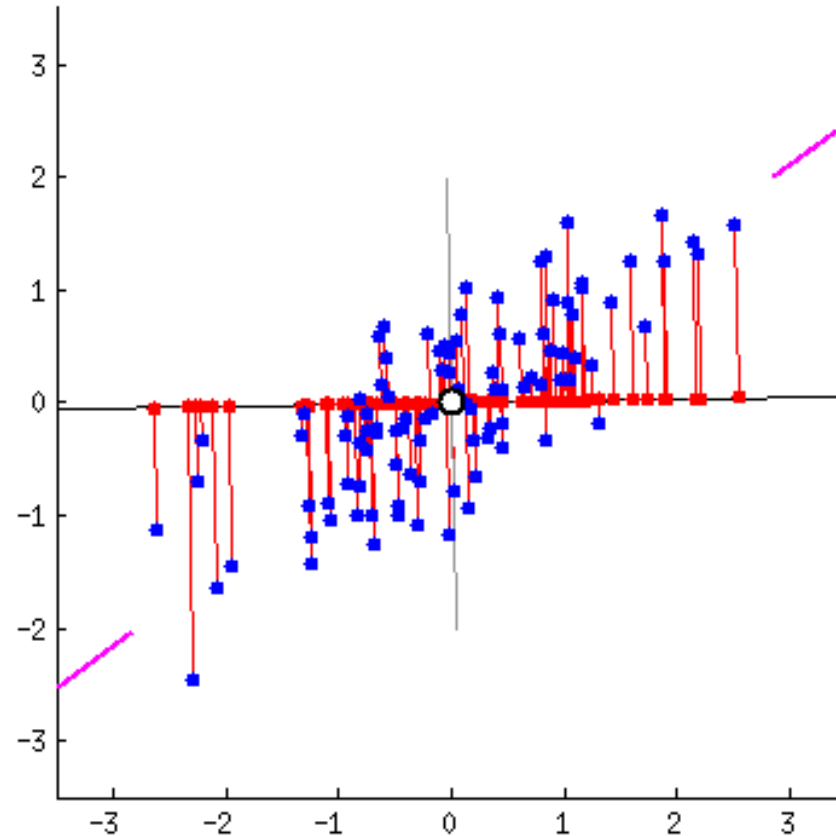


Principle Component Analysis

Jing Qin

03/05/2023

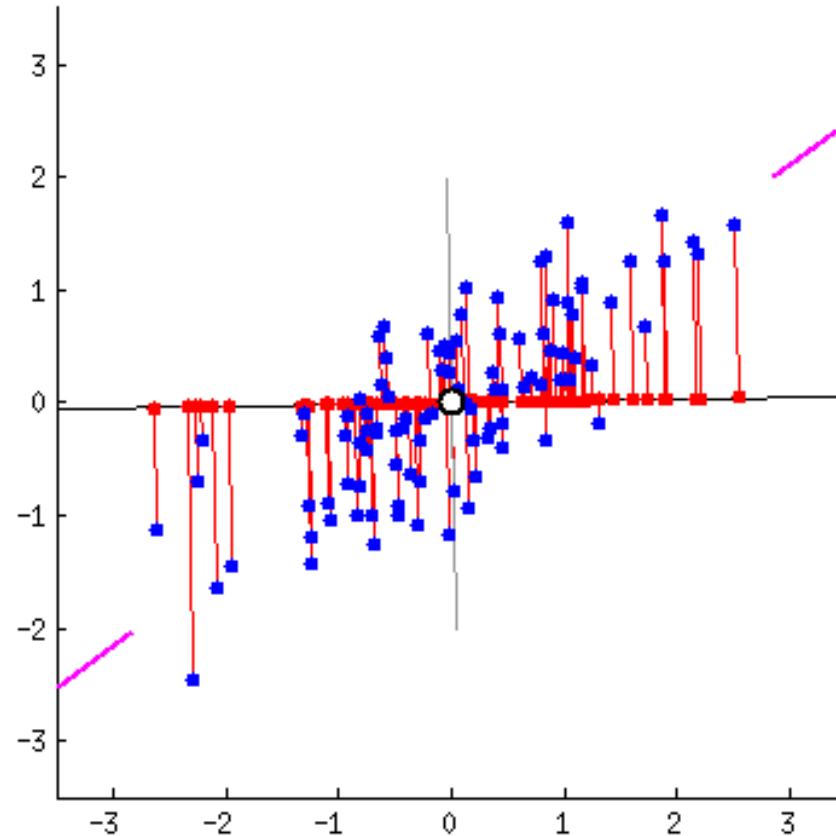
What is PCA? Start with a movie...



The **first** principal component accounts for the **largest possible variance** in the data set.

The second principal component is calculated as the one that is **uncorrelated with** (perpendicular to) the first principal component and that it accounts for the **next highest variance**...

What is PCA? Start with a movie...



How to do
this?

The **first** principal component accounts for the **largest possible variance** in the data set.

The second principal component is calculated as the one that is **uncorrelated with** (perpendicular to) the first principal component and that it accounts for the **next highest variance**...

Main result:

This is the answer to the 'How' question.
Step-by-step gets into it ~

Principle components are **linear combinations** of the random variables, and the coefficients are determined by the **eigenvectors** of covariance matrix.

Result 8.1. Let Σ be the covariance matrix associated with the random vector $\mathbf{X}' = [X_1, X_2, \dots, X_p]$. Let Σ have the eigenvalue-eigenvector pairs $(\lambda_1, \mathbf{e}_1), (\lambda_2, \mathbf{e}_2), \dots, (\lambda_p, \mathbf{e}_p)$ where $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0$. Then the i th principal component is given by

$$Y_i = \mathbf{e}_i' \mathbf{X} = e_{i1}X_1 + e_{i2}X_2 + \dots + e_{ip}X_p, \quad i = 1, 2, \dots, p \quad (8-4)$$

With these choices,

$$\begin{aligned} \text{Var}(Y_i) &= \mathbf{e}_i' \Sigma \mathbf{e}_i = \lambda_i & i = 1, 2, \dots, p \\ \text{Cov}(Y_i, Y_k) &= \mathbf{e}_i' \Sigma \mathbf{e}_k = 0 & i \neq k \end{aligned} \quad (8-5)$$

If some λ_i are equal, the choices of the corresponding coefficient vectors, \mathbf{e}_i , and hence Y_i , are not unique.

All the eigenvalues of Σ are non-negative (Σ is positive semi-definite)

- Σ has eigenvalues $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0$
- Alternative definition: For all \mathbf{y} , the quadratic form $\mathbf{y}' \cdot \Sigma \cdot \mathbf{y} \geq 0$.
- Proof of Σ is positive semi-definite:

Assume that $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n$ is a random sample and $\bar{\mathbf{X}}$ is the sample mean vector

$$\begin{aligned}\Sigma &= \frac{1}{n} \sum_{j=1}^n (\mathbf{X}_j - \bar{\mathbf{X}})(\mathbf{X}_j - \bar{\mathbf{X}})' \\ \mathbf{y}' \cdot \Sigma \cdot \mathbf{y} &= \frac{1}{n} \sum_{j=1}^n \mathbf{y}' (\mathbf{X}_j - \bar{\mathbf{X}})(\mathbf{X}_j - \bar{\mathbf{X}})' \mathbf{y} \\ &= \frac{1}{n} \sum_{j=1}^n ((\mathbf{X}_j - \bar{\mathbf{X}})' \mathbf{y})^2 \geq 0\end{aligned}$$

Inner product (2-4)

In sum...

Original components

Principle components are p linear combinations of X_1, X_2, \dots, X_p

$$Y_1 = \mathbf{a}'_1 \mathbf{X} = a_{11}X_1 + a_{12}X_2 + \dots + a_{1p}X_p$$

Principle components

$$Y_2 = \mathbf{a}'_2 \mathbf{X} = a_{21}X_1 + a_{22}X_2 + \dots + a_{2p}X_p$$

...

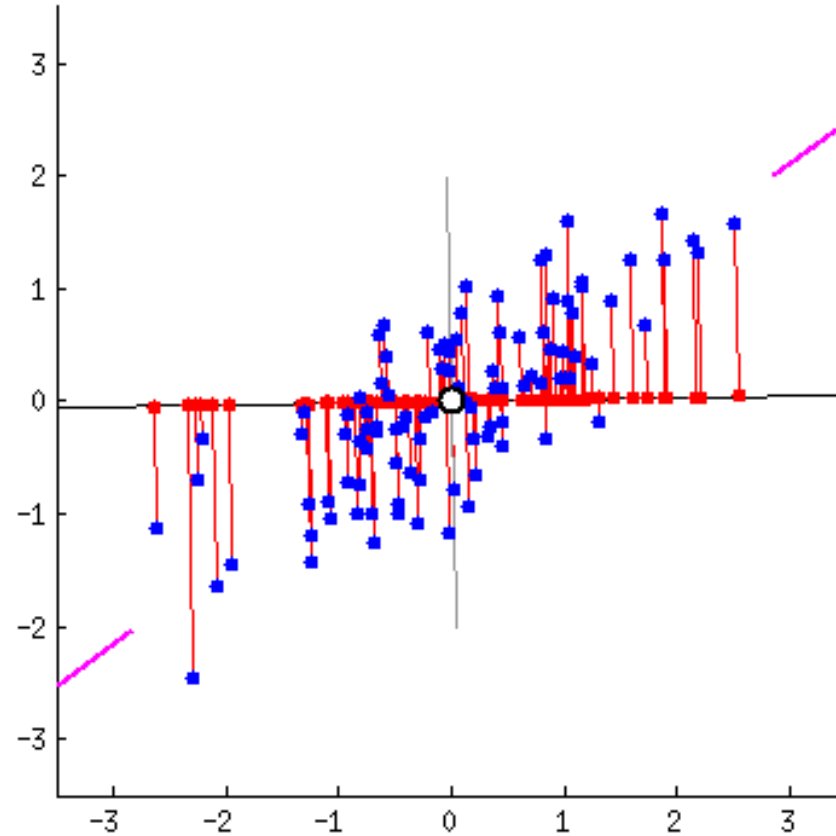
$$Y_p = \mathbf{a}'_p \mathbf{X} = a_{p1}X_1 + a_{p2}X_2 + \dots + a_{pp}X_p$$

satisfying

- variance as large as possible $Var(Y_1) \geq Var(Y_2) \geq \dots \geq Var(Y_p)$
- uncorrelated $Cov(Y_j, Y_k) = 0$ for any pair (j, k)
- “standardized” $\mathbf{a}'_i \cdot \mathbf{a} = 1$, where $\mathbf{a}'_i = (a_{i1}, a_{i2}, \dots, a_{ip})$

(Result 8.1) $\mathbf{a}_i = \mathbf{e}_i$, where $(\lambda_i, \mathbf{e}_i)$ is eigenvalue-eigenvector pair of Σ .

The movie, again



The **first** principal component accounts for the **largest possible variance** in the data set.

The second principal component is calculated as the one that is **uncorrelated with** (perpendicular to) the first principal component and that it accounts for the **next highest variance...**

How many PCs we need (not all of them)

Result 8.2. Let $\mathbf{X}' = [X_1, X_2, \dots, X_p]$ have covariance matrix Σ , with eigenvalue-eigenvector pairs $(\lambda_1, \mathbf{e}_1), (\lambda_2, \mathbf{e}_2), \dots, (\lambda_p, \mathbf{e}_p)$ where $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0$. Let $Y_1 = \mathbf{e}_1' \mathbf{X}, Y_2 = \mathbf{e}_2' \mathbf{X}, \dots, Y_p = \mathbf{e}_p' \mathbf{X}$ be the principal components. Then

$$\sigma_{11} + \sigma_{22} + \dots + \sigma_{pp} = \sum_{i=1}^p \text{Var}(X_i) = \lambda_1 + \lambda_2 + \dots + \lambda_p = \sum_{i=1}^p \text{Var}(Y_i)$$

Result 8.2 says that

$$\begin{aligned} \text{Total population variance} &= \sigma_{11} + \sigma_{22} + \dots + \sigma_{pp} \\ &= \lambda_1 + \lambda_2 + \dots + \lambda_p \end{aligned} \quad (8-6)$$

and consequently, the proportion of total variance due to (explained by) the k th principal component is

$$\left(\begin{array}{c} \text{Proportion of total} \\ \text{population variance} \\ \text{due to } k\text{th principal} \\ \text{component} \end{array} \right) = \frac{\lambda_k}{\lambda_1 + \lambda_2 + \dots + \lambda_p} \quad k = 1, 2, \dots, p \quad (8-7)$$

≥ 95%

How many PCs we need (not all of them)

Result 8.2 says that

$$\begin{aligned}\text{Total population variance} &= \sigma_{11} + \sigma_{22} + \cdots + \sigma_{pp} \\ &= \lambda_1 + \lambda_2 + \cdots + \lambda_p\end{aligned}\quad (8-6)$$

and consequently, the proportion of total variance due to (explained by) the k th principal component is

$$\left(\begin{array}{c} \text{Proportion of total} \\ \text{population variance} \\ \text{due to } k\text{th principal} \\ \text{component} \end{array} \right) = \frac{\lambda_k}{\lambda_1 + \lambda_2 + \cdots + \lambda_p} \quad k = 1, 2, \dots, p \quad (8-7)$$

- Proportion of total variance explained by the first k component

$$\frac{\lambda_1 + \cdots + \lambda_k}{\lambda_1 + \lambda_2 + \cdots + \lambda_p}$$

How to do PCA in R with `prcomp()`

```
> pca <- prcomp(df)
```

```
> pca
```

Standard deviations (1, ..., p=3):

```
[1] 0.15265434 0.02446027 0.01896934
```

Rotation (n x k) = (3 x 3):

	PC1	PC2	PC3
V1	0.6831023	-0.1594791	0.7126974
V2	0.5102195	-0.5940118	-0.6219534
V3	0.5225392	0.7884900	-0.3244015

```
> summary(pca)
```

Importance of components:

	PC1	PC2	PC3
Standard deviation	0.1527	0.02446	0.01897
Proportion of Variance	0.9605	0.02466	0.01483
Cumulative Proportion	0.9605	0.98517	1.00000

How many PCs we need (scree plot) with
`screeplot()`

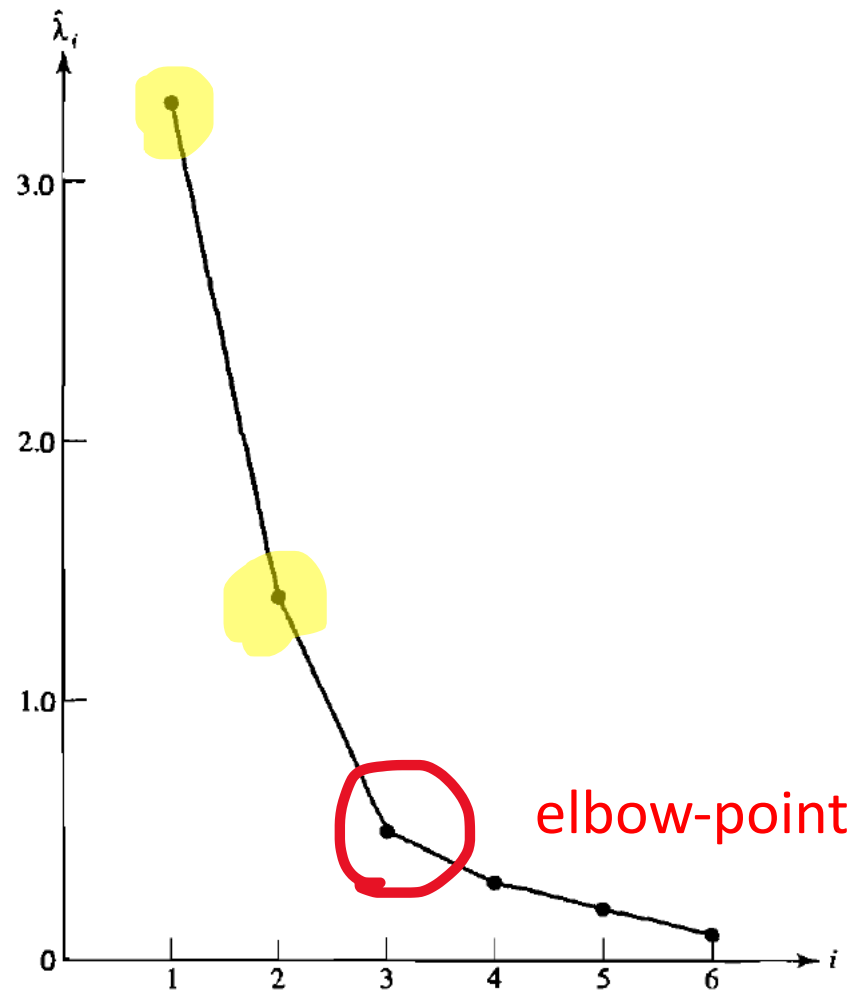


Figure 8.2 A scree plot.



Standardization (PCA on correlation matrix $\boldsymbol{\rho}$ not $\boldsymbol{\Sigma}$)

```
pca2 <- prcomp(df, scale=T)
```

by (2-37). The principal components of \mathbf{Z} may be obtained from the eigenvectors of the *correlation* matrix $\boldsymbol{\rho}$ of \mathbf{X} . All our previous results apply, with some simplifications, since the variance of each Z_i is unity. We shall continue to use the notation Y_i to refer to the i th principal component and $(\lambda_i, \mathbf{e}_i)$ for the eigenvalue–eigenvector pair from either $\boldsymbol{\rho}$ or $\boldsymbol{\Sigma}$. *However, the $(\lambda_i, \mathbf{e}_i)$ derived from $\boldsymbol{\Sigma}$ are, in general, not the same as the ones derived from $\boldsymbol{\rho}$.*

Result 8.4. The i th principal component of the standardized variables $\mathbf{Z}' = [Z_1, Z_2, \dots, Z_p]$ with $\text{Cov}(\mathbf{Z}) = \boldsymbol{\rho}$, is given by

$$Y_i = \mathbf{e}_i' \mathbf{Z} = \mathbf{e}_i' (\mathbf{V}^{1/2})^{-1} (\mathbf{X} - \boldsymbol{\mu}), \quad i = 1, 2, \dots, p$$

Moreover,

$$\sum_{i=1}^p \text{Var}(Y_i) = \sum_{i=1}^p \text{Var}(Z_i) = p$$

(8-11) that the total (standardized variables) population variance is the sum of the diagonal elements of the matrix $\boldsymbol{\rho}$. Using (8-7) with \mathbf{Z} in

place of \mathbf{X} , we find that the proportion of total variance explained by the k th principal component of \mathbf{Z} is

$$\left(\begin{array}{c} \text{Proportion of (standardized)} \\ \text{population variance due} \\ \text{to } k\text{th principal component} \end{array} \right) = \frac{\lambda_k}{p}, \quad k = 1, 2, \dots, p \quad (8-12)$$

where the λ_k 's are the eigenvalues of $\boldsymbol{\rho}$.

MANGE TUSIND
AF HJERTET
TIL DIG
FRA MIG

TAK tak

FOR OPMÆRKSOMHEDEN

ER KUN ET FATTIGT ORD

FOR HJÆLPEN

FOR SIDST

FORDI DU ER DIG

FOR ALT

FOR INDBYDELSEN

TAK

LARS LILHOLT BAND



HELD OG LYKKE