



西南财经大学  
SOUTHWESTERN UNIVERSITY OF FINANCE AND ECONOMICS



经世济民 孜孜以求

机器学习基础

主成分分析



- 基本思想
- 总体主成分分析
- 样本主成分分析
- 应用：基因组分析



- 基本思想
- 总体主成分分析
- 样本主成分分析
- 应用：基因组分析



## 例子

- 采集了房屋的价格和面积，可以看出两者完全正相关，有一列其实是多余的

	房价（百万元）	面积（百平米）
<i>a</i>	10	10
<i>b</i>	2	2
<i>c</i>	1	1
<i>d</i>	7	7
<i>e</i>	3	3



## 例子

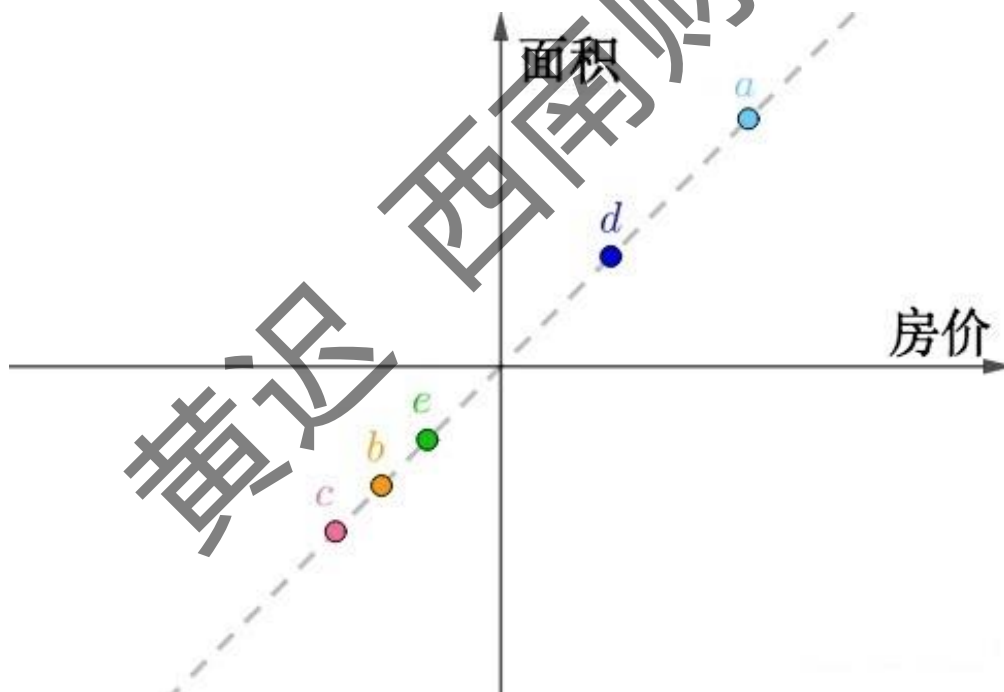
- 求出房价样本、面积样本的均值，分别对房价样本、面积样本进行“中心化”后得到：

	房价（百万元）	面积（百平米）
$a$	5.4	5.4
$b$	-2.6	-2.6
$c$	-3.6	-3.6
$d$	2.4	2.4
$e$	-1.6	-1.6



## 例子

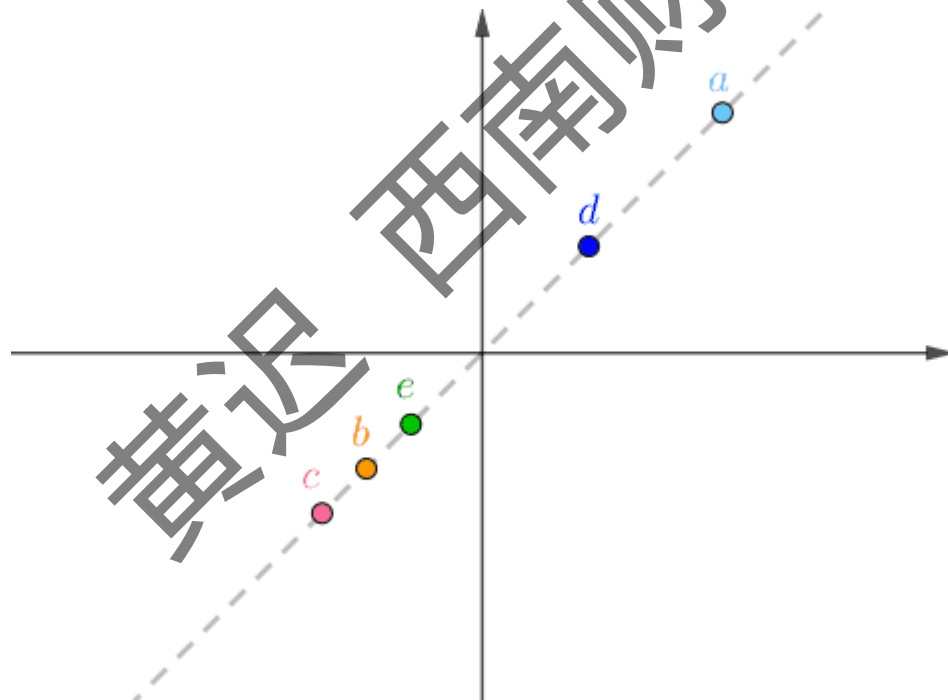
- 把样本数据画在坐标系上





## 例子

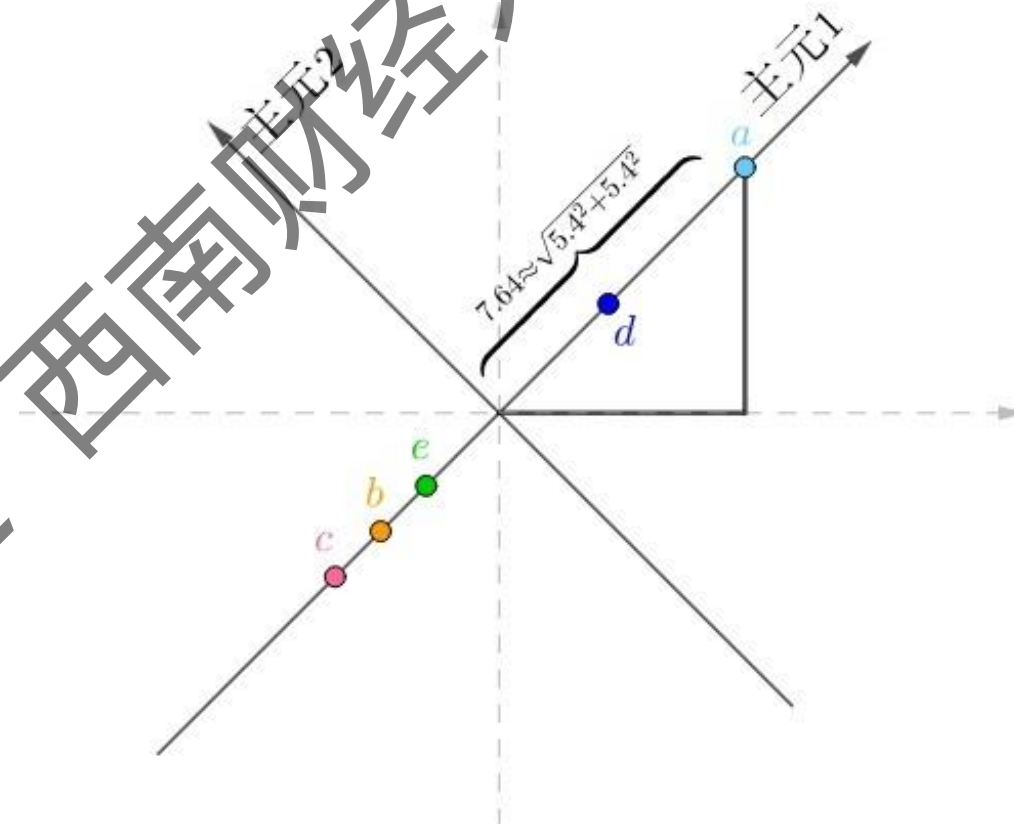
- 旋转坐标系，让横坐标和这条直线重合





## 例子

- 旋转后的坐标系，横纵坐标不再代表“房价”、“面积”了，而是两者的混合（线性组合），即“主元1”、“主元2”







## 例子

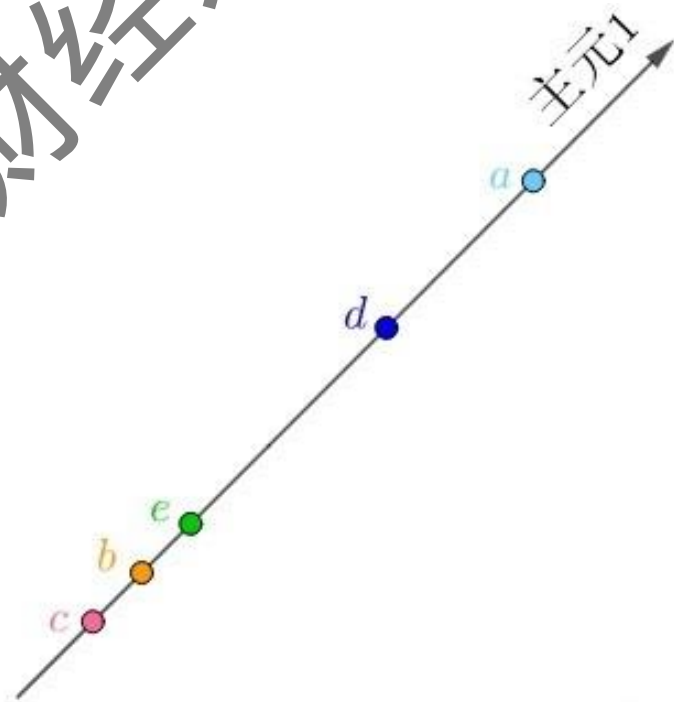
- 重新计算样本在新坐标系下的坐标有

	主元 1	主元 2
$a$	7.64	0
$b$	-3.68	0
$c$	-5.09	0
$d$	3.39	0
$e$	-2.26	0



## 例子

- 所有样本在主元2上的坐标都是0，所以只需要保留主元1，这样就把数据降为一维，而且没有丢失任何信息





## 更加真实的例子

- 样本的房价和面积不再完全相同

	房价 (百万元)	面积 (百平米)
$a$	10	9
$b$	2	3
$c$	1	2
$d$	7	6.5
$e$	3	2.5



## 更加真实的例子

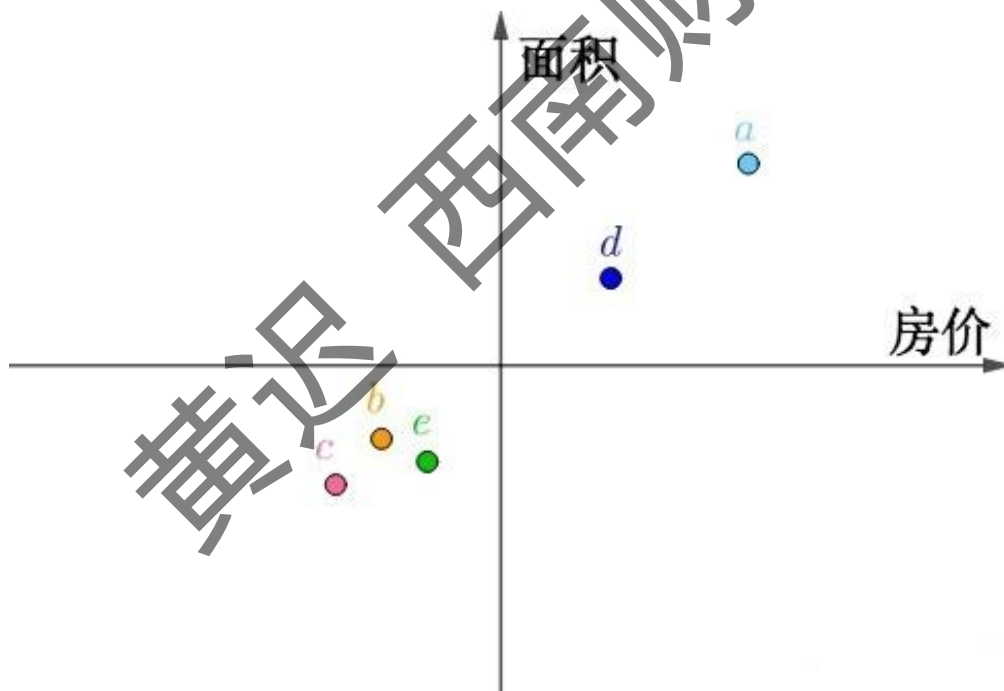
- 中心化

	房价 (百万元)	面积 (百平米)
<i>a</i>	5.4	4.4
<i>b</i>	-2.6	-1.6
<i>c</i>	-3.6	-2.6
<i>d</i>	2.4	1.9
<i>e</i>	-1.6	-2.1



## 更加真实的例子

- 把样本数据画在坐标系上





## 更加真实的例子

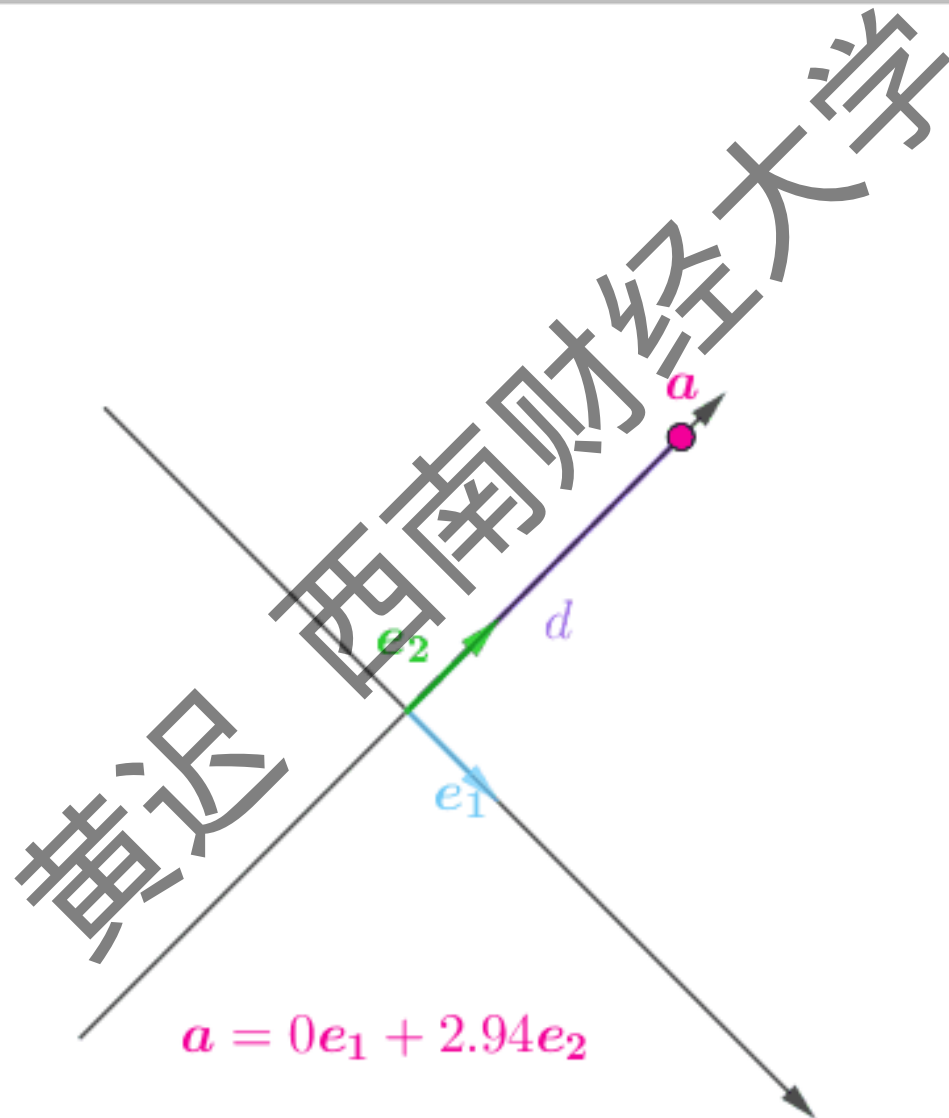
- 如何降维？

- 假设  $e_1, e_2$  是二维平面中任意一组标准正交基 为什么要求正交？
- 点  $a$  在  $e_1, e_2$  下的坐标为  $(x, y)$ , 即  $a = xe_1 + ye_2$
- 当  $e_1, e_2$  变化时, 坐标会发生改变, 但点  $a$  到坐标原点的距离不会发生改变, 即  $d = x^2 + y^2$  是常数
- 若  $x$  变大,  $y$  会变小; 反之亦然



## 如何降维？

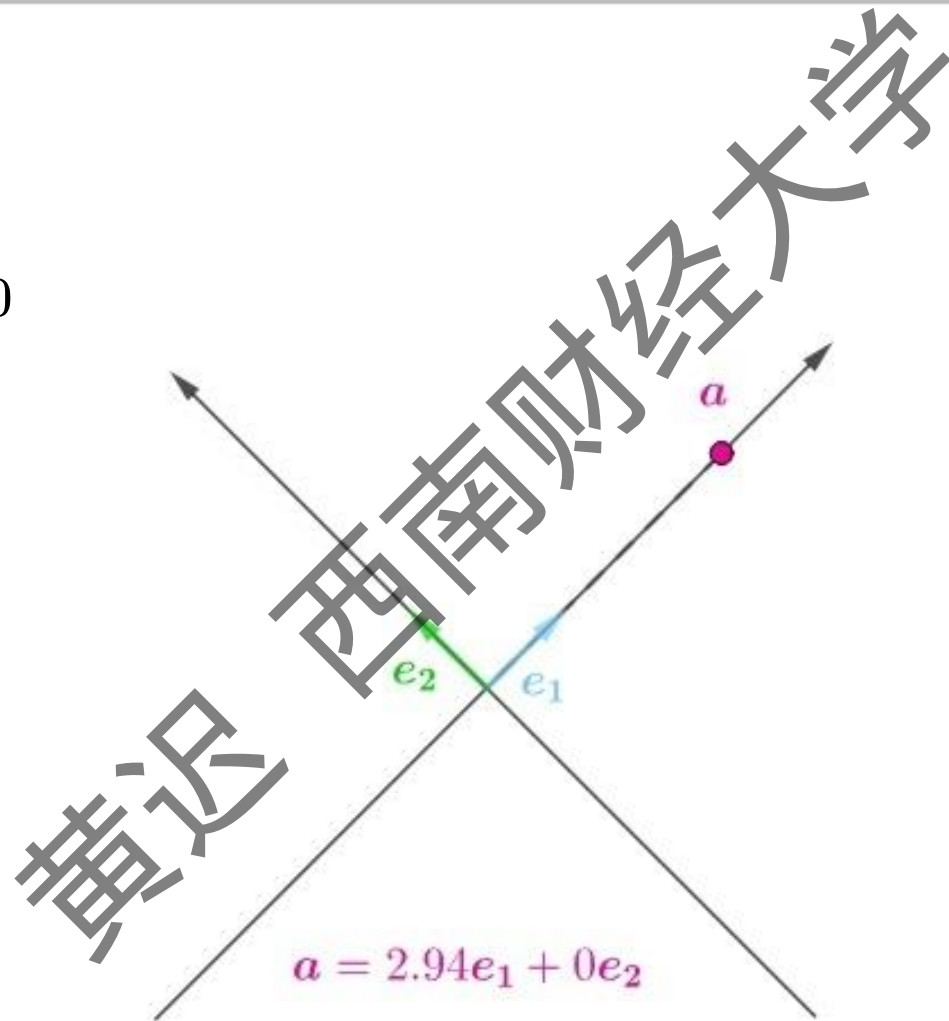
- $x$  和  $y$  的变化





## 如何降维？

- 极端情况下， $y=0$

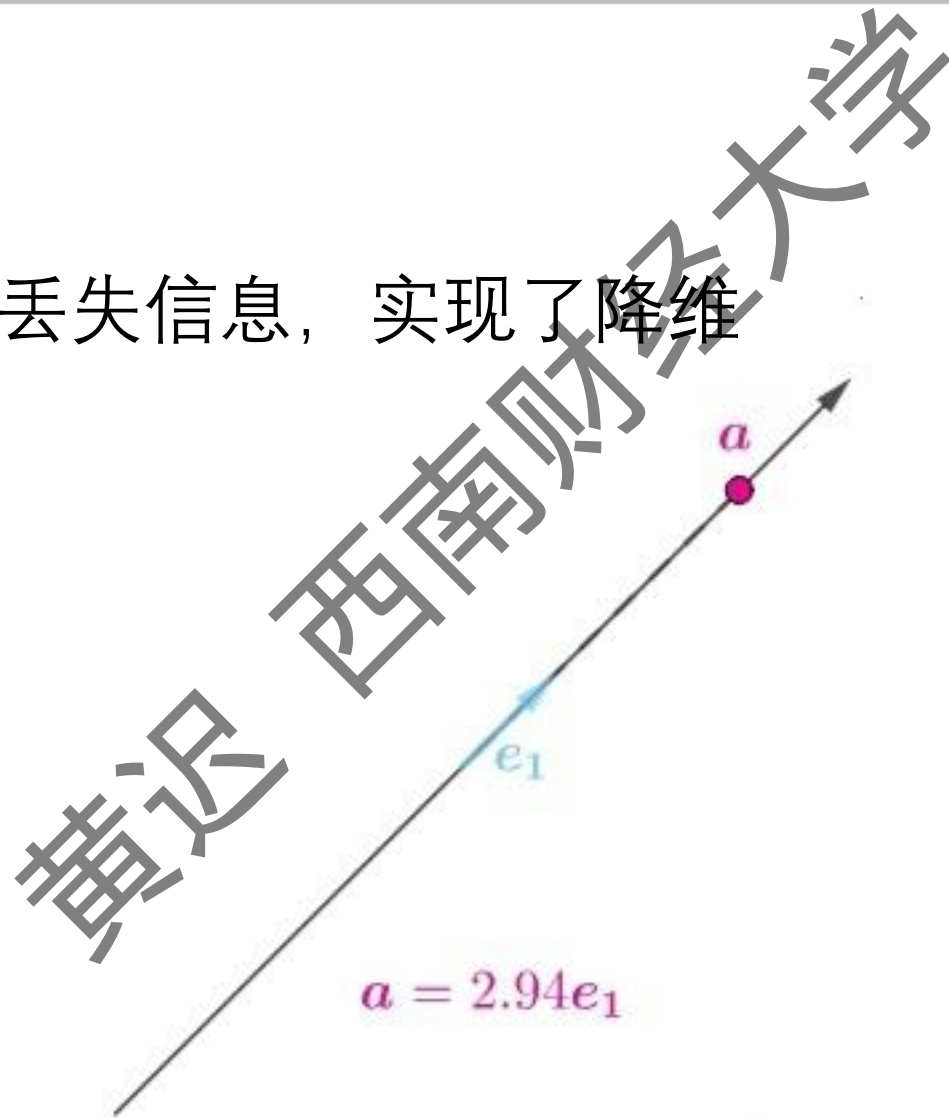






## 如何降维？

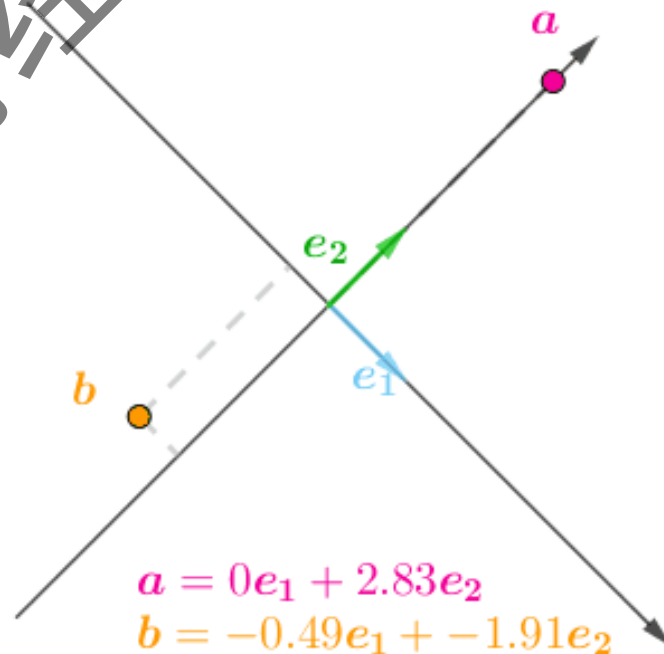
- 因此去掉  $e_2$  不会丢失信息，实现了降维





## 如何降维？

- 如果是两个点  $a = (x_1, y_1)$ ,  $b = (x_2, y_2)$ , 则很难使  $y_1, y_2$  同时为 0
- 为了降维, 应该尽量增大  $x_1, x_2$ , 相应的  $y_1, y_2$  就会减小



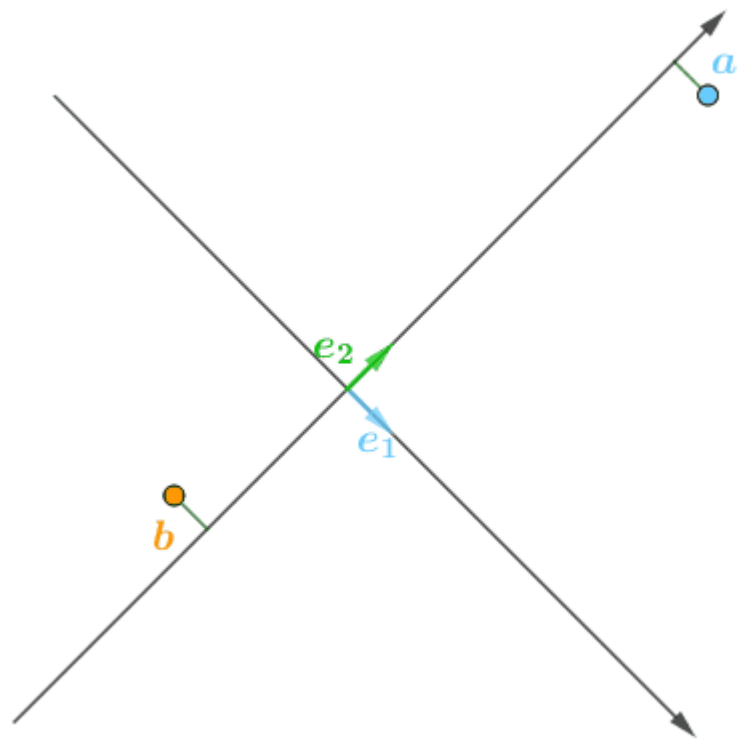


## 如何降维？

- 令样本在不同的标准正交基下坐标为

	$\mathbf{X} = (1, 0)$	$\mathbf{Y} = (0, 1)$	$\mathbf{e}_1$	$\mathbf{e}_2$
$\mathbf{a}$	$a_x$	$a_y$	$x_1$	$y_1$
$\mathbf{b}$	$b_x$	$b_y$	$x_2$	$y_2$

- 随着坐标系的变化,  $x_1, x_2$  也在不断变化





## 如何降维？

- 为了实现降维，应该要让  $x_1, x_2$  尽量大，即实现  $\max x_1^2 + x_2^2$
- 由于  $x_1 = \mathbf{a} \cdot \mathbf{e}_1$ ,  $x_2 = \mathbf{b} \cdot \mathbf{e}_1$ , 有

$$\begin{aligned} x_1^2 + x_2^2 &= \mathbf{e}_1^\top \mathbf{a}^\top \mathbf{a} \mathbf{e}_1 + \mathbf{e}_1^\top \mathbf{b}^\top \mathbf{b} \mathbf{e}_1 \\ &= \mathbf{e}_1^\top \begin{pmatrix} a_x^2 + b_x^2 & a_x a_y + b_x b_y \\ a_x a_y + b_x b_y & a_y^2 + b_y^2 \end{pmatrix} \mathbf{e}_1 \triangleq \mathbf{e}_1^\top P \mathbf{e}_1 \end{aligned}$$

- 矩阵  $P$  是二次型矩阵，存在正交分解  $P = U \Sigma U^\top = U \begin{bmatrix} \sigma_1 & 0 \\ 0 & \sigma_2 \end{bmatrix} U^\top$   
其中  $\sigma_i$  是特征值（奇异值），且  $\sigma_1 \geq \sigma_2$



## 如何降维？

- 将分解代回，有  $x_1^2 + x_2^2 = \mathbf{e}_1^\top U \Sigma U^\top \mathbf{e}_1$
- 令  $\mathbf{n} = \begin{pmatrix} n_1 \\ n_2 \end{pmatrix} = U^\top \mathbf{e}_1$ ，则  $\mathbf{n}$  也是单位向量，即  $n_1^2 + n_2^2 = 1$
- 继续回代有  $x_1^2 + x_2^2 = \mathbf{n}^\top \Sigma \mathbf{n} = \mathbf{n}^\top \begin{pmatrix} \sigma_1 & 0 \\ 0 & \sigma_2 \end{pmatrix} \mathbf{n} = \sigma_1 n_1^2 + \sigma_2 n_2^2$
- 最终优化问题变为

$$\max x_1^2 + x_2^2$$



$$\max \sigma_1 n_1^2 + \sigma_2 n_2^2$$

$$s.t. \quad n_1^2 + n_2^2 = 1$$

$$\sigma_1 \geq \sigma_2$$



## 如何降维？

- 当  $n_1=1, n_2=0$  时，该优化问题达到最优，此时有

$$\mathbf{n} = \begin{pmatrix} 1 \\ 0 \end{pmatrix} = U^\top \mathbf{e}_1 \Rightarrow \mathbf{e}_1 = U \begin{pmatrix} 1 \\ 0 \end{pmatrix} = \text{Col}_1(U)$$

即  $\mathbf{e}_1$  是最大特征值  $\sigma_1$  的特征向量

- 由于  $\mathbf{e}_1, \mathbf{e}_2$  相互正交，因此  $\mathbf{e}_2 = \text{Col}_2(U)$

即  $\mathbf{e}_2$  是第二特征值  $\sigma_2$  的特征向量



## 统计分析

- 回顾样本数据:

	$\mathbf{X} = (1, 0)$	$\mathbf{Y} = (0, 1)$	$\mathbf{e}_1$	$\mathbf{e}_2$
$\mathbf{a}$	$a_x$	$a_y$	$x_1$	$y_1$
$\mathbf{b}$	$b_x$	$b_y$	$x_2$	$y_2$

- 令  $\mathbf{x} = \begin{pmatrix} a_x \\ b_x \end{pmatrix}$ ,  $\mathbf{y} = \begin{pmatrix} a_y \\ b_y \end{pmatrix}$ , 分别表示两个特征在样本上的分布, 有

$$P = \begin{pmatrix} a_x^2 + b_x^2 & a_x a_y + b_x b_y \\ a_x a_y + b_x b_y & a_y^2 + b_y^2 \end{pmatrix} = \begin{pmatrix} \mathbf{x} \cdot \mathbf{x} & \mathbf{x} \cdot \mathbf{y} \\ \mathbf{x} \cdot \mathbf{y} & \mathbf{y} \cdot \mathbf{y} \end{pmatrix}$$



## 统计分析

- 由于数据已经中心化，所以有

- 样本方差  $Var(\mathbf{x}) = \mathbf{x} \cdot \mathbf{x}$ ,  $Var(\mathbf{y}) = \mathbf{y} \cdot \mathbf{y}$

- 样本协方差  $Cov(\mathbf{x}, \mathbf{y}) = \mathbf{x} \cdot \mathbf{y}$

- 矩阵  $P$  就是样本协方差矩阵，即

$$P = \begin{pmatrix} \mathbf{x} \cdot \mathbf{x} & \mathbf{x} \cdot \mathbf{y} \\ \mathbf{x} \cdot \mathbf{y} & \mathbf{y} \cdot \mathbf{y} \end{pmatrix} = \frac{1}{n-1} \begin{pmatrix} Var(\mathbf{x}) & Cov(\mathbf{x}, \mathbf{y}) \\ Cov(\mathbf{x}, \mathbf{y}) & Var(\mathbf{y}) \end{pmatrix}$$

- 主元是样本协方差矩阵的单位特征向量





## 回到例子

	房价	面积
$a$	5.4	4.4
$b$	-2.6	-1.6
$c$	-3.6	-2.6
$d$	2.4	1.9
$e$	-1.6	-2.1

$$\mathbf{x} = \begin{pmatrix} 5.4 \\ -2.6 \\ -3.6 \\ 2.4 \\ -1.6 \end{pmatrix} \quad \mathbf{y} = \begin{pmatrix} 4.4 \\ -1.6 \\ -2.6 \\ 1.9 \\ -2.1 \end{pmatrix}$$

$$P = \frac{1}{n-1} \begin{pmatrix} \text{Var}(\mathbf{x}) & \text{Cov}(\mathbf{x}, \mathbf{y}) \\ \text{Cov}(\mathbf{x}, \mathbf{y}) & \text{Var}(\mathbf{y}) \end{pmatrix} \\ = \frac{1}{4} \begin{pmatrix} 57.2 & 45.2 \\ 45.2 & 36.7 \end{pmatrix}$$

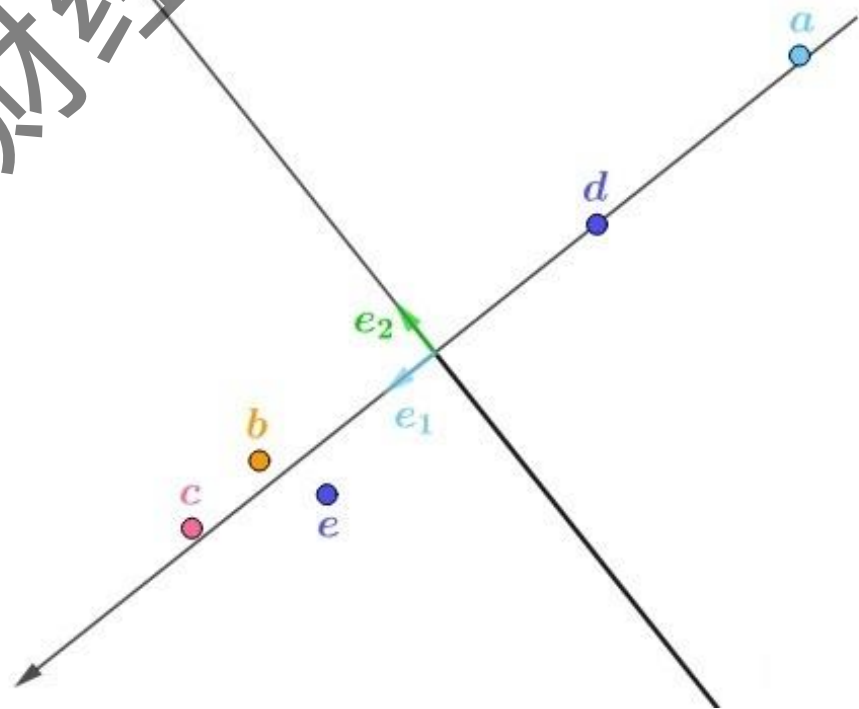


## 回到例子

- 求得样本协方差矩阵的特征值，即主成分分别为

$$e_1 = \begin{pmatrix} -0.78 \\ -0.62 \end{pmatrix}, e_2 = \begin{pmatrix} -0.62 \\ 0.78 \end{pmatrix}$$

- 在主成分  $e_1, e_2$  下，样本如图所示：





## 回到例子

- 样本在主成分下的坐标为

	主成分 1	主成分 2
$a$	-6.94	0.084
$b$	3.02	0.364
$c$	4.42	0.204
$d$	-3.05	-0.006
$e$	2.55	-0.646

- 样本在主成分2中的数值很小，损失信息也很小，降维效果好



- 基本思想
- 总体主成分分析
- 样本主成分分析
- 应用：基因组分析



## 定义

- 假设  $\mathbf{x} = (x_1, x_2, \dots, x_m)^\top$  是  $m$  维随机变量 (原始坐标)
  - 均值  $\boldsymbol{\mu} = \mathbb{E}(\mathbf{x})$ ; 协方差矩阵  $\Sigma = \text{Cov}(\mathbf{x}, \mathbf{x}) = \mathbb{E}[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^\top]$
- 定义  $m$  维随机变量  $\mathbf{y} = (y_1, y_2, \dots, y_m)^\top$  (旋转后的坐标, 即主成分) :

$$y_i = \alpha_i^\top \mathbf{x} = \alpha_{1i}x_1 + \alpha_{2i}x_2 + \dots + \alpha_{mi}x_m$$

- $\mathbb{E}(y_i) = \mathbb{E}(\alpha_i^\top \mathbf{x}) = \alpha_i^\top \boldsymbol{\mu}$ ,  $\text{Var}(y_i) = \alpha_i^\top \Sigma \alpha_i$

$$\text{Cov}(y_i, y_j) = \alpha_i^\top \Sigma \alpha_j, \quad i = 1, 2, \dots, m, \quad j = 1, 2, \dots, m$$



## 总体主成分

- 对  $y_i, i=1, 2, \dots, m$ , 满足以下条件
  - 系数向量  $\alpha_i^\top$  是单位向量,  $\alpha_i^\top \alpha_i = 1, i=1, 2, \dots, m$
  - $y_i$  与  $y_j$  互不相关,  $Cov(y_i, y_j) = 0, i \neq j$
  - 第一主成分  $y_1$ : 是  $\mathbf{x}$  所有线性组合中方差最大的, 即

$$y = \arg \max_y Var(y) = \arg \max_{\alpha_1} \alpha_1^\top \Sigma \alpha_1$$
$$s.t. \quad \alpha_1^\top \alpha_1 = 1$$



## 总体主成分

- 第二主成分  $y_2$ : 是与  $y_1$  不相关的  $x$  所有线性组合中方差最大的, 即

$$y = \arg \max_y \text{Var}(y) = \arg \max_{\alpha_2} \alpha_2^\top \Sigma \alpha_2$$

$$s.t. \quad \alpha_1^\top \Sigma \alpha_2 = 0, \quad \alpha_2^\top \alpha_2 = 1$$

- 第  $i$  主成分  $y_i$ : 是与  $y_1, y_2, \dots, y_{i-1}$  都不相关的  $x$  所有线性组合中方差最大的, 即

$$y = \arg \max_y \text{Var}(y) = \arg \max_{\alpha_i} \alpha_i^\top \Sigma \alpha_i$$

$$s.t. \quad \alpha_1^\top \Sigma \alpha_i = 0, \dots, \alpha_{i-1}^\top \Sigma \alpha_i = 0, \quad \alpha_i^\top \alpha_i = 1$$



## 主要性质

- 设  $\Sigma$  的特征值分别是  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_m \geq 0$ , 特征值对应的单位特征向量分别是  $\alpha_1, \alpha_2, \dots, \alpha_m$ , 则  $\mathbf{x}$  第  $k$  主成分是

$$y_k = \alpha_k^\top \mathbf{x}$$

$\mathbf{x}$  的第  $k$  主成分的方差是

$$\text{Var}(y_k) = \alpha_k^\top \Sigma \alpha_k = \lambda_k$$





## 主要性质

- 令  $A = [\alpha_1, \alpha_2, \dots, \alpha_m]$ , 则  $A$  为正交矩阵, 且  $\mathbf{y} = (y_1, y_2, \dots, y_m)^\top = A^\top \mathbf{x}$
- $Cov(\mathbf{y}) = \Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_m)$
- 总体主成分  $\mathbf{y}$  的方差之和等于随机变量  $\mathbf{x}$  的方差之和  $\sum_{i=1}^m \lambda_i = \sum_{i=1}^m \sigma_{ii}$
- 因子负荷量  $\rho(y_k, x_i)$  表示主成分  $y_k$  与变量  $x_i$  之间的相关关系

$$\rho(y_k, x_i) = \frac{Cov(y_k, x_i)}{\sqrt{Var(y_k)Var(x_i)}} = \frac{\sqrt{\lambda_k} \alpha_{ik}}{\sqrt{\sigma_{ii}}}, \quad k, i = 1, 2, \dots, m$$

- 一般的, 小于0.4是低, 大于0.6是高



## 主要性质

- 主成分  $y_k$  与  $m$  个变量的因子负荷量满足

$$\sum_{k=1}^m \sigma_{ii} \rho^2(y_k, x_i) = \sum_{k=1}^m \lambda_{ik} \alpha_{ik}^2 = \lambda_k \alpha_k^\top \alpha_k = \lambda_k$$

- $m$  个主成分与变量  $x_i$  的因子负荷量满足

$$\sum_{k=1}^m \rho^2(y_k, x_i) = \sum_{k=1}^m \frac{\lambda_{ik} \alpha_{ik}^2}{\sigma_{ii}} = 1$$

- $\sigma_{ii} = [\Sigma]_{ii} = [A \Lambda A^\top]_{ii} = \sum_{k=1}^m \lambda_k \alpha_{ik}^2$



## 主成分的个数

- **定理16.2:** 对任意  $q$ ,  $1 \leq q \leq m$ , 考虑正交线性变化  $\mathbf{y} = B^\top \mathbf{x}$ , 其中  $B^\top$  是  $q \times m$  矩阵, 令  $\mathbf{y}$  的协方差矩阵为  $\Sigma_{\mathbf{y}} = B^\top \Sigma B$ , 则  $\text{trace}(\Sigma_{\mathbf{y}})$  在  $B = A_q$  时取得最大值, 其中矩阵  $A_q$  由正交矩阵  $A$  的前  $q$  列组成
- 定理结论表明选择前  $q$  个主成分能保留原有变量的最多信息
  - 方差越大保留的信息越多
- **定理16.3:**  $\text{trace}(\Sigma_{\mathbf{y}})$  在  $B = A_p$  时取得最小值, 其中矩阵  $A_p$  由正交矩阵  $A$  的后  $p$  列组成



## 方差贡献率

- 主成分  $y_k$  的方差贡献率是  $y_k$  的方差与所有方差之和的比

$$\eta_k = \frac{\lambda_k}{\sum_{i=1}^m \lambda_i}$$

- 前  $k$  个主成分的累计方差贡献率是前  $k$  个方差之和与所有方差之和的比

$$\sum_{i=1}^k \eta_i = \frac{\sum_{i=1}^k \lambda_i}{\sum_{i=1}^m \lambda_i}$$



## 规范化变量的总体主成分

- 对变量  $x_i$  进行规范化  $x_i^* = \frac{x_i - \mathbb{E}(x_i)}{\sqrt{\text{Var}(x_i)}}$  (均值为0, 方差为1)
- 规范化随机变量的协方差矩阵就是相关矩阵
- 规范化变量主成分的协方差矩阵  $\Lambda^* = \text{diag}(\lambda_1^*, \lambda_2^*, \dots, \lambda_m^*)$ 
  - $\lambda_1^* \geq \lambda_2^* \geq \dots \geq \lambda_m^* \geq 0$ , 且  $\lambda_i^*$  是相关矩阵  $R$  的特征值
- 协方差矩阵的特征值之和为  $m$

$$\sum_{k=1}^m \lambda_k^* = \sum_{i=1}^m \sigma_{ii}^* = m$$



## 规范化变量的总体主成分

- 规范化变量与主成分的因子负荷量是

$$\rho(y_k^*, x_i^*) = \frac{Cov(y_k^*, x_i^*)}{\sqrt{Var(y_k^*)Var(x_i^*)}} = \frac{\sqrt{\lambda_k^*} e_{ik}^*}{\sqrt{\sigma_{ii}^*}} = \sqrt{\lambda_k^*} e_{ik}^*, k, i = 1, 2, \dots, m$$

- $e_k^* = (e_{1k}^*, e_{2k}^*, \dots, e_{mk}^*)^\top$  是相关矩阵对应于特征值  $\lambda_k^*$  的特征向量
- 规范化主成分与  $m$  个规范化变量的因子负荷量满足

$$\sum_{i=1}^m \rho(y_k^*, x_i^*)^2 = \sum_{i=1}^m \lambda_k^* e_{ik}^{*2} = \lambda_k^* \|e_k^*\|^2 = \lambda_k^*, k = 1, 2, \dots, m$$

- 规范化变量与规范化主成分的因子负荷量满足  $\sum_{k=1}^m \rho^2(y_k^*, x_i^*) = 1$



- 基本思想
- 总体主成分分析
- 样本主成分分析
- 应用：基因组分析



## 样本定义

- 对随机向量  $\mathbf{x} = (x_1, x_2, \dots, x_m)^\top$  进行  $n$  次独立观测
  - 得到观测样本  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ , 其中  $\mathbf{x}_j = (x_{1j}, x_{2j}, \dots, x_{mj})^\top$
- 样本矩阵  $\mathbf{X} = [\mathbf{x}_1 \ \mathbf{x}_2 \ \dots \ \mathbf{x}_n]$
- 样本矩阵向量  $\bar{\mathbf{x}} = (\bar{x}_1, \bar{x}_2, \dots, \bar{x}_m)^\top = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i$
- 样本协方差矩阵  $S = [s_{ij}]_{m \times m}$ , 其中  $s_{ij} = \frac{1}{n-1} \sum_{k=1}^n (x_{ik} - \bar{x}_i)(x_{jk} - \bar{x}_j)$
- 样本相关矩阵  $R = [r_{ij}]_{m \times m}$ , 其中  $r_{ij} = \frac{s_{ij}}{\sqrt{s_{ii}s_{jj}}}$





## 样本主成分的定义

- 令  $m$  维向量之间的线性变换为  $\mathbf{y} = A^\top \mathbf{x}$ , 其中  $A = [a_1 \ a_2 \ \cdots \ a_m] = [a_{ij}]_{m \times m}$
- 样本  $\mathbf{x}_i$  在变换下为  $\mathbf{y}_i = A^\top \mathbf{x}_i$ , 变换后的均值为

$$\bar{\mathbf{y}} = (\bar{y}_1, \bar{y}_2, \cdots, \bar{y}_m)^\top = \frac{1}{n} \sum_{i=1}^n \mathbf{y}_i = \frac{1}{n} \sum_{i=1}^n A^\top \mathbf{x}_i$$

- 其中  $\bar{y}_j = \frac{1}{n} \sum_{i=1}^n a_j^\top \mathbf{x}_i = a_j^\top \bar{\mathbf{x}}$ ,  $\bar{\mathbf{x}}$  为样本均值
- 变换后样本的第  $i$  个分量  $Y_i = (y_{i1}, y_{i2}, \cdots, y_{in})$ ,  $i = 1, 2, \cdots, m$  的方差为

$$\text{Var}(Y_i) = \frac{1}{n-1} \sum_{j=1}^n (y_{ij} - \bar{y}_i)^2 = \frac{1}{n-1} \sum_{j=1}^n (a_i^\top \mathbf{x}_j - a_i^\top \bar{\mathbf{x}})^2 = a_i^\top S a_i$$



## 样本主成分的定义

- 变换后两个分量的协方差为

$$\begin{aligned} \text{Var}(Y_i, Y_k) &= \frac{1}{n-1} \sum_{j=1}^n (y_{ij} - \bar{y}_i) (y_{kj} - \bar{y}_k) \\ &= \frac{1}{n-1} \sum_{j=1}^n (a_i^\top \mathbf{x}_j - a_i^\top \bar{\mathbf{x}}) (a_k^\top \mathbf{x}_j - a_k^\top \bar{\mathbf{x}}) \\ &= a_i^\top \left[ \frac{1}{n-1} (\mathbf{x}_j - \bar{\mathbf{x}}) (\mathbf{x}_j - \bar{\mathbf{x}})^\top \right] a_k \\ &= a_i^\top S a_k \end{aligned}$$



## 样本主成分

- 给定样本矩阵  $\mathbf{x}$
- 样本第一主成分  $y_1 = \mathbf{a}_1^\top \mathbf{x}$  是在  $\mathbf{a}_1^\top \mathbf{a}_1 = 1$  下使得  $\text{Var}(Y_1) = \mathbf{a}_1^\top \mathbf{S} \mathbf{a}_1$  最大的线性组合
- 样本第二主成分  $y_2 = \mathbf{a}_2^\top \mathbf{x}$  是在  $\mathbf{a}_2^\top \mathbf{a}_2 = 1, \text{Var}(Y_2, Y_1) = \mathbf{a}_2^\top \mathbf{S} \mathbf{a}_1 = 0$  下使得  $\text{Var}(Y_2)$  最大的线性组合
- 样本第  $i$  主成分  $y_i = \mathbf{a}_i^\top \mathbf{x}$  是在  $\mathbf{a}_i^\top \mathbf{a}_i = 1, \text{Var}(Y_i, Y_j) = \mathbf{a}_i^\top \mathbf{S} \mathbf{a}_j = 0, j = 1, 2, \dots, i-1$  下使得  $\text{Var}(Y_i)$  最大的线性组合



## 样本主成分与总体主成分

- 总体主成分：在数据总体上进行的主成分分析
- 样本主成分：在有限样本上进行的主成分分析
  - 用样本协方差矩阵  $S$  替代总体协方差矩阵  $\Sigma$
  - 用样本均值  $\bar{x}$  替代总体均值  $\mu$
- 但结论相同，主成分都是协方差矩阵的特征值所对应的特征向量所构成的线性组合



## 规范后的样本主成分

- 样本数据规范化:  $x_{ij}^* = \frac{x_{ij} - \bar{x}_i}{\sqrt{s_{ii}}}$ 
  - 其中  $\bar{x}_i = \frac{1}{n} \sum_{j=1}^n x_{ij}$ ,  $s_{ii} = \frac{1}{n-1} \sum_{j=1}^n (x_{ij} - \bar{x}_i)^2$ ,  $i=1, 2, \dots, m$
  - 规范化后, 样本第  $i$  个分量  $X_i^* = (x_{i1}^*, x_{i2}^*, \dots, x_{in}^*)$  的均值为0, 方差为1
- 令规范化后的样本矩阵为  $X^* = [x_{ij}^*]_{m \times m}$ , 则规范化后的样本协方差矩阵就是样本相关矩阵  $R = \frac{1}{n-1} X^* (X^*)^\top$



## 相关矩阵的特征值分解

- 样本主成分分析是通过求相关矩阵的特征值实现，步骤如下：
  1. 将样本进行规范化，得到规范化后的样本矩阵  $X \in \mathbb{R}^{m \times n}$
  2. 计算相关矩阵  $R = [r_{ij}]_{m \times m} = \frac{1}{n-1} X X^\top$ ，其中  $r_{ij} = \frac{1}{n-1} \sum_{l=1}^n x_{il} x_{lj}$
  3. 求相关矩阵  $R$  的特征值  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_m$  和对应的特征向量
  4. 求方差贡献率达到预定值的主成分个数  $k$
  5. 求出前  $k$  个样本主成分  $y_i = a_i^\top \mathbf{x}$ ,  $i = 1, 2, \dots, k$



## 样本矩阵的奇异值分解

- 传统的主成分分析通过样本的相关矩阵的特征值分解进行，现在常用的方法是通过样本矩阵的奇异值分解进行，步骤如下

1. 构造矩阵  $X' = \frac{1}{\sqrt{n-1}} X^\top \in \mathbb{R}^{n \times m}$
2. 对  $X'$  进行  $k$  阶截断奇异值分解，得到  $X' = U \Sigma V^\top$
3. 求前  $k$  个主成分的矩阵  $Y = V^\top X$





## 奇异值分解的优点

- 计算样本的相关矩阵  $R = \frac{1}{n-1} XX^T$  及其特征值需要消耗大量的计算量，特别是当样本量很大时
- 注意到PCA仅仅使用了SVD的右奇异矩阵，没有使用的左奇异矩阵
- 有些SVD的实现算法可以不用求出  $R$  也能计算出右奇异矩阵  $V$ ，从而极大的节约了计算成本



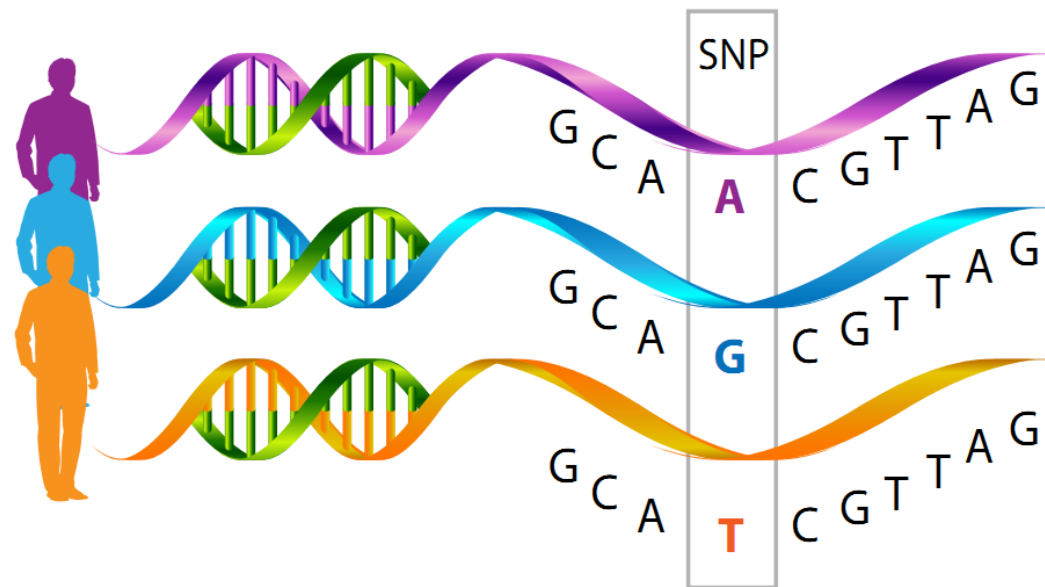


- 基本思想
- 总体主成分分析
- 样本主成分分析
- 应用：基因组分析



## 人类基因组的SNP数据

- 1064个志愿者，并把他们的SNP数据数字化，也就是把每个位置上可能出现的10种碱基对用数字来代表
- SNP（单核苷酸多态性）是人类可遗传的变异中最常见的一种，估计其总数可达300万个甚至更多广泛用于基于群体的基因识别等方面的研究



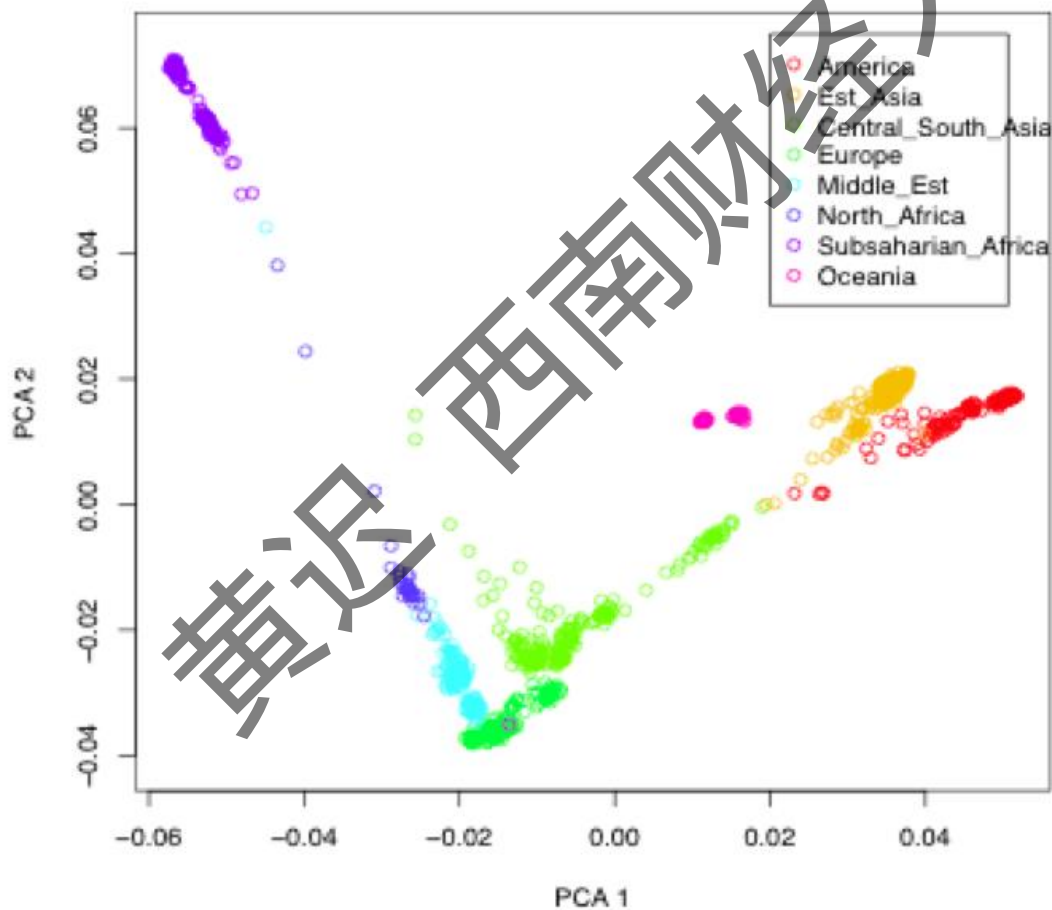


## 基因数据

	SNP 1	SNP 2	...	SNP m
志愿者1	0	1	...	4
志愿者2	1	0	...	2
志愿者3	0	2	...	1
...	...	...	...	...
志愿者n	0	9	...	0



## 主成分分析的结果





## 主成分分析的结果

