar{b}_3 = \begin{pmatrix} -\sqrt{2}/2 \\ 0 \\ \sqrt{2}/2 \end{pmatrix}$$

→ step 1: estimate a center C for the data set with missing values; i.e., compute averages of the known values, coordinate-wise:

$$C \approx \begin{pmatrix} 1/3 & (1-1+1) \\ 1/2 & (1+1) \\ 1/2 & (-1-1) \end{pmatrix} = \begin{pmatrix} 1/3 \\ 1 \\ -1 \end{pmatrix}$$

→ step 2: use the coordinate values of C to replace the '?' values; i.e., "repair" the x_i data with missing values using values of C:

$$x_1 \approx \begin{pmatrix} 1 \\ 1/3 \\ 1 \end{pmatrix}, x_2 \approx \begin{pmatrix} 1/3 \\ 1 \\ 1 \end{pmatrix}, x_3 \approx \begin{pmatrix} -1 \\ 1 \\ 1 \end{pmatrix}, x_4 \approx \begin{pmatrix} 1 \\ 1 \\ -1 \end{pmatrix}$$

Stratovan

■ DIMENSION REDUCTION - Cont' d.

• Missing data II: Example continued...

→ step 3: compute a PCA-based orthonormal basis $\{\bar{b}_j\}$;

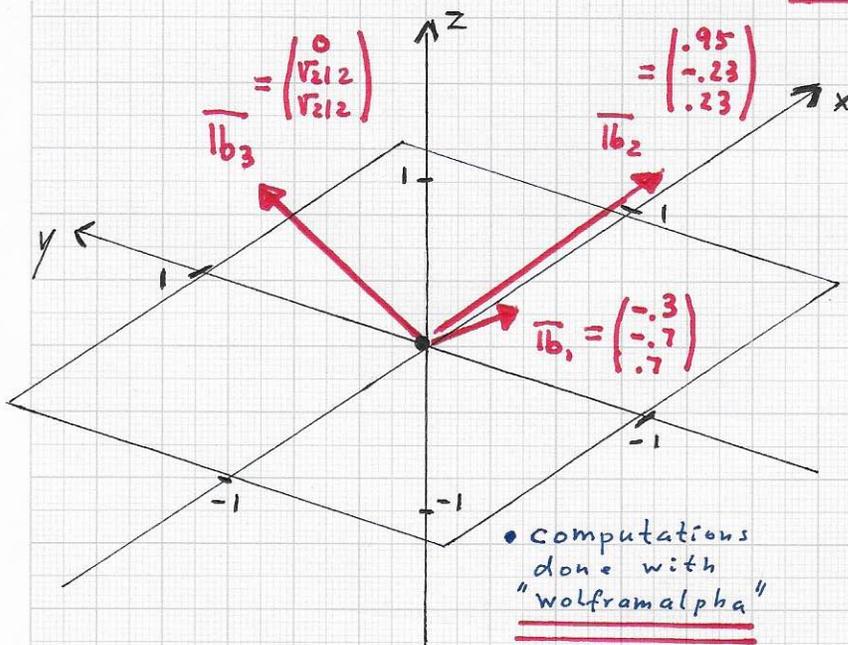
i.e., use the current estimates for $\{x_i\}$ to perform PCA:

$$C = \begin{pmatrix} 3/9 & 1/3 & -1/3 \\ 1/3 & 4 & -4 \\ -1/3 & -4 & 4 \end{pmatrix}; \quad \lambda_1 = \frac{2}{9} (25 + \sqrt{193}) \approx 8.64, \\ \lambda_2 = \frac{2}{9} (25 - \sqrt{193}) \approx 2.47, \\ (\lambda_3 = 0)$$

⇒ not normalized eigenvectors: $\phi_1 = \begin{pmatrix} \frac{1}{6} (11 - \sqrt{193}) \\ -1 \\ 1 \end{pmatrix} \approx \dots,$

$\phi_2 = \begin{pmatrix} \frac{1}{6} (11 + \sqrt{193}) \\ -1 \\ 1 \end{pmatrix} \approx \dots, \quad \phi_3 = \begin{pmatrix} 0 \\ 1 \\ 1 \end{pmatrix}$

⇒ orthonormal basis vectors: $\bar{b}_1 \approx \begin{pmatrix} -.32 \\ -.67 \\ .67 \end{pmatrix}, \bar{b}_2 = \begin{pmatrix} .95 \\ -.23 \\ .23 \end{pmatrix}, \quad \bar{b}_3 = \begin{pmatrix} 0 \\ \sqrt{2}/2 \\ \sqrt{2}/2 \end{pmatrix}$



• computations done with "wolframalpha"

→ step 4: use the current \bar{b}_j estimates to perform best approximation of the given $\{x_i\}$ to improve the estimates for all '?' coord. values;

i.e., perform the method for "Missing data I" for all x_i data with '?' values:

⇒ best approx. of $x_i = \left(\langle \bar{b}_i, \bar{b}_j \rangle \right) \begin{pmatrix} \bar{x}_i \\ \bar{z}_i \end{pmatrix} = \left(\langle x_i, \bar{b}_i \rangle \right)$ with '0' indicating that only the first two coords. are used.

Stratovan■ DIMENSION REDUCTION - Cont'd.• Missing data II: Example continued...

$$\dots \Rightarrow \text{best approx. of } \underline{x}_1: \begin{pmatrix} \langle \bar{b}_1, \bar{b}_1 \rangle & \dots & \langle \bar{b}_1, \bar{b}_3 \rangle \\ \vdots & & \vdots \\ \langle \bar{b}_3, \bar{b}_1 \rangle & \dots & \langle \bar{b}_3, \bar{b}_3 \rangle \end{pmatrix} \begin{pmatrix} \bar{x}_1 \\ \bar{y}_1 \\ \bar{z}_1 \end{pmatrix} = \begin{pmatrix} \langle \underline{x}_1, \bar{b}_1 \rangle \\ \vdots \\ \langle \underline{x}_1, \bar{b}_3 \rangle \end{pmatrix}$$

$$\Leftrightarrow \begin{pmatrix} .5513 & -.1499 & -.4738 \\ -.1499 & .9554 & -.1626 \\ -.4738 & -.1626 & .5 \end{pmatrix} \begin{pmatrix} \bar{x}_1 \\ \bar{y}_1 \\ \bar{z}_1 \end{pmatrix} = \begin{pmatrix} -.99 \\ .72 \\ .7071 \end{pmatrix}$$

$$\Rightarrow \begin{pmatrix} \bar{x}_1 \\ \bar{y}_1 \\ \bar{z}_1 \end{pmatrix} = \begin{pmatrix} -.0261 \\ 1.0438 \\ 1.7289 \end{pmatrix}$$

$$\Rightarrow \underline{x}_1 = -.0261 \begin{pmatrix} -.32 \\ -.67 \\ .67 \end{pmatrix} + 1.0438 \begin{pmatrix} .95 \\ -.23 \\ .23 \end{pmatrix} + 1.7289 \begin{pmatrix} 0 \\ \sqrt{2}/2 \\ \sqrt{2}/2 \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \\ 1.445 \end{pmatrix}$$

(only computation of third coord. value, 1.445, needed!)

\Rightarrow best approx. of $\underline{x}_2, \underline{x}_3, \underline{x}_4$: perform same computations!

- \rightarrow Iteration:
- perform step 3, using $\{\underline{x}_i\}$ to improve $\{\bar{b}_j\}$;
 - perform step 4, using $\{\bar{b}_j\}$ to best-approximate $\{\underline{x}_i\}$;
 - stop when orthonormal basis $\{\bar{b}_j\}$ has "converged."

-
- Note: It is unclear whether one should estimate a new data set center \mathbb{C} and perform mean-subtraction after every computation of updated \underline{x}_i -values.

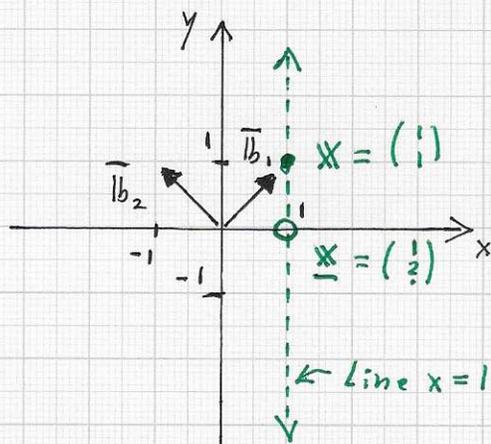
Intuitively, it seems to make sense and be desirable to do this, since the mean indeed changes after every iteration.

- Reference: Erverson and Sirovich, The Karhunen-Loève transform for incomplete data, J. Opt. Soc. Am., A, pp. 1657-1664, 1995.

Stratovan

■ DIMENSION REDUCTION - Cont'd.

- Missing data I and missing data II: simple examples



"Computing only missing ?-values when basis $\{\bar{b}_i\}$ is given, known"

• Ex.: missing data I

→ known orthonormal basis: $\bar{b}_1 = \begin{pmatrix} \sqrt{2}/2 \\ \sqrt{2}/2 \end{pmatrix}$, $\bar{b}_2 = \begin{pmatrix} -\sqrt{2}/2 \\ \sqrt{2}/2 \end{pmatrix}$

→ datum $\underline{x} = \begin{pmatrix} 1 \\ ? \end{pmatrix}$ misses y-coord. value

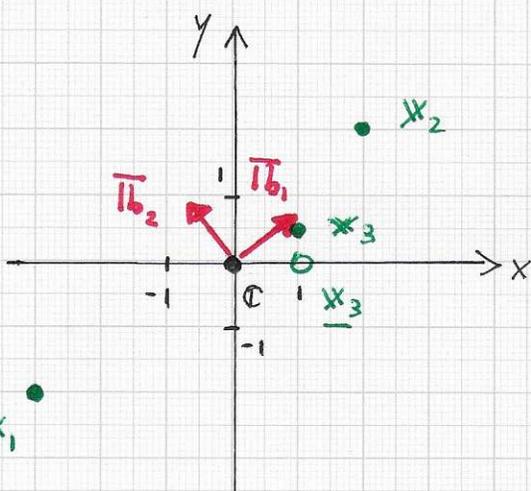
→ perform best approximation, using only the x-values for inner products:

$$\begin{pmatrix} \langle \bar{b}_1, \bar{b}_1 \rangle & \langle \bar{b}_1, \bar{b}_2 \rangle \\ \langle \bar{b}_2, \bar{b}_1 \rangle & \langle \bar{b}_2, \bar{b}_2 \rangle \end{pmatrix} \begin{pmatrix} \bar{x} \\ \bar{y} \end{pmatrix} = \begin{pmatrix} \langle \underline{x}, \bar{b}_1 \rangle \\ \langle \underline{x}, \bar{b}_2 \rangle \end{pmatrix}$$

$$\Leftrightarrow \begin{pmatrix} 1/2 & -1/2 \\ -1/2 & 1/2 \end{pmatrix} \begin{pmatrix} \bar{x} \\ \bar{y} \end{pmatrix} = \begin{pmatrix} \sqrt{2}/2 \\ -\sqrt{2}/2 \end{pmatrix}$$

⇒ solution family: $\bar{y} = \sqrt{2} - \bar{x}$

$$\Rightarrow \underline{x} = \bar{x} \begin{pmatrix} \sqrt{2}/2 \\ \sqrt{2}/2 \end{pmatrix} + (\sqrt{2} - \bar{x}) \begin{pmatrix} -\sqrt{2}/2 \\ \sqrt{2}/2 \end{pmatrix} = \begin{pmatrix} \bar{x}\sqrt{2} - 1 \\ 1 \end{pmatrix} \stackrel{\bar{x}=1}{=} \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$



"Iterative computation of unknown basis vectors and \underline{x}_i data with ?-values"

• Ex.: missing data II

→ orthonormal basis NOT known

→ given data: $\underline{x}_1 = \begin{pmatrix} -3 \\ -2 \end{pmatrix}$, $\underline{x}_2 = \begin{pmatrix} 2 \\ 2 \end{pmatrix}$, $\underline{x}_3 = \begin{pmatrix} 1 \\ ? \end{pmatrix}$

→ Center (mean): $\underline{c} = \begin{pmatrix} 1/3 \cdot 0 \\ 1/2 \cdot 0 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$

(⇒ subsequent mean-subtraction step would not change data!)

→ replace ?-value by mean-value:

$$\underline{x}_1 = \begin{pmatrix} -3 \\ -3 \end{pmatrix}, \underline{x}_2 = \begin{pmatrix} 2 \\ 2 \end{pmatrix}, \underline{x}_3 = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$$

Stratovan■ DIMENSION REDUCTION - Cont'd.• Ex.: missing data II ...

→ first approximation of orthonormal basis:

$$C = \begin{pmatrix} -3 & 2 & 1 \\ -3 & 2 & 0 \end{pmatrix} \begin{pmatrix} -3 & -3 \\ 2 & 2 \\ 1 & 0 \end{pmatrix} = \begin{pmatrix} 14 & 10 \\ 10 & 8 \end{pmatrix} \Rightarrow \begin{aligned} \lambda_1 &= 11 + \sqrt{109} = 21.44 \\ \lambda_2 &= 11 - \sqrt{109} = 0.56 \end{aligned}$$

$$\Rightarrow \bar{b}_1 \approx \begin{pmatrix} .8 \\ .6 \end{pmatrix}, \bar{b}_2 \approx \begin{pmatrix} -.6 \\ .8 \end{pmatrix}$$

→ HERE: $\lambda_1 \gg \lambda_2 \Rightarrow \bar{b}_1$ "important", \bar{b}_2 "irrelevant"!→ compute next approximation for \bar{x} -value of \underline{x}_3 :perform inner product computations only for x-coordinate, considering basis vector \bar{b}_1 only:

$$\left(\langle \bar{b}_1, \bar{b}_1 \rangle \right) (\bar{x}) = \left(\langle \underline{x}_3, \bar{b}_1 \rangle \right)$$

$$\Leftrightarrow .64 \bar{x} = 1.8$$

$$\Leftrightarrow \bar{x} = .8$$

$$\Rightarrow \underline{x}_3 = \bar{x} \bar{b}_1 = .8 \begin{pmatrix} .8 \\ .6 \end{pmatrix} = \begin{pmatrix} .64 \\ .48 \end{pmatrix}$$

$$\underline{x} = \begin{pmatrix} 1 \\ .48 \end{pmatrix}$$

this x-value is NOT considered; it is given as $x=1$

→ perform iterative improvement of basis vectors and originally unknown values for (?).

≈
BH