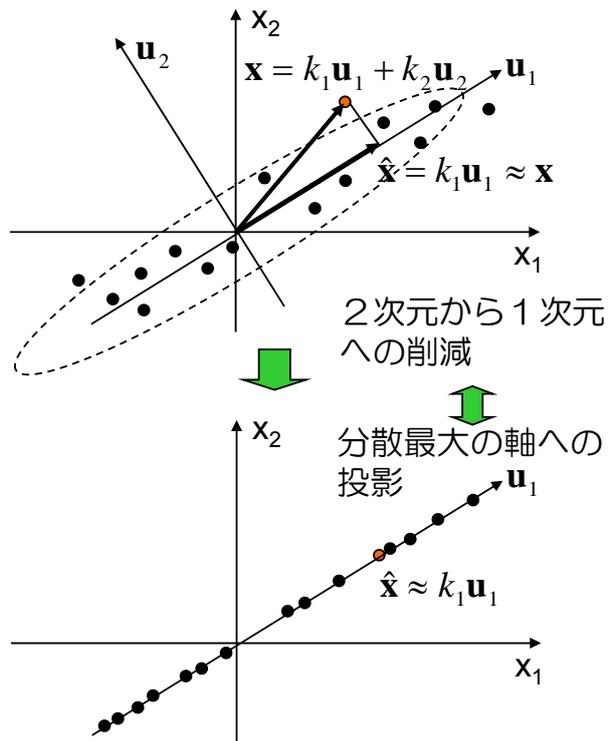


「互いに相関のある多種類の変数を、互いに無相関な少数個の変数に要約する。」



少ない次元数で解析，圧縮などを行う。

平均ベクトルが0の2変数分布の場合



$u_1$  : サンプルの分散が最大の方向  
(単位ベクトル)

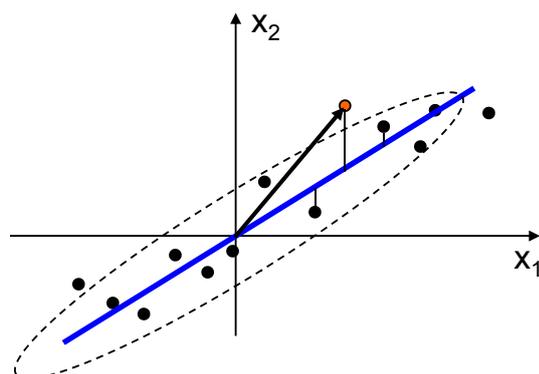
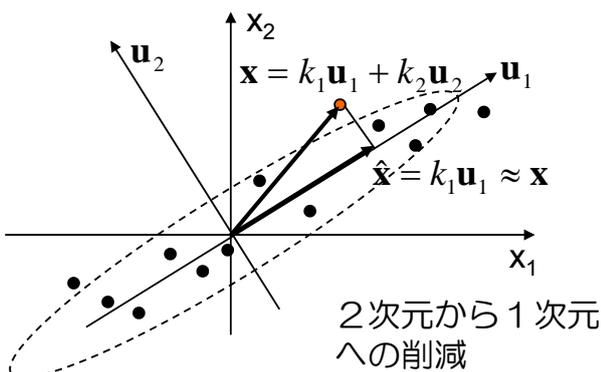
$u_2$  : 2番目に分散が大きい方向  
(単位ベクトル)  
(2変数の場合は，分散が最小の方向)

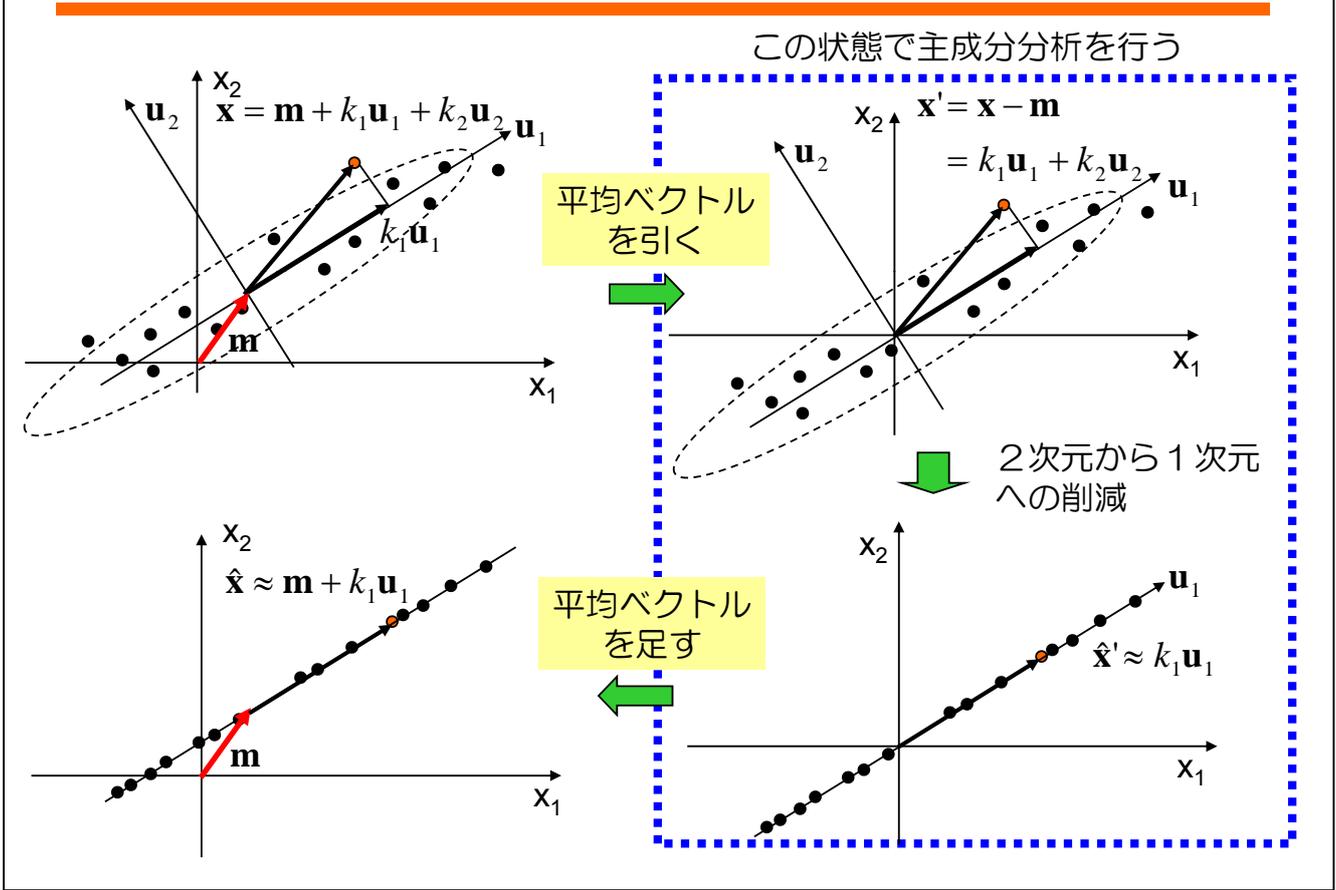
## 回帰分析とはどこが違う？

主成分分析：  
複数個の変数から互いに無相関な変数  
(もとの変数の線形結合で表現される)  
に集約すること。  
(回帰分析のような主従関係はない)

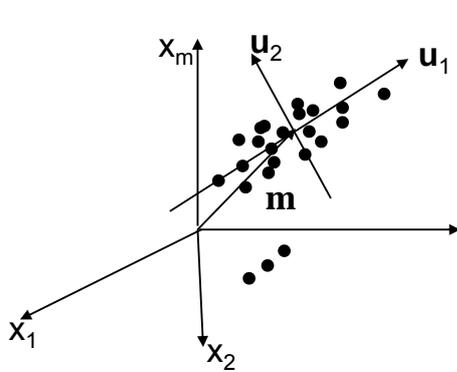
回帰分析：  
1個以上の説明変数を用いて1個(以上)  
の目的変数を近似的な関数で表現  
すること。

2変数に対する1次回帰分析との違い

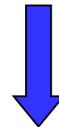




m次元から2次元への要約



$$\mathbf{x} = \mathbf{m} + k_1 \mathbf{u}_1 + k_2 \mathbf{u}_2 + \dots + k_m \mathbf{u}_m$$



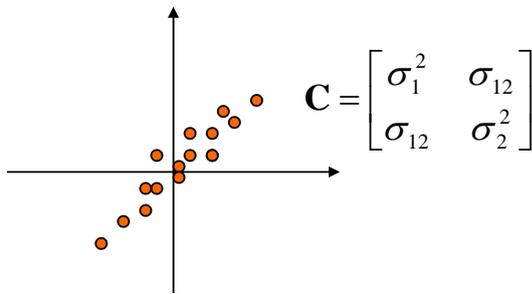
$$\hat{\mathbf{x}} = \mathbf{m} + k_1 \mathbf{u}_1 + k_2 \mathbf{u}_2$$

共分散行列 各変数とも、平均を0にしてから  
相関を計算して得られる行列

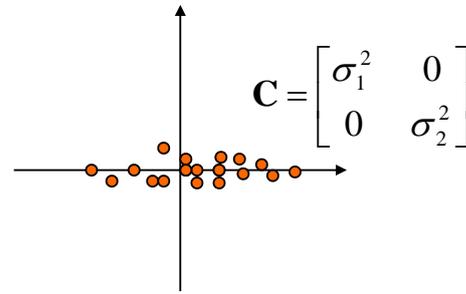
$$\mathbf{C} = \begin{bmatrix} \frac{1}{n} \sum_{i=1}^n (x_{i,1} - m_1)^2 & \frac{1}{n} \sum_{i=1}^n (x_{i,1} - m_1)(x_{i,2} - m_2) \\ \frac{1}{n} \sum_{i=1}^n (x_{i,1} - m_1)(x_{i,2} - m_2) & \frac{1}{n} \sum_{i=1}^n (x_{i,2} - m_2)^2 \end{bmatrix}$$

$$= \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \mathbf{m})(\mathbf{x}_i - \mathbf{m})^T = \begin{bmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{12} & \sigma_2^2 \end{bmatrix}$$

(例1) 2変数間に相関がある場合



(例2) 2変数間に相関がない場合



主成分の方向 $\mathbf{u}_1, \mathbf{u}_2$ を求める

方針

2変数間で相関が0になる方向、  
すなわち共分散行列が対角行列になる  
方向を求めればよい。

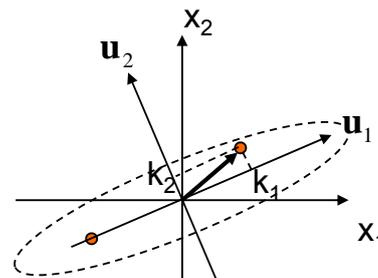
上記の条件を満足する、正規直交基底  
ベクトルを $\mathbf{u}_1, \mathbf{u}_2$ とする。

この2つの基底ベクトルを用いた座標  
変換は以下のように表される。

$$\begin{bmatrix} k_1 \\ k_2 \end{bmatrix} = \begin{bmatrix} u_{11} & u_{12} \\ u_{21} & u_{22} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \quad \dots(1)$$

またはベクトル表現で

$$\mathbf{k} = \mathbf{U}^T \mathbf{x} \quad \text{ただし} \quad \mathbf{U} = [\mathbf{u}_1 \quad \mathbf{u}_2]$$



$\mathbf{u}_1, \mathbf{u}_2$ : 正規直交基底ベクトル

$\mathbf{U}$ で変換した後の共分散行列 $\mathbf{C}_k$ は

$$\mathbf{C}_k = \langle \mathbf{k} \mathbf{k}^T \rangle = \langle \mathbf{U}^T \mathbf{x} (\mathbf{U}^T \mathbf{x})^T \rangle$$

$$= \mathbf{U}^T \langle \mathbf{x} \mathbf{x}^T \rangle \mathbf{U} = \mathbf{U}^T \mathbf{C} \mathbf{U} \quad \dots(2)$$

両辺に $\mathbf{U}$ を左からかけて  
左右入れ替えると

$$\mathbf{U} \mathbf{U}^T \mathbf{C} \mathbf{U} = \mathbf{U} \mathbf{C}_k \quad \dots(3)$$

## 主成分の方向を求める (つづき)

$$\mathbf{U}\mathbf{U}^T\mathbf{C}\mathbf{U} = \mathbf{U}\mathbf{C}_k \quad \dots(3)$$

$\mathbf{U}$ の正規直交性より

$$\mathbf{U}^T = \mathbf{U}^{-1} \quad \dots(4)$$

よって (3) 式は

$$\mathbf{C}\mathbf{U} = \mathbf{U}\mathbf{C}_k$$

共分散の対角化が $\mathbf{U}$ によってなされるとすると

$$\mathbf{C}_k = \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{bmatrix} \quad \dots(5)$$

成分に分けて表すと

$$\mathbf{C}[\mathbf{u}_1 \quad \mathbf{u}_2] = [\mathbf{u}_1 \quad \mathbf{u}_2] \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{bmatrix}$$

または

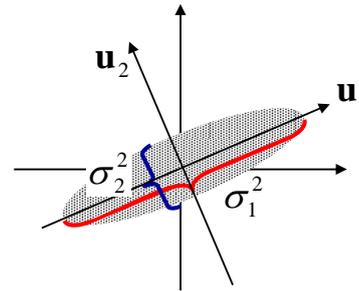
$$\mathbf{C}\mathbf{u}_1 = \sigma_1^2\mathbf{u}_1, \mathbf{C}\mathbf{u}_2 = \sigma_2^2\mathbf{u}_2$$

これは、もとのデータの共分散行列  $\mathbf{C}$  に対する固有値問題に他ならない。

すなわち、 $\mathbf{u}_1, \mathbf{u}_2$  は行列  $\mathbf{C}$  の固有ベクトル、 $\sigma_1^2, \sigma_2^2$  は固有値として求められる。

$\mathbf{C}$  は実対称行列  $\Rightarrow \sigma_i^2 \geq 0$  (for all  $i$ )

固有値  $\sigma_i^2$  は対角化されたデータの各変数の分散を与える。



## 主成分の方向を求める—一般の高次元データ—

2次元での議論をそのまま $m$ 次元へ拡張すればよい

サンプルデータ：

$$\begin{aligned} \mathbf{X} &= [\mathbf{x}_1 \quad \dots \quad \mathbf{x}_n] \\ &= \begin{bmatrix} x_{1,1} & \dots & x_{n,1} \\ \vdots & \dots & \vdots \\ x_{1,m} & \dots & x_{n,m} \end{bmatrix} \end{aligned}$$

サンプルデータの共分散行列：

$$\mathbf{C} = \frac{1}{n} \mathbf{X}\mathbf{X}^T$$

対角化する基底ベクトルのセット：

$$\begin{aligned} \mathbf{U} &= [\mathbf{u}_1 \quad \dots \quad \mathbf{u}_m] \\ &= \begin{bmatrix} u_{1,1} & \dots & u_{m,1} \\ \vdots & \dots & \vdots \\ u_{1,m} & \dots & u_{m,m} \end{bmatrix} \end{aligned}$$

座標変換後の共分散行列：

$$\mathbf{C}_k = \begin{bmatrix} \sigma_1^2 & & 0 \\ & \ddots & \\ 0 & & \sigma_m^2 \end{bmatrix}$$

固有値問題の表現：

$$\mathbf{C}\mathbf{U} = \mathbf{U}\mathbf{C}_k$$

または

$$\mathbf{C}\mathbf{u}_i = \sigma_i^2\mathbf{u}_i, \text{ for } i = 1, \dots, m$$

## 主成分の方向を求めるー計算例ー

共分散行列が以下の式で表されるサンプルの集合を考える。

$$\mathbf{C} = \begin{bmatrix} 1 & \alpha \\ \alpha & 1 \end{bmatrix}$$

ただし平均ベクトルは0とする。

問題1： $\alpha$ が以下の3種類の場合について、サンプルの分布を模式的に描きなさい。

- (i)  $\alpha = 0$
- (ii)  $\alpha = 0.5$
- (iii)  $\alpha = 1.0$

問題2： $0 < \alpha < 1$ の場合について、主成分を計算で求めなさい。

ヒント

- ①固有値問題  $\mathbf{C}\mathbf{u} = \lambda\mathbf{u} \Leftrightarrow (\mathbf{C} - \lambda\mathbf{I})\mathbf{u} = \mathbf{0} \dots (1)$  を解くことと同値である。
- ②正規性の条件  $\|\mathbf{u}\| = \sqrt{u_1^2 + u_2^2} = 1$  を用いる。

## 低次元主成分による近似と誤差

もとの変数と主成分ベクトルの係数との間には以下の関係がある。

$$\mathbf{k} = \mathbf{U}^T \mathbf{x}$$

また、 $\mathbf{U}$ の正規直交性より

$$\mathbf{U}\mathbf{k} = \mathbf{U}\mathbf{U}^T \mathbf{x} \Rightarrow \mathbf{x} = \mathbf{U}\mathbf{k}$$

または

$$\mathbf{x} = \sum_{i=1}^m k_i \mathbf{u}_i$$

いま、 $m=2$ とし、1次元主成分による近似とそれによる誤差を考える。

ある $j$ 番目のサンプルについて、もとデータおよび低次元（1次元）による近似表現は以下のように書ける。

$$\mathbf{x}^{(j)} = k_1^{(j)} \mathbf{u}_1 + k_2^{(j)} \mathbf{u}_2$$

$$\hat{\mathbf{x}}^{(j)} = k_1^{(j)} \mathbf{u}_1$$

誤差ベクトルは

$$\begin{aligned} \mathbf{e}^{(j)} &= \mathbf{x}^{(j)} - \hat{\mathbf{x}}^{(j)} \\ &= (k_1^{(j)} \mathbf{u}_1 + k_2^{(j)} \mathbf{u}_2) - k_1^{(j)} \mathbf{u}_1 \\ &= k_2^{(j)} \mathbf{u}_2 \end{aligned}$$

誤差ベクトルの大きさ（ノルムの2乗）は

$$|\mathbf{e}^{(j)}|^2 = |k_2^{(j)} \mathbf{u}_2|^2 = k_2^{(j)2}$$

サンプル全体での誤差の平均は

$$E = \langle \mathbf{e} \rangle = \frac{1}{n} \sum_{j=1}^n |\mathbf{e}^{(j)}|^2 = \frac{1}{n} \sum_{j=1}^n k_2^{(j)2} = \sigma_2^2$$

すなわち誤差は、用いなかった第2主成分の残差（分散）に等しい。

## 低次元主成分による近似と誤差（つづき）

より一般に、 $m$ 次元データに対する $r$ 次元主成分による近似とそれによる誤差を考える。

もとのデータおよび近似データを

$$\mathbf{x} = \sum_{i=1}^m k_i \mathbf{u}_i, \quad \hat{\mathbf{x}} = \sum_{i=1}^r k_i \mathbf{u}_i \quad r < m$$

と書く。誤差ベクトルは

$$\begin{aligned} \mathbf{e}^{(j)} &= \mathbf{x}^{(j)} - \hat{\mathbf{x}}^{(j)} \\ &= \sum_{i=1}^m k_i^{(j)} \mathbf{u}_i - \sum_{i=1}^r k_i^{(j)} \mathbf{u}_i \\ &= \sum_{i=r+1}^m k_i^{(j)} \mathbf{u}_i \end{aligned}$$

であり、その2ノルムは

$$\begin{aligned} |\mathbf{e}^{(j)}|^2 &= k_{r+1}^{(j)2} + k_{r+2}^{(j)2} + \cdots + k_m^{(j)2} \\ &= \sum_{i=r+1}^m k_i^{(j)2} \end{aligned}$$

このとき、サンプル全体での誤差の平均は

$$\begin{aligned} E(r) &\equiv \langle |\mathbf{e}|^2 \rangle \\ &= \frac{1}{n} \sum_{j=1}^n |\mathbf{e}^{(j)}|^2 \\ &= \frac{1}{n} \sum_{j=1}^n \sum_{i=r+1}^m k_i^{(j)2} \\ &= \sum_{i=r+1}^m \frac{1}{n} \sum_{j=1}^n k_i^{(j)2} \\ &= \sum_{i=r+1}^m \sigma_i^2 \end{aligned}$$

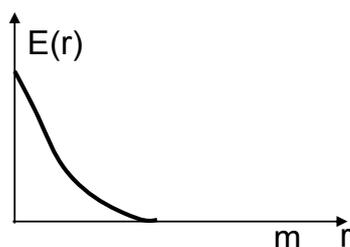
すなわち、誤差は、用いなかった主成分の残差（分散）の和に等しい。

## 近似による誤差と累積寄与率

主成分を、分散の大きい順に番号付けしたものは、誤差

$$E(r) = \sum_{i=r+1}^m \langle k_i^2 \rangle = \sum_{i=r+1}^m \sigma_i^2$$

は、 $r$ に関して単調減少関数となる



逆に、はじめの $r$ 個の成分でどのくらい正確に、もとの分布を表せるかの尺度として、以下に示す累積寄与率がある。

$$\text{累積寄与率}(r) = \frac{\sum_{i=1}^r \sigma_i^2}{\sum_{i=1}^m \sigma_i^2}$$

