

Lecture 3: Principle Component Analysis (PCA)

Reading: Section 14.5

GU4241/GR5241 Statistical Machine Learning

Linxí Liu
January 26, 2018

Quadratic Forms 二次型

In applications, symmetric matrices often occur in quadratic forms.

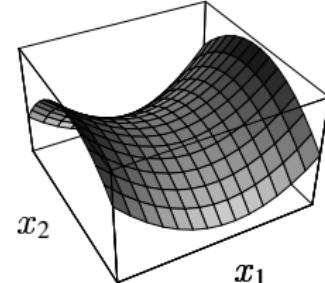
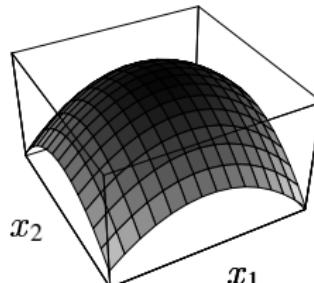
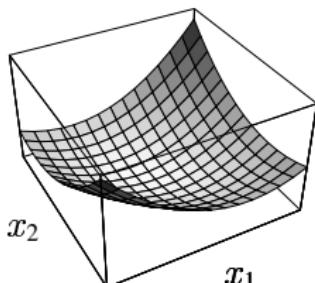
Definition

The **quadratic form** defined by a matrix A is the function

$$\boxed{q_A : \mathbb{R}^m \rightarrow \mathbb{R}} \\ x \mapsto \langle x, Ax \rangle$$

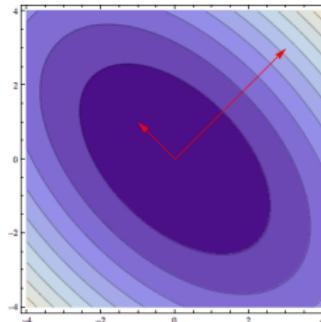
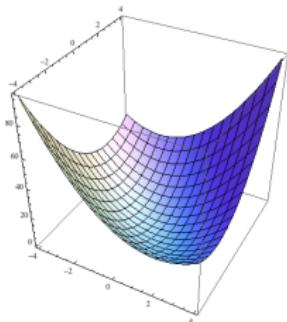
Intuition

A quadratic form is the m -dimensional analogue of a quadratic function ax^2 , with a vector substituted for the scalar x and the matrix A substituted for the scalar $a \in \mathbb{R}$.



Quadratic Forms

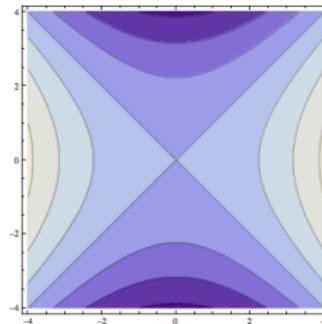
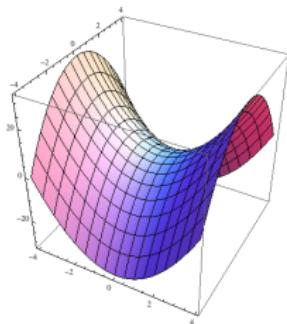
Here is the quadratic form for the matrix $A = \begin{pmatrix} 2 & 1 \\ 1 & 2 \end{pmatrix}$:



- ▶ Left: The function value q_A is graphed on the vertical axis.
- ▶ Right: Each line corresponds to a constant function value of q_A . Dark color = small values.
- ▶ The red lines are eigenvector directions of A . Their lengths represent the (absolute) values of the eigenvalues.
- ▶ In this case, both eigenvalues are positive. If all eigenvalues are positive, the contours are ellipses. So:
positive definite matrices \leftrightarrow elliptic quadratic forms

Quadratic Forms

In this plot, the eigenvectors are axis-parallel, and one eigenvalue is negative:



The matrix here is $A = \begin{pmatrix} 2 & 0 \\ 0 & -2 \end{pmatrix}$.

Intuition

- ▶ If we change the sign of one of the eigenvalue, the quadratic function along the corresponding eigen-axis flips.
- ▶ There is a point which is a minimum of the function along one axis direction, and a maximum along the other. Such a point is called a *saddle point*.

Application: Covariance Matrix

Recall: Covariance

The covariance of two random variables X_1, X_2 is

$$\text{Cov}[X_1, X_2] = \mathbb{E}[(X_1 - \mathbb{E}[X_1])(X_2 - \mathbb{E}[X_2])].$$

If $X_1 = X_2$, the covariance is the variance: $\text{Cov}[X, X] = \text{Var}[X]$.

Covariance matrix

If $X = (X_1, \dots, X_m)$ is a random vector with values in \mathbb{R}^m , the matrix of all covariances

$$\text{Cov}[X] := (\text{Cov}[X_i, X_j])_{i,j} = \begin{pmatrix} \text{Cov}[X_1, X_1] & \cdots & \text{Cov}[X_1, X_m] \\ \vdots & & \vdots \\ \text{Cov}[X_m, X_1] & \cdots & \text{Cov}[X_m, X_m] \end{pmatrix}$$

is called the **covariance matrix** of X .

协方差阵

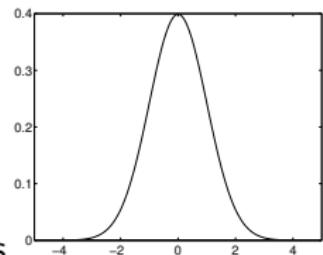
Notation

It is customary to denote the covariance matrix $\text{Cov}[X]$ by Σ .

Gaussian Distribution

Gaussian density in one dimension

$$p(x; \mu, \sigma) := \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right)$$



- ▶ μ = expected value of x , σ^2 = variance, σ = standard deviation
- ▶ The quotient $\frac{x-\mu}{\sigma}$ measures deviation of x from its expected value in units of σ (i.e. σ defines the length scale)

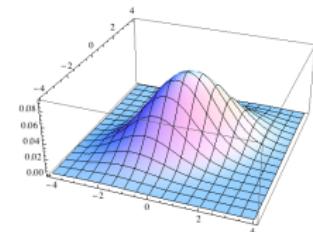
Gaussian density in m dimensions

The quadratic function

$$-\frac{(x - \mu)^2}{2\sigma^2} = -\frac{1}{2}(x - \mu)(\sigma^2)^{-1}(x - \mu)$$

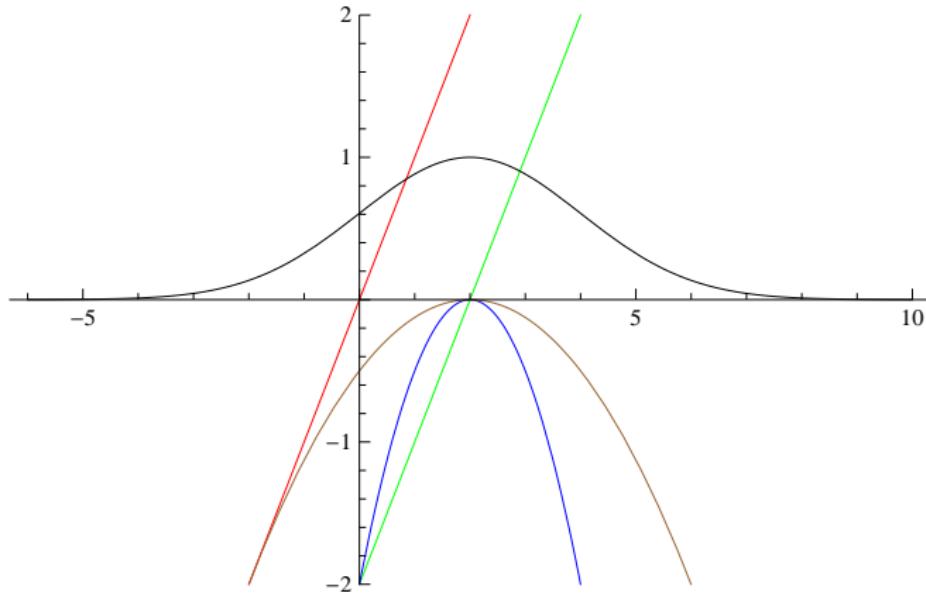
is replaced by a quadratic form:

$$p(\mathbf{x}; \boldsymbol{\mu}, \Sigma) := \frac{1}{\sqrt{2\pi \det(\Sigma)}} \exp\left(-\frac{1}{2} \langle (\mathbf{x} - \boldsymbol{\mu}), \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu}) \rangle\right)$$



Components of a 1D Gaussian

$$\mu = 2, \sigma = 2$$



- ▶ Red: $x \mapsto x$
- ▶ Green: $x \mapsto x - \mu$
- ▶ Blue: $x \mapsto -\frac{1}{2}(x - \mu)^2$

- ▶ Brown: $x \mapsto -\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2$
- ▶ Black: $x \mapsto \exp\left(-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right)$

Geometry of Gaussians

Covariance matrix of a Gaussian

If a random vector $X \in \mathbb{R}^m$ has Gaussian distribution with density $p(\mathbf{x}; \mu, \Sigma)$, its covariance matrix is $\text{Cov}[X] = \Sigma$. In other words, a Gaussian is parameterized by its covariance.

Observation

Since $\text{Cov}[X_i, X_j] = \text{Cov}[X_j, X_i]$, the covariance matrix is symmetric.

What is the eigenstructure of Σ ?

- We know: Σ symmetric \Rightarrow there is an eigenvector ONB
任何一个矩阵的特征向量都可以组成一组正交基
- Call the eigenvectors in this ONB ξ_1, \dots, ξ_m and their eigenvalues $\lambda_1, \dots, \lambda_m$
因为对称
- We can rotate the coordinate system to ξ_1, \dots, ξ_m . In the new coordinate system, Σ has the form
称序

$$\Sigma_{[\xi_1, \dots, \xi_m]} = \begin{pmatrix} \lambda_1 & 0 & \cdots & 0 \\ 0 & \lambda_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \lambda_m \end{pmatrix} = \text{diag}(\lambda_1, \dots, \lambda_m)$$

Example

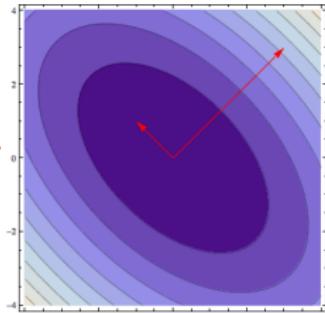
Quadratic form

$$\langle \mathbf{x}, \Sigma \mathbf{x} \rangle \quad \text{with} \quad \Sigma = \begin{pmatrix} 2 & 1 \\ 1 & 2 \end{pmatrix}$$

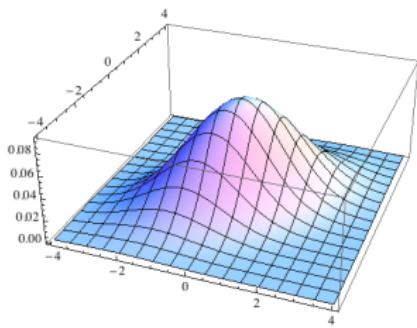
The eigenvectors are $(1, 1)$ and $(-1, 1)$ with eigenvalues 3 and 1.

Gaussian density

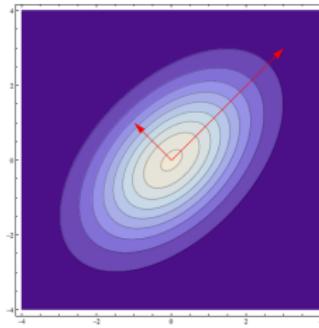
$p(\mathbf{x}; \boldsymbol{\mu}, \Sigma)$ with $\boldsymbol{\mu} = (0, 0)$.



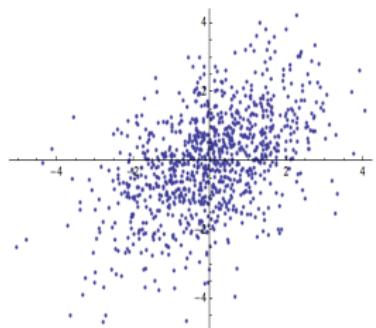
二次型下特征向量
总能指向根据变动最大的方向



Density graph



Density contour



1000 sample points 9 / 30

Interpretation

The ξ_i as random variables

Write e_1, \dots, e_m for the ONB of axis vectors. We can represent each ξ_i as

$$\xi_i = \sum_{j=1}^m \alpha_{ij} e_j$$

新政基为原政基向量的线性组合

Then $O = (\alpha_{ij})$ is the orthogonal transformation matrix between the two bases.

We can represent random vector $X \in \mathbb{R}^m$ sampled from the Gaussian in the eigen-ONB as

$$X_{[\xi_1, \dots, \xi_m]} = (X'_1, \dots, X'_m) \quad \text{with} \quad X'_i = \sum_{j=1}^m \alpha_{ij} X_j$$

Since the X_j are random variables (and the α_{ij} are fixed), each X'_i is a scalar random variable.

Interpretation

Meaning of the random variables ξ_i

For any Gaussian $p(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$, we can

1. shift the origin of the coordinate system into $\boldsymbol{\mu}$
2. rotate the coordinate system to the eigen-ONB of $\boldsymbol{\Sigma}$.

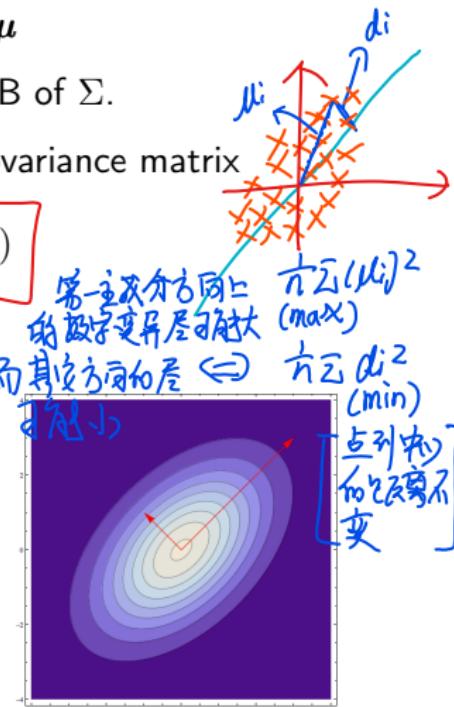
In this new coordinate system, the Gaussian has covariance matrix

$$\boldsymbol{\Sigma}_{[\xi_1, \dots, \xi_m]} = \text{diag}(\lambda_1, \dots, \lambda_m)$$

where λ_i are the eigenvalues of $\boldsymbol{\Sigma}$.

Gaussian in the new coordinates

A Gaussian vector $\mathbf{X}_{[\xi_1, \dots, \xi_m]}$ represented in the new coordinates consists of m independent 1D Gaussian variables X'_i . Each X'_i has mean 0 and variance λ_i .



Principal Component Analysis

$$Z(X) = \mu$$

$$\bar{Z} = \begin{pmatrix} \bar{b}_{11} & \bar{b}_{12} & \dots & \bar{b}_{1p} \\ \vdots & \ddots & \ddots & \vdots \\ \bar{b}_{p1} & \dots & \dots & \bar{b}_{pp} \end{pmatrix}$$

标准化之前

(利用 2)

$$\begin{cases} \text{每个主成分是} \\ \text{原 Base 的} \\ \text{线性组合} \end{cases} \begin{cases} Z_1 = a_1^T X \\ Z_2 = a_2^T X \\ \dots \\ Z_p = a_p^T X \end{cases}$$

- This is the most popular unsupervised procedure ever.
- Invented by Karl Pearson (1901).
- Developed by Harold Hotelling (1933).
- What does it do? It provides a way to visualize high dimensional data, summarizing the most important information. 又： Z 是对称阵 \Rightarrow 又有所谓特征值

$$Q = [a_1, a_2, \dots, a_p] \Rightarrow Q^T \bar{Z} Q = \Lambda = [\lambda_1, \dots, \lambda_p]$$

又的特点： $Z(X) = \bar{Z}(Q^T X) = Q^T \mu$

$$\left\{ \text{Var}(Z) = \text{Var}(Q^T X) = Q^T X Q = \Lambda \right.$$

$$\left. \text{tr}(\Lambda) = \text{tr}(Q^T \bar{Z} Q) = \text{tr}(\bar{Z}) \Rightarrow \sum_{i=1}^p \lambda_i = \sum_{i=1}^p \bar{b}_{ii} \right.$$

$$\left. \sum_i \text{Var}(Z_i) = \sum_i \text{Var}(X_i) \right.$$

$$\text{Var}(Z_i) = a_i^T \bar{Z} a_i$$

$$\text{Cov}(Z_i, Z_j) = a_i^T \bar{Z} a_j$$

目标找到 a_1 使 $\text{Var}(Z_1)$ 最大

$$a_1^T \bar{Z} a_1 = \lambda_1$$

$$a_1^T \bar{Z} a_1 = a_1^T \bar{Z} a_1 = a_1 \lambda_1$$

$$\Rightarrow \bar{Z} a_1 = a_1 \lambda_1$$

特征向量 特征值

$\frac{\lambda_i}{\sum \lambda_i} \Rightarrow$ 叫做贡献率

$$X = QZ \Rightarrow x_j = q_{j1}z_1 + q_{j2}z_2 + \dots + q_{jp}z_p$$

$$\text{cov}(x_j, z_i) = q_{ji} \lambda_i$$

$$P(x_j, z_i) = \frac{\sqrt{\lambda_i}}{\sqrt{b_{ii}}} q_{ji}$$

$$b_{ij} = q_{j1}^2 \lambda_1 + q_{j2}^2 \lambda_2 + \dots + q_{jp}^2 \lambda_p$$

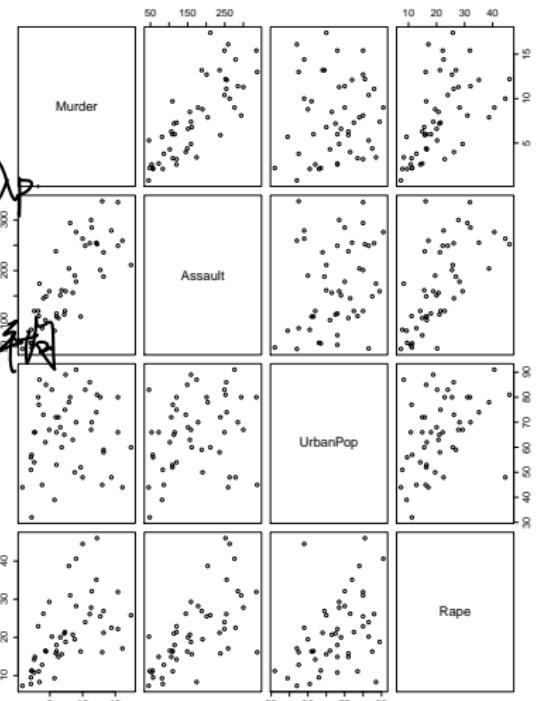
$$q_{j1}^2 + \dots + q_{jp}^2 = 1$$

故 b_{ij} 是 $\lambda_1, \dots, \lambda_p$ 的加权平均数

$$P_j^2 | 1 \dots p = \sum_{i=1}^p P(x_j, z_i)^2$$

$$= \sum_{i=1}^p \frac{\lambda_i q_{ji}^2}{b_{ii}} = 1$$

所以每个 z_i 都对 x_j 进行了一部分的解释，加起来为 1.



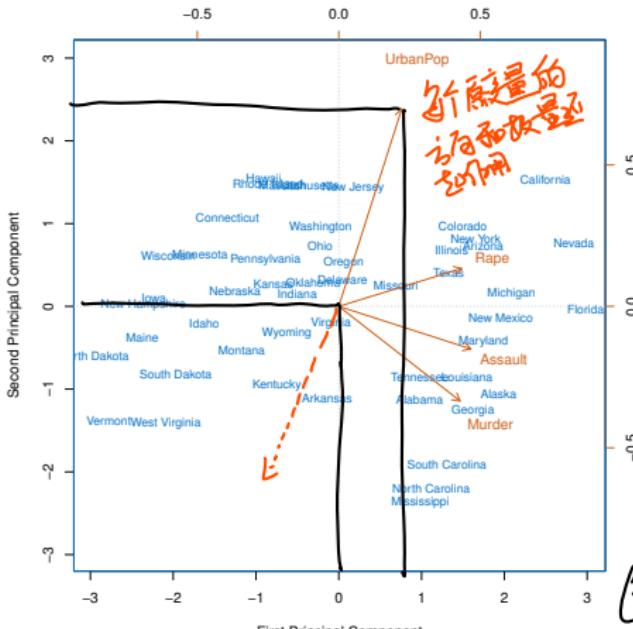
What is PCA good for?

标准化后: 利用 R(相关系数)

$$\begin{cases} E(Z) = 0 & b_{ii} = 1 \\ \sum \lambda_i = p & \text{恒成立} \end{cases}$$

$$P(x_j, z_i) = \sqrt{\lambda_i} q_{ji}$$

What is PCA good for?



ISL Figure 10.1

① 行主成分的特征组合向量成了新变量

$$\text{如: } \text{UrbanPop} = 2.5Z_1 + 0.8Z_2$$

② 行主成分是之前变量的线性组合.

$$Z = Q^T X$$

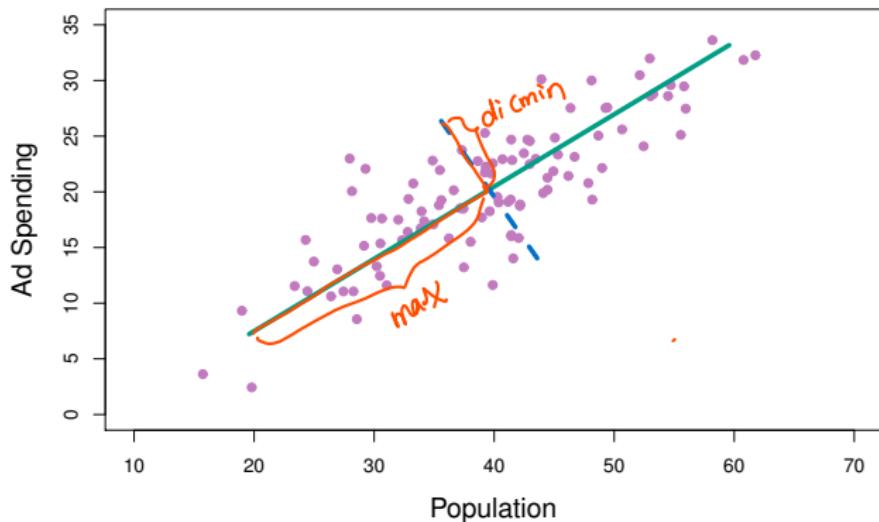
$$X = Q Z$$

$$Q = [q_1, \dots, q_p]$$

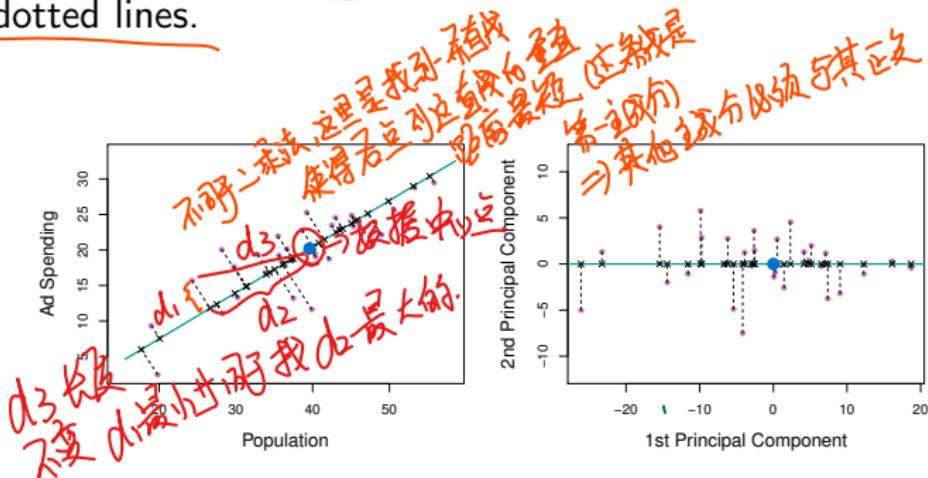
Q 是又怎樣組合為 X

What is the first principal component?

It is the vector which passes the closest to a cloud of samples, in terms of squared Euclidean distance.



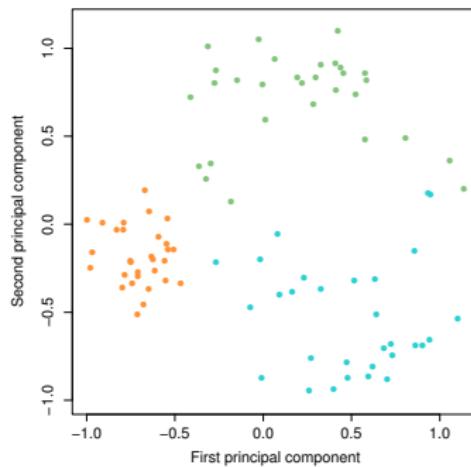
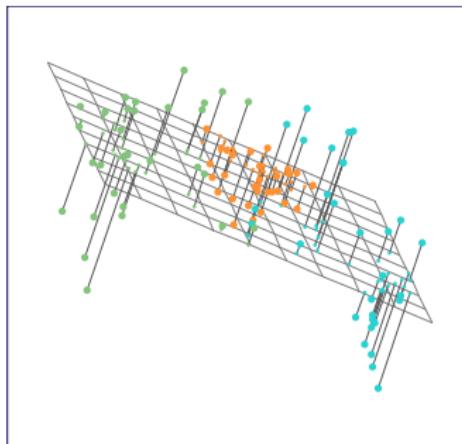
i.e. The green direction minimizes the average squared length of the dotted lines.



ISL Figure 6.15

What does this look like with 3 variables?

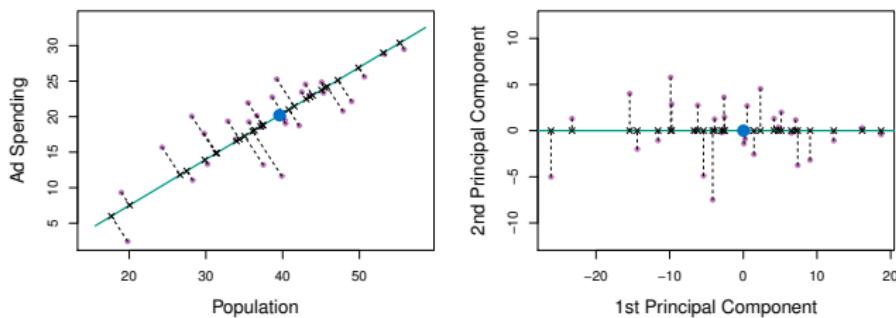
The first two principal components span a plane which is closest to the data.



ISL Figure 10.2

A second interpretation

The projection onto the first principal component is the one with the highest variance.

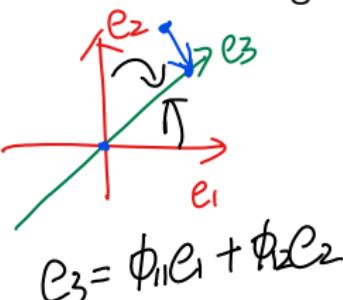


ISL Figure 6.15

How do we say this in math?

Let \mathbf{X} be a data matrix with n samples, and p variables. From each variable, we subtract the mean of the column; i.e. we center the variables.

To find the first principal component $\phi_1 = (\phi_{11}, \dots, \phi_{p1})$, we solve the following optimization



* 改变后的新生元

$$\max_{\phi_{11}, \dots, \phi_{p1}} \left\{ \frac{1}{n} \sum_{i=1}^n \left(\sum_{j=1}^p \phi_{j1} x_{ij} \right)^2 \right\} \quad \max \mathbf{a}_1^T \mathbf{X} \mathbf{a}_1$$
$$\text{subject to } \sum_{j=1}^p \phi_{j1}^2 = 1. \quad \mathbf{a}_1^T \mathbf{a}_1 = 1$$

Projection of the i th sample onto ϕ_1 . Also known as the score z_{i1}

How do we say this in math?

Let \mathbf{X} be a data matrix with n samples, and p variables. From each variable, we subtract the mean of the column; i.e. we **center** the variables.

To find the first principal component $\phi_1 = (\phi_{11}, \dots, \phi_{p1})$, we solve the following optimization

$$\max_{\phi_{11}, \dots, \phi_{p1}} \left\{ \frac{1}{n} \sum_{i=1}^n \left(\sum_{j=1}^p \phi_{j1} x_{ij} \right)^2 \right\}$$

subject to $\sum_{j=1}^p \phi_{j1}^2 = 1$.

Variance of the n samples projected onto ϕ_1 .

How do we say this in math?

To find the second principal component $\phi_2 = (\phi_{12}, \dots, \phi_{p2})$, we solve the following optimization

$$\max_{\phi_{12}, \dots, \phi_{p2}} \left\{ \frac{1}{n} \sum_{i=1}^n \left(\sum_{j=1}^p \phi_{j2} x_{ij} \right)^2 \right\}$$

subject to $\sum_{j=1}^p \phi_{j2}^2 = 1$ and $\sum_{j=1}^p \phi_{j1} \phi_{j2} = 0$.

$\alpha_i^T \alpha_j = 0$

First and second principal components must be orthogonal.

How do we say this in math?

To find the second principal component $\phi_2 = (\phi_{12}, \dots, \phi_{p2})$, we solve the following optimization

$$\max_{\phi_{12}, \dots, \phi_{p2}} \left\{ \frac{1}{n} \sum_{i=1}^n \left(\sum_{j=1}^p \phi_{j2} x_{ij} \right)^2 \right\}$$

subject to $\sum_{j=1}^p \phi_{j2}^2 = 1$ and $\sum_{j=1}^p \phi_{j1} \phi_{j2} = 0$.

First and second principal components must be orthogonal.

Equivalent to saying that the scores (z_{11}, \dots, z_{n1}) and (z_{12}, \dots, z_{n2}) are uncorrelated.

?

Solving the optimization

This optimization is fundamental in linear algebra. It is satisfied by either:

- The singular value decomposition (SVD) of \mathbf{X} :

$$u_i = \langle \mathbf{x}, \phi_i \rangle$$

$$u_1, \dots, u_n$$

$$\mathbf{X} = \mathbf{U}\Sigma\Phi^T$$

$$\sum u_i^2 = (\mathbf{x}\phi_i)^T(\mathbf{x}\phi_i)$$

where the *i*th column of Φ is the *i*th principal component ϕ_i ,
 and the *i*th column of $\mathbf{U}\Sigma$ is the *i*th vector of scores
 (z_{1i}, \dots, z_{ni}) .

$$= \sum_{i=1}^n u_i^2 \lambda_i$$

- The eigendecomposition of $\mathbf{X}^T \mathbf{X}$:

$$\text{set } a_1 = 1$$

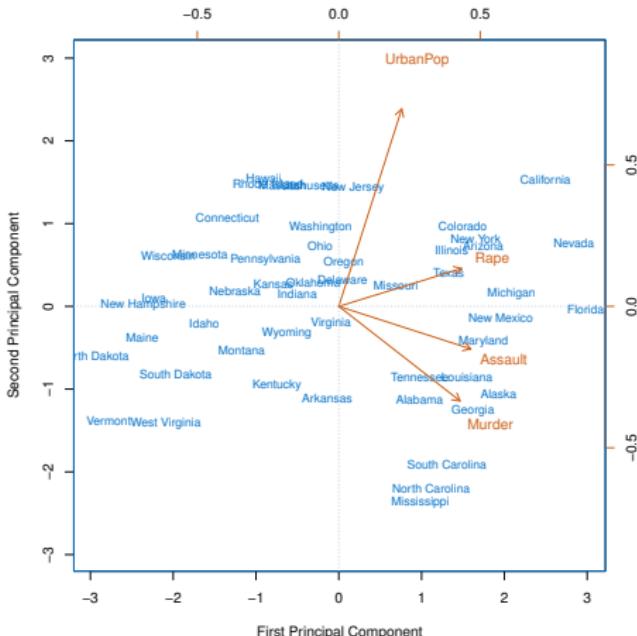
$$a_2 = a_3 = \dots = a_p = 0$$

$$\boxed{\begin{aligned} \bar{\Sigma} &= \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{pmatrix} \\ \mathbf{X}^T \mathbf{X} &= \mathbf{V} \bar{\Sigma} \mathbf{V}^T \end{aligned}}$$

$$\mathbf{X}\bar{\Phi} = \begin{pmatrix} \langle \mathbf{x}, \bar{\phi}_1 \rangle \\ \vdots \\ \langle \mathbf{x}, \bar{\phi}_n \rangle \end{pmatrix} \quad \mathbf{X}^T \mathbf{X} = \Phi \Sigma^2 \Phi^T$$

after centered
 $\mathbf{X}^T \mathbf{X}$ is correlation matrix

PCA in practice: The biplot



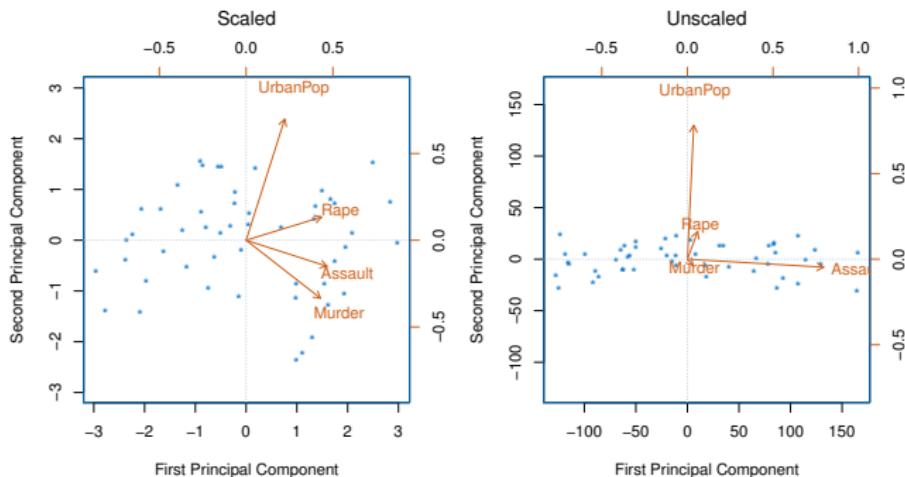
ISL Figure 10.1

Scaling the variables

Most of the time, we don't care about the absolute numerical value of a variable. We care about the value relative to the spread observed in the sample.

Before PCA, in addition to centering each variable, we also multiply it times a constant to make its variance equal to 1. 标准化

Example: scaled vs. unscaled PCA



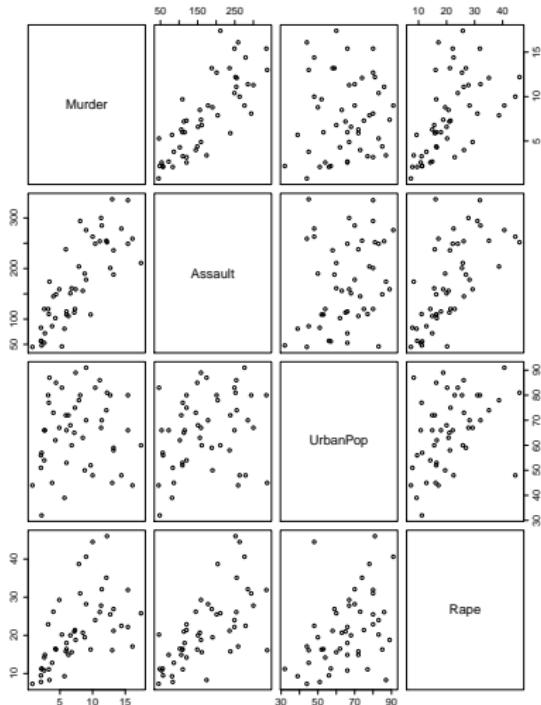
ISL Figure 10.3

Scaling the variables

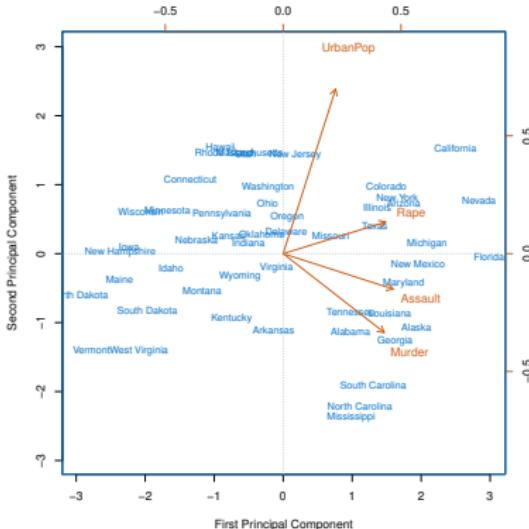
In special cases, we have variables measured in the same unit; e.g. gene expression levels for different genes.

Therefore, we care about the absolute value of the variables and we can perform PCA without scaling.

How many principal components are enough?



How many principal components are enough?



We said 2 principal components capture most of the relevant information. But how can we tell?

The proportion of variance explained

We can think of the top **principal components** as directions in space in which the data vary the most.

The i th **score vector** (z_{1i}, \dots, z_{ni}) can be interpreted as a *new* variable. The variance of this variable decreases as we take i from 1 to p . However, the total variance of the score vectors is the same as the total variance of the original variables:

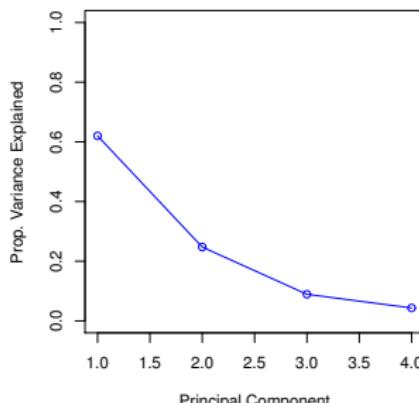
$$\sum_{i=1}^p \frac{1}{n} \sum_{j=1}^n z_{ji}^2 = \sum_{k=1}^p \text{Var}(x_k).$$

We can quantify how much of the variance is captured by the first m principal components/score variables.

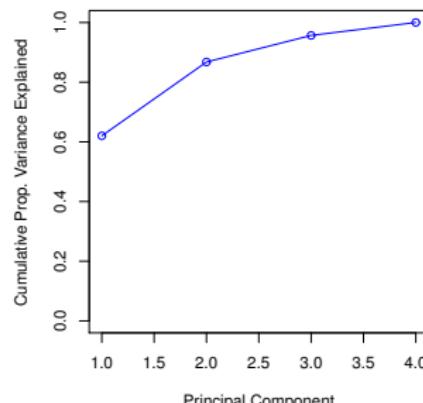
The proportion of variance explained

The variance of the m th score variable is:

$$\frac{1}{n} \sum_{i=1}^n z_{im}^2 = \frac{1}{n} \sum_{i=1}^n \left(\sum_{j=1}^p \phi_{jm} x_{ij} \right)^2 = \frac{1}{n} \Sigma_{mm}^2.$$



Scree plot



Generalizations of PCA

PCA works under a Euclidean geometry in the space of variables. Often, the natural geometry is different:

- ▶ We expect some variables to be “closer” to each other than to other variables.
- ▶ Some correlations between variables would be more surprising than others.

Examples:

- ▶ Variables are pixel values, samples are different images of the brain. We expect neighboring pixels to have stronger correlations.
- ▶ Variables are rainfall measurements at different regions. We expect neighboring regions to have higher correlations.

Generalizations of PCA

There are ways to include this knowledge in a PCA. See:

1. Susan Holmes. *Multivariate Analysis, the French way*. (2006).
2. Omar de la Cruz and Susan Holmes. *An introduction to the duality diagram*. (2011).
3. Stéphane Dray and Thibaut Jombart. *Revisiting Guerry's data: Introducing spatial constraints in multivariate analysis*. (2011).
4. Genevera Allen, Logan Grosenick, and Jonathan Taylor. *A Generalized Least Squares Matrix Decomposition*. (2011).

Thanks to Sergio Bacallado and Peter Orbanz
for sharing the slides.