## Eigen-decomposition of the covariance matrix

Example data



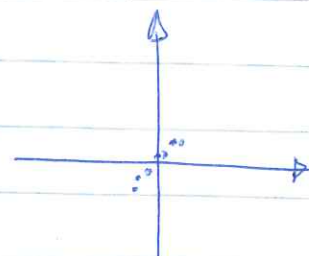| Systolic BP | Diastolic BP |
|---|---|
| 126 | 78 |
| 128 | 80 |
| 128 | 82 |
| 130 | 82 |
| 130 | 84 |
| 132 | 86 |

We will use PCA to combine the two blood pressure variables into just one variable based on data from 6 individuals.

1. Center the data
2. Calculate the covariance matrix (CM)
3. Calculate eigen values of the CM
4.      "      eigen vectors of the CM
5. Order the eigen vectors
6. Calculate the principal components (PCs)

Note that sometimes we have to standardize the data by

$$\frac{x - \mu}{\sigma}$$

Step 1:

| SBP | DBP |
|---|---|
| 126 − 129 = −3 | 78 − 82 = −4 |
| 128 − 129 = −1 | 80 − 82 = −2 |
| 128 − 129 = −1 | 82 − 82 = 0 |
| 130 − 129 = 1 | 82 − 82 = 0 |
| 130 − 129 = 1 | 84 − 82 = 2 |
| 132 − 129 = 3 | 86 − 82 = 4 |

⇒



data centered around (0,0).

SBP $\Rightarrow$ cSBP

DBP $\Rightarrow$ cDBP

Step 2: Calculate the CM

$$\hat{\sigma} = \begin{array}{c} \\ cSBP \\ cDBP \end{array} \overset{\displaystyle \overset{cSBP \qquad cDBP}{}}{\begin{pmatrix} 4.4 & 5.6 \\ 5.6 & 8.0 \end{pmatrix}}$$

$\hookrightarrow$ variances (spread in cDBP is higher than in cSBP)

covariances

$$Var(cSBP) = \frac{1}{M-1} \sum_{i=1}^{M} (cSBP_i - \overline{cSBP})^2$$

$$Var(cDBP) = \frac{1}{M-1} \sum_{i=1}^{M} (cDBP_i - \overline{cDBP})^2$$

$$\sigma_{xy} = \frac{1}{M-1} \sum_{i=1}^{M} (cSBP_i - \overline{cSBP})(cDBP_i - \overline{cDBP})$$

Step 3. $\det | \hat{\sigma} - \lambda \hat{I} | = 0$

$$\det \begin{vmatrix} 4.4 - \lambda & 5.6 \\ 5.6 & 8 - \lambda \end{vmatrix} = 0 \qquad \Rightarrow \quad 3.84 - 12.4\lambda + \lambda^2 = 0$$

$$\boxed{\begin{array}{l} \lambda_1 = 0.32 \\ \lambda_2 = 12.08 \end{array}}$$

$\hookrightarrow$ eigen values of the CM.

→ eigenvectors!

Step 4.    $\hat{\sigma} \, v = \lambda \, v$

$$\begin{bmatrix} 4.4 & 5.6 \\ 5.6 & 8.0 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = \overset{\lambda_2}{\underset{\downarrow}{12.08}} \begin{bmatrix} x \\ y \end{bmatrix}$$

$\Rightarrow \quad 5.6y = 7.68x \qquad$ solving for $y \Rightarrow y = 1.37x$

$\qquad 5.6x = 4.08y \qquad\qquad\qquad$ for $x=1 \Rightarrow y = 1.37$

Therefore:  $\quad v_2 = \begin{bmatrix} 1 \\ 1.37 \end{bmatrix}$



After normalization:  $v_2 = \begin{bmatrix} 0.59 \\ 0.81 \end{bmatrix}$

For the other eigenvalue $\boxed{\lambda_1 = 0.32}$, we get $v_1 = \begin{bmatrix} -0.81 \\ 0.59 \end{bmatrix}$

(normalized)

Since the CM is a symmetric matrix, the eigenvectors will be
<u>ORTHOGONAL</u> :  $v_1 \perp v_2$

Step 5.  Ordering eigen vector

The eigen vector with the largest eigen value becomes our first eigenvector.

$V_2 \rightarrow V_{PC1}$     We order these eigenvectors in a matrix
$V_1 \rightarrow V_{PC2}$     called $\hat{V}$:

define a matrix $\hat{D}$ that has our centered data:

$$\hat{V} = \begin{bmatrix} 0.59 & -0.81 \\ 0.81 & 0.59 \end{bmatrix}$$

$V_{PC1}$   $V_{PC2}$

⇑ Principal Components

$$\hat{D} = \begin{bmatrix} -3 & -4 \\ -1 & -2 \\ -1 & 0 \\ 1 & 0 \\ 1 & 2 \\ 3 & 4 \end{bmatrix}$$
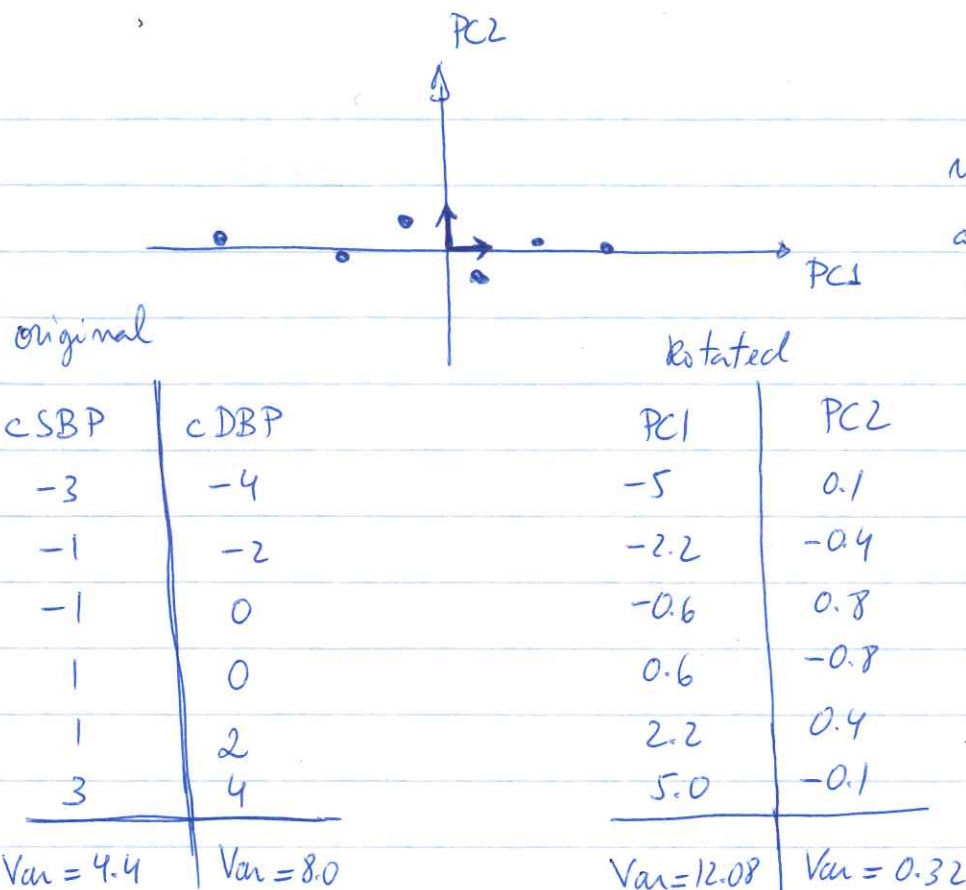
and make $\hat{D}\hat{V} =$

PC1   PC2
$$\begin{bmatrix} -5.0 & 0.1 \\ -2.2 & -0.4 \\ -0.6 & 0.8 \\ 0.6 & -0.8 \\ 2.2 & 0.4 \\ 5.0 & -0.1 \end{bmatrix}$$

transformed data! ⇒ Principal component Scores!

This represents the original centered data in the PC space !

# ROTATION

PC2

PC1

rotated plot
also named
score plot!

| original | | | Rotated | |
|---|---|---|---|---|
| cSBP | cDBP | | PC1 | PC2 |
| −3 | −4 | | −5 | 0.1 |
| −1 | −2 | | −2.2 | −0.4 |
| −1 | 0 | | −0.6 | 0.8 |
| 1 | 0 | | 0.6 | −0.8 |
| 1 | 2 | | 2.2 | 0.4 |
| 3 | 4 | | 5.0 | −0.1 |
| Var = 4.4 | Var = 8.0 | | Var = 12.08 | Var = 0.32 |

eigen values

and variance
difference
increased!

$$\% \, Var = \frac{\lambda_{PC1}}{\lambda_{PC1} + \lambda_{PC2}} = \frac{12.08}{12.08 + 0.32} = 97.4\%$$

PC1 captures 97.4 % of the total variance of the data!

$$\hat{\sigma}_{PC} = \begin{pmatrix} 12.08 & 0 \\ 0 & 0.32 \end{pmatrix}$$

Rotation

Transformed data

$$PC1 = 0.59 \, cSBP + 0.81 \, cDBP$$
$$PC2 = -0.81 \, cSBP + 0.59 \, cDBP$$

centred data

For example: PC score for person #6: $PC1_6 = 0.59 \times \boxed{3} + 0.81 \times \boxed{4} = 5$

But how variable reduction play a role in PCA?
We still have the same # of ~~variables~~ variables as of PCs.

Since the 1º PC captures ≈97% of all variance (carries most of the information about the data), we can neglect PC2.
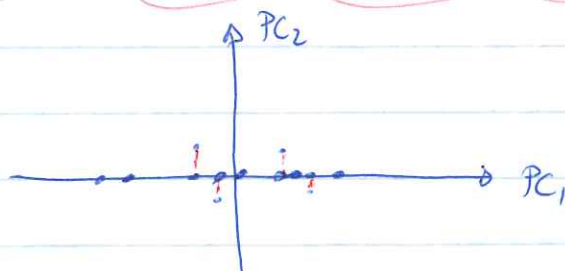
$$PC1 = 0.59 \, cSBP + 0.81 \, cDBP$$

We are combining the two variables cSBP and cDBP into 1 variable, the PC1, in a way that maximize the variance of the linear combination.

The weights tell how much each variable contributes to the PC.

$$W_{cDBP} > W_{cSBP} \; : \quad \text{PCA puts more weight into cDBP}$$
when the 2 variables are combined.

(rotates)
NOTE: The covariance matrix transforms any vector into the direction of the eigenvector of largest eigenvalue or variance!



ignoring PC2, data is projected into PC1.