

# 多変量データの分析

## 機械学習の観点で言うと

教師あり学習: 学習データに正解を与えた状態で学習

教師なし学習: 学習データに正解を与えない状態で学習

教師なしデータ



次元圧縮

**主成分分析**, 特異値分解

クラスタリング

**K-means**, 階層クラスタリング

教師ありデータ

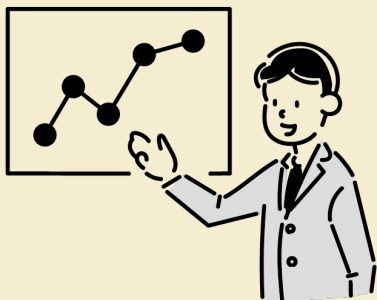


回帰

線形回帰, 決定木, ニューラルネットワーク

分類

SVM, 判別分析, 単純ベイズ



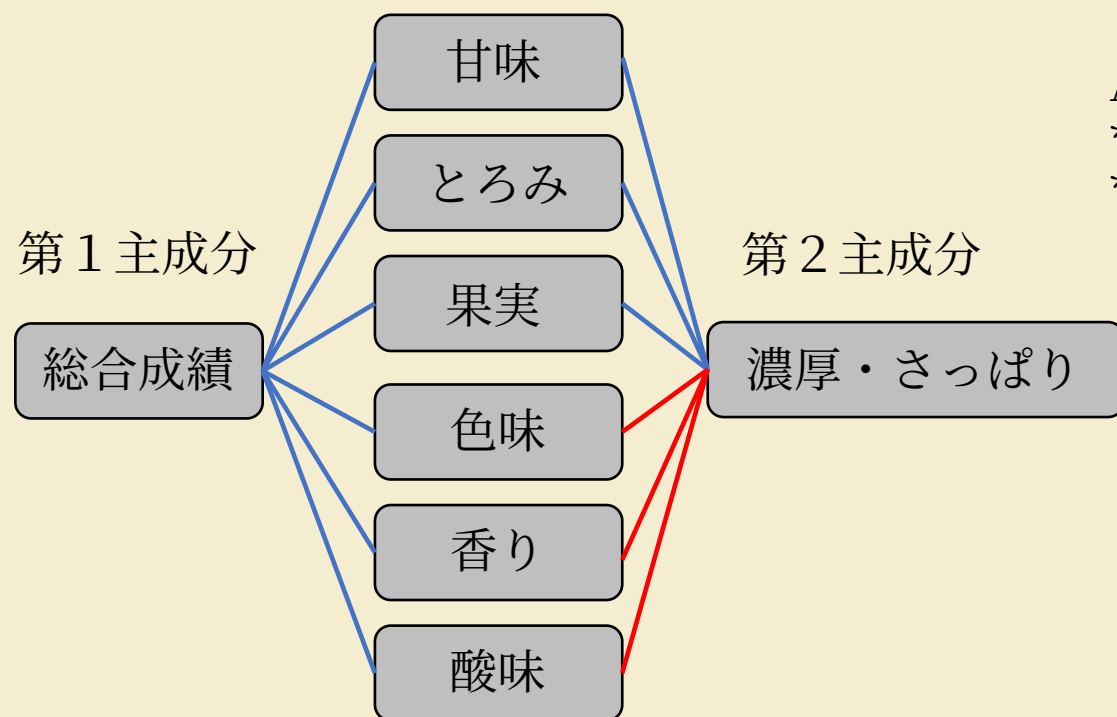
では、主成分分析とK-means法の原理と応用例を紹介していきます。

# 主成分分析 (PCA)

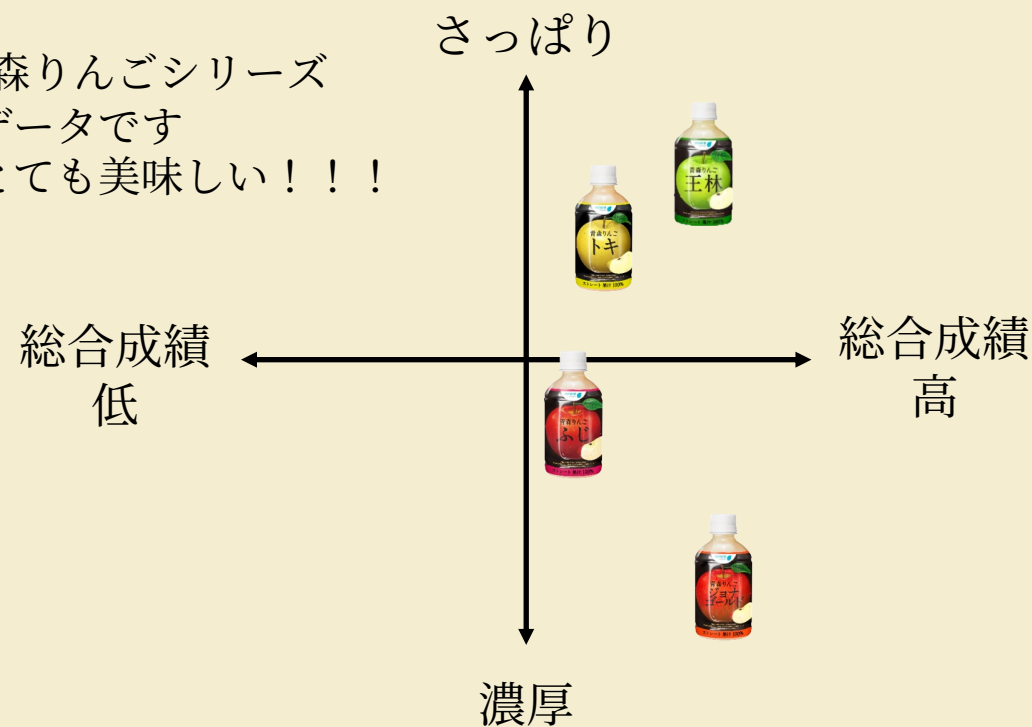
## ✓ 次元の圧縮に関する手法

多次元データから合成変数を用いて低次元空間に情報を縮約・圧縮

2・3次元へ圧縮 → データの視覚化 → 解釈のしやすさ

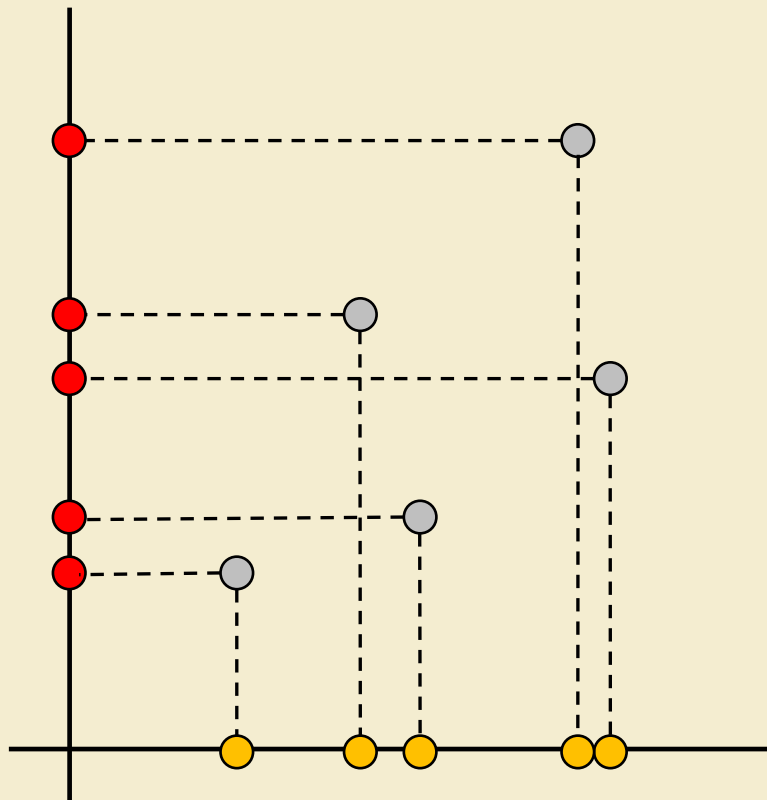


Acure青森りんごシリーズ  
\*架空のデータです  
\*どれもとても美味しい!!!



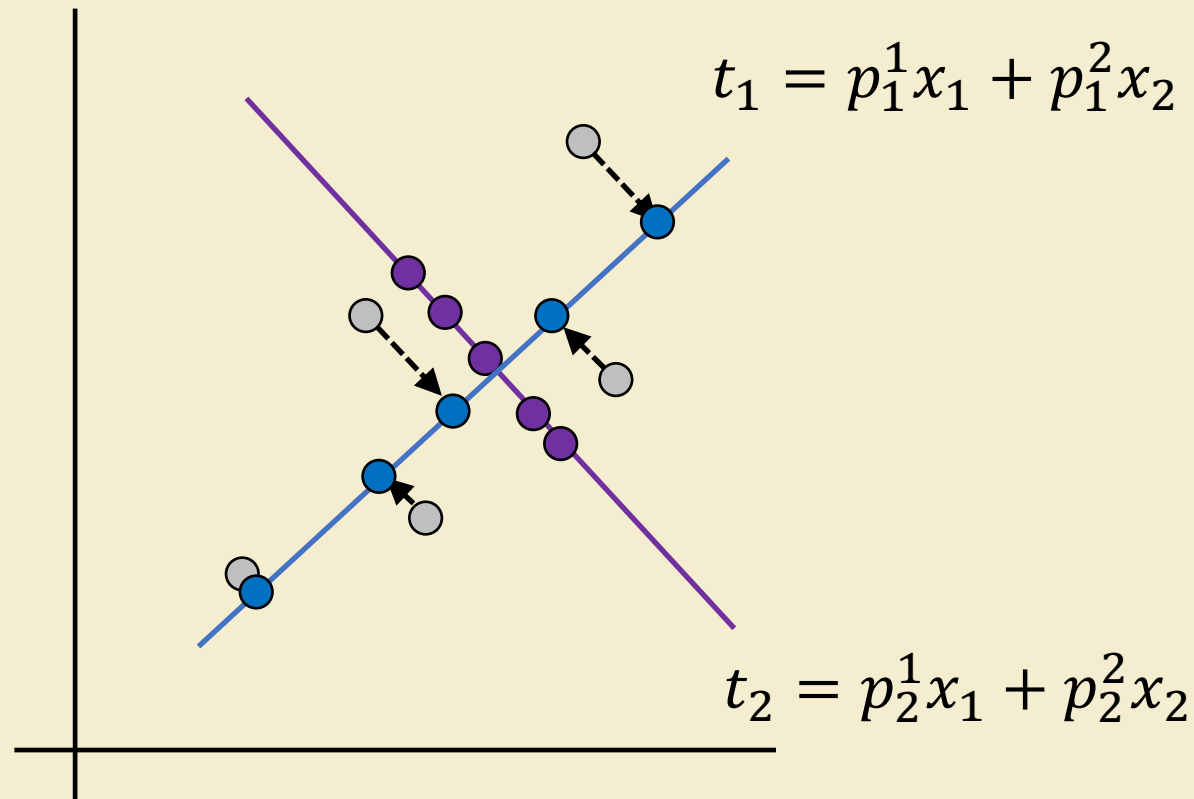
# 主成分分析の図解

## 1次元への圧縮



もちろん次元圧縮は情報量の減少を意味する  
データのばらつきが大きいほど情報量が多い

## 射影したデータの分散の最大化



$t_i$  第  $i$  主成分

$p_i^j$  第  $i$  主成分に対応する  $j$  番目の重み

# じゃあ、どう解くの？

N個サンプルがある場合

$x_1$	$x_2$
$x_1^1$	$x_2^1$
$x_1^2$	$x_2^2$
...	...
$x_1^n$	$x_2^n$

$$1 \text{ 個目 } t_1 = p_1^1 x_1 + p_1^2 x_2^1 \quad t_2 = p_2^1 x_1 + p_2^2 x_2$$

$$2 \text{ 個目 } t_1 = p_1^1 x_1 + p_1^2 x_2^1 \quad t_2 = p_2^1 x_1 + p_2^2 x_2$$

$$\vdots \quad \quad \quad \vdots \quad \quad \quad \vdots$$

$$n \text{ 個目 } t_1 = p_1^1 x_1 + p_1^2 x_2^1 \quad t_2 = p_2^1 x_1 + p_2^2 x_2$$

行列で書くと

$$\begin{bmatrix} t_1^1 & t_2^1 \\ \vdots & \vdots \\ t_1^n & t_2^n \end{bmatrix} = \begin{bmatrix} x_1^1 & x_2^1 \\ \vdots & \vdots \\ x_1^n & x_2^n \end{bmatrix} \begin{bmatrix} p_1^1 & p_2^1 \\ p_1^2 & p_2^2 \end{bmatrix}$$

$t_i^k$  第  $k$  個目のサンプルの  $i$  主成分の値

$x_i^k$  第  $k$  個目のサンプルの  $i$  番目の変数の値

$p_i^j$  第  $i$  主成分に対応する  $j$  番目の重み

# 主成分の値の二乗和を最大化

$$\begin{bmatrix} t_1^1 & t_2^1 \\ \vdots & \vdots \\ t_1^n & t_2^n \end{bmatrix} = \begin{bmatrix} x_1^1 & x_2^1 \\ \vdots & \vdots \\ x_1^n & x_2^n \end{bmatrix} \begin{bmatrix} p_1^1 & p_2^1 \\ p_1^2 & p_2^2 \end{bmatrix}$$

既知                      変数

- $t_i^k$  第 k 個目のサンプルの i 主成分の値
- $x_i^k$  第 k 個目のサンプルの i 番目の変数の値
- $p_i^j$  第 i 主成分に対応する j 番目の重み

## 主成分の分散の最大化

主成分の値の二乗和を最大化させることに対応 (変数のセンタリング後)

$$S = \sum_{i=1}^n (t_1^i)^2 = (p_1^1)^2 \sum_{i=1}^n (x_1^i)^2 + 2p_1^1 p_1^2 \sum_{i=1}^n x_1^i x_2^i + (p_1^2)^2 \sum_{i=1}^n (x_2^i)^2$$

この時、規格化条件  $(p_1^1)^2 + (p_1^2)^2 = 1$

Lagrangeの未定乗数法でSの最大化を行う

# Lagrangeの未定乗数法

$\lambda$ を未知の定数としてGが最大になる $\lambda, p_1^1, p_1^2$ を決める

$$S = (p_1^1)^2 \sum_{i=1}^n (x_1^i)^2 + 2p_1^1 p_1^2 \sum_{i=1}^n x_1^i x_2^i + (p_1^2)^2 \sum_{i=1}^n (x_2^i)^2 \quad \text{制約条件} \quad (p_1^1)^2 + (p_1^2)^2 = 1$$

$$G = S - \lambda \{(p_1^1)^2 + (p_1^2)^2 - 1\}$$

$$= (p_1^1)^2 \sum_{i=1}^n (x_1^i)^2 + 2p_1^1 p_1^2 \sum_{i=1}^n x_1^i x_2^i + (p_1^2)^2 \sum_{i=1}^n (x_2^i)^2 - \lambda \{(p_1^1)^2 + (p_1^2)^2 - 1\}$$

Gの最大値  $\rightarrow$  Gの極大値  $\rightarrow \lambda, p_1^1, p_1^2$ の偏微分が0となるとき

# 偏微分

偏微分すると…

$$\left(\sum_{i=1}^n (x_1^i)^2 - \lambda\right) p_1^1 + \left(\sum_{i=1}^n x_1^i x_2^i\right) p_1^2 = 0$$

$$\left(\sum_{i=1}^n x_1^i x_2^i\right) p_1^1 + \left(\sum_{i=1}^n (x_2^i)^2 - \lambda\right) p_1^2 = 0$$



$$\begin{bmatrix} \sum_{i=1}^n (x_1^i)^2 - \lambda & \sum_{i=1}^n x_1^i x_2^i \\ \sum_{i=1}^n x_1^i x_2^i & \sum_{i=1}^n (x_2^i)^2 - \lambda \end{bmatrix} \begin{bmatrix} p_1^1 \\ p_1^2 \end{bmatrix} = 0$$

書き換えると…

$$(\mathbf{X}^T \mathbf{X} - \lambda \mathbf{E}) \mathbf{p}_1 = 0 \quad \mathbf{X} = \begin{bmatrix} x_1^1 & x_2^1 \\ \vdots & \vdots \\ x_1^n & x_2^n \end{bmatrix} \quad \mathbf{E} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \quad \mathbf{p}_1 = \begin{bmatrix} p_1^1 \\ p_1^2 \end{bmatrix}$$

分散の最大化には,  $p_1^1 = p_1^2 = 0$ 以外の解が必要!!! という問題になる

# 固有値問題へ

$p_1^1 = p_1^2 = 0$ 以外の解が必要



逆行列を持たない  $\rightarrow (X^T X - \lambda E)$  の行列式が0である必要

$\lambda$  を固有値,  $p_1, p_2$  を固有ベクトルとする固有値問題

$p_1, p_2$  を求め, 対応する主成分が得られる!

今までやっていたのは分散共分散行列の固有値問題を解いていることに

$$\Sigma = \begin{bmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{12} & \sigma_2^2 \end{bmatrix} = \begin{bmatrix} \frac{1}{n} \sum_{i=1}^n (x_1^i - m)^2 & \frac{1}{n} \sum_{i=1}^n (x_1^i - m)(x_2^i - m) \\ \frac{1}{n} \sum_{i=1}^n (x_1^i - m)(x_2^i - m) & \frac{1}{n} \sum_{i=1}^n (x_2^i - m)^2 \end{bmatrix}$$



# 主成分分析の実装

使用するデータ: 1973年の車に関するデータセット (32 × 11)

	mpg	cyl	disp	hp	drat	wt	qsec	vs	am	gear	carb
<b>Mazda RX4</b>	21.0	6	160.0	110	3.90	2.620	16.46	0	1	4	4
<b>Mazda RX4 Wag</b>	21.0	6	160.0	110	3.90	2.875	17.02	0	1	4	4
<b>Datsun 710</b>	22.8	4	108.0	93	3.85	2.320	18.61	1	1	4	1
<b>Hornet 4 Drive</b>	21.4	6	258.0	110	3.08	3.215	19.44	1	0	3	1
<b>Hornet Sportabout</b>	18.7	8	360.0	175	3.15	3.440	17.02	0	0	3	2
<b>Valiant</b>	18.1	6	225.0	105	2.76	3.460	20.22	1	0	3	1
<b>Duster 360</b>	14.3	8	360.0	245	3.21	3.570	15.84	0	0	3	4

mpg	Miles/(US) gallon	qsec	1/4 mile time
cyl	Number of cylinders	vs	Engine (0 = V-shaped, 1 = straight)
disp	Displacement (cu.in.)	am	Transmission (0 = automatic, 1 = manual)
hp	Gross horsepower	gear	Number of forward gears
drat	Rear axle ratio	carb	Number of carburetors
wt	Weight (1000 lbs)		

# Princomp vs. Prcomp

	Princomp	Prcomp
分析対象となるデータ	R mode (行数 >> 列数)	R mode (行数 >> 列数) Q mode (行数 << 列数)
目的	サンプルに対する変数の解析 (特徴量の抽出)	変数に対してサンプルの解析 (データ行列Xの転置行列に)
計算アルゴリズム	共分散行列の固有値	特異値分解
固有値	主成分得点の不偏分散	主成分得点の標本分散

- データボックス: イギリスの心理学者レイモンド・キャッテルが導入  
データを3つの次元 (サンプル、変数、時間) から考えるもの
- **迷ったら, prcomp()**

なぜ共分散行列に基づくアルゴリズムでは, 行数 >> 列数が問題となるのか...

- データ行列の次元:  $n \times v$  (変数:  $v$  個, 観測値:  $n$  個)
- 共分散行列を用いてRモードで解析: 分散共分散行列  $C$  の次元は  $v \times v$ , ランクは  $n$  と  $v$  の小さい方
- 共分散行列を用いてQモードで解析: 分散共分散行列  $C$  の次元は  $n \times n$ , ランクは  $n$  と  $v$  の小さい方  
 $n > v$  の時 → 固有値および固有ベクトルが  $v$  個しか求まらない → 計算ストップ



# データの前処理

## centering

各変数の平均を0にする

$$x_i - \bar{\mu}$$

## standardization

各変数の平均を0、分散を1にする

$$\frac{x_i - \bar{\mu}}{\sigma}$$

## scaling

各変数の分散を1にする

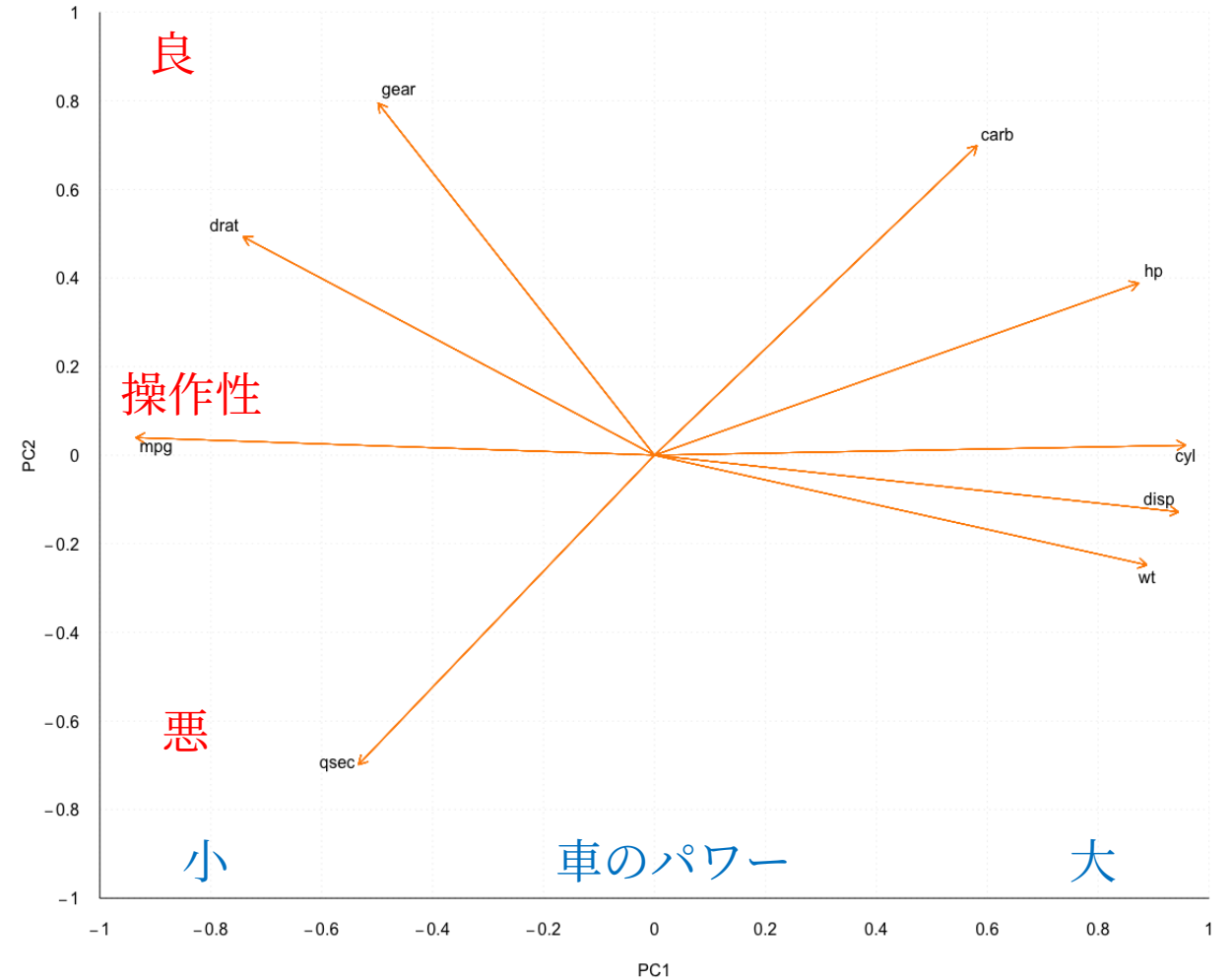
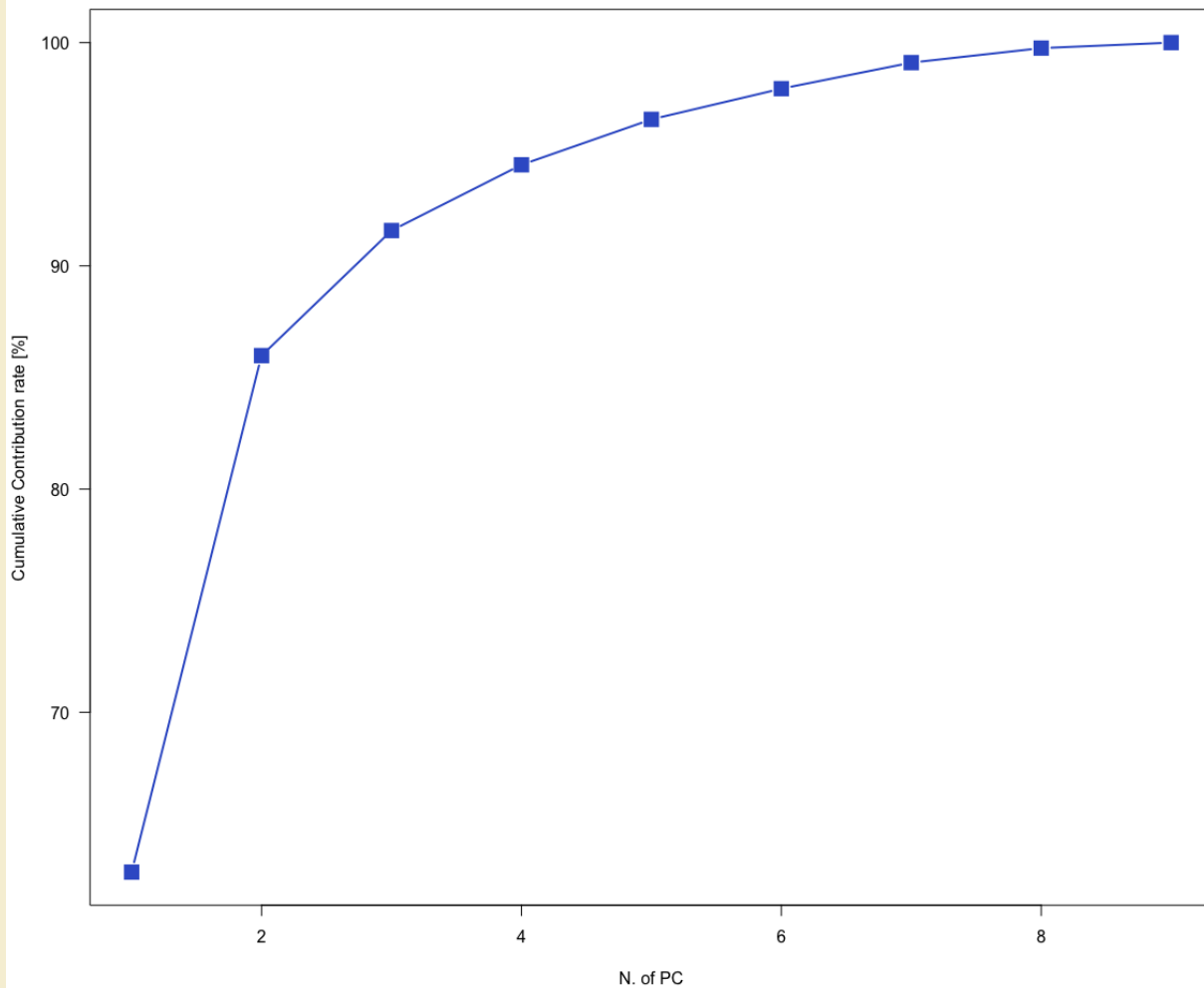
$$x_i / \sigma$$



単位系が異なる場合なども、  
各変数が同等の重みを持つようになる

# 主成分分析の実装

## 結果の視覚化



# 主成分分析のまとめ

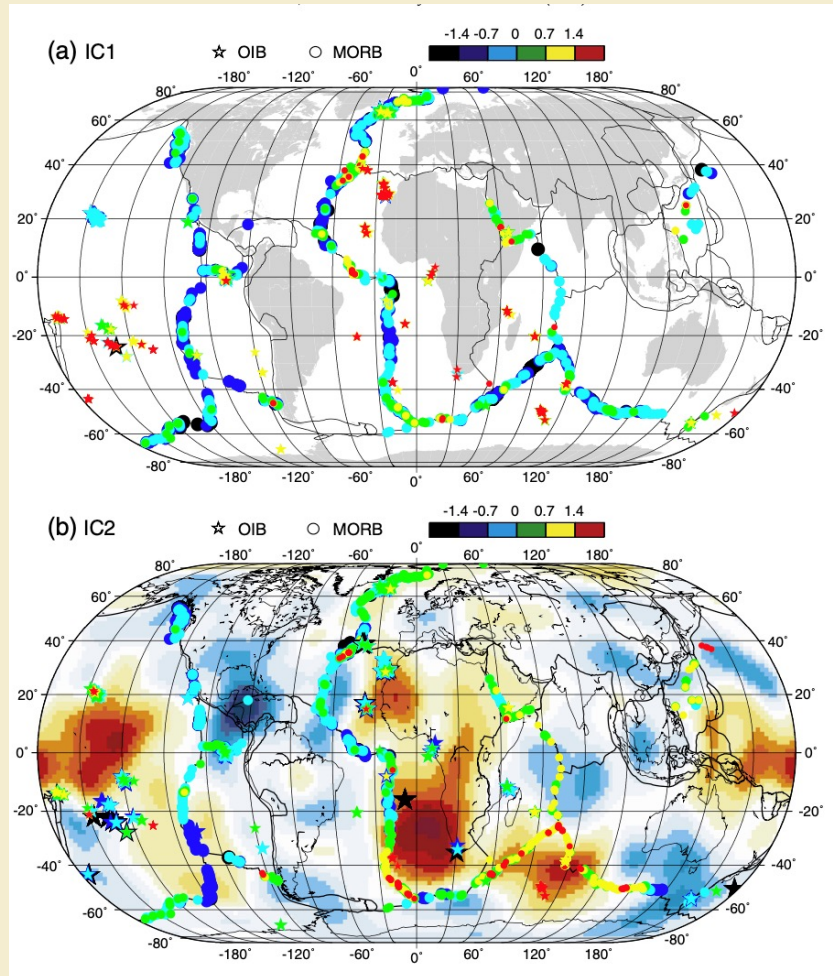
- 次元を圧縮, 互いに無相関な成分を計算
- 分散共分散行列の固有値固有ベクトル分解に相当
- N次元のデータからはN個の主成分  
(少ない成分で説明できるほど嬉しい)
- データ構造の本質的な要因を抽出できる可能性



# 地球化学への応用: 主成分分析

Iwamori + (2010) *EPSL*

独立成分分析(PCAを拡張したもの)を海嶺・海洋島玄武岩の同位体組成解析 (Pb, Nd, Sr, Hf) に応用



IC1: MORBと海洋島玄武岩OIB を分ける

IC2: 地理的分布, 特にDUPAL異常を識別

そこから得られた知見

DMM, EMのような地球化学的マントル端成分の相互作用ではなく, 2つの独立な物質分化プロセスの重複を示唆

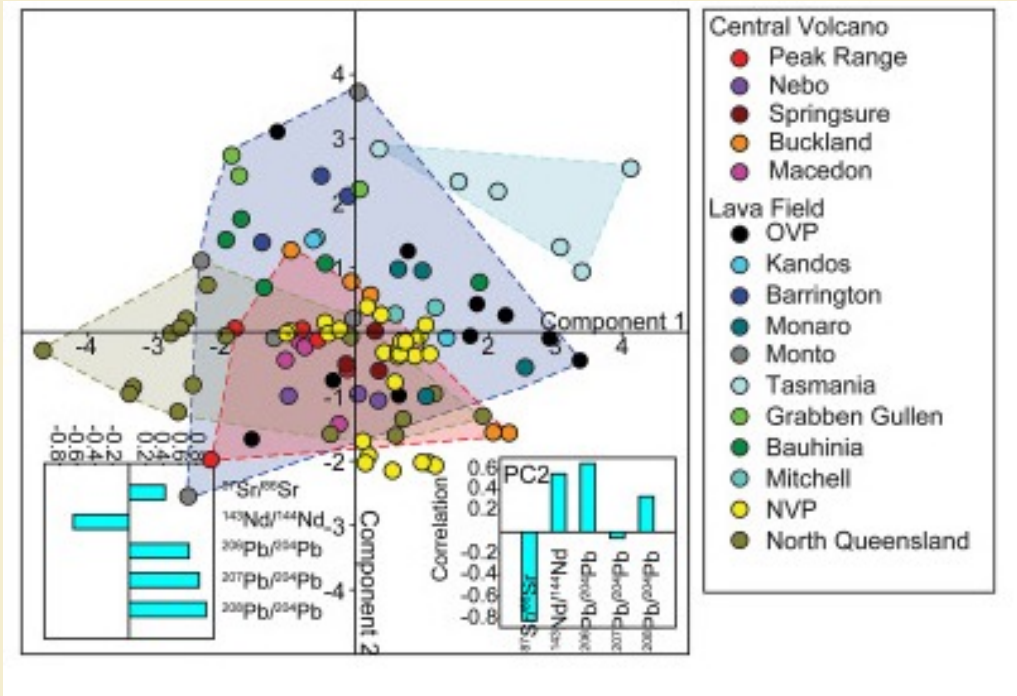
IC1: マントル溶融プロセスで元素分別, 多様性形成

IC2: 沈み込み帯における脱水・加水でIC2方向の多様性

これらが再循環し, 海嶺で再溶融するという, 2つの独立な分別作用の相互重複プロセスがマントルの不均質の原因と推定

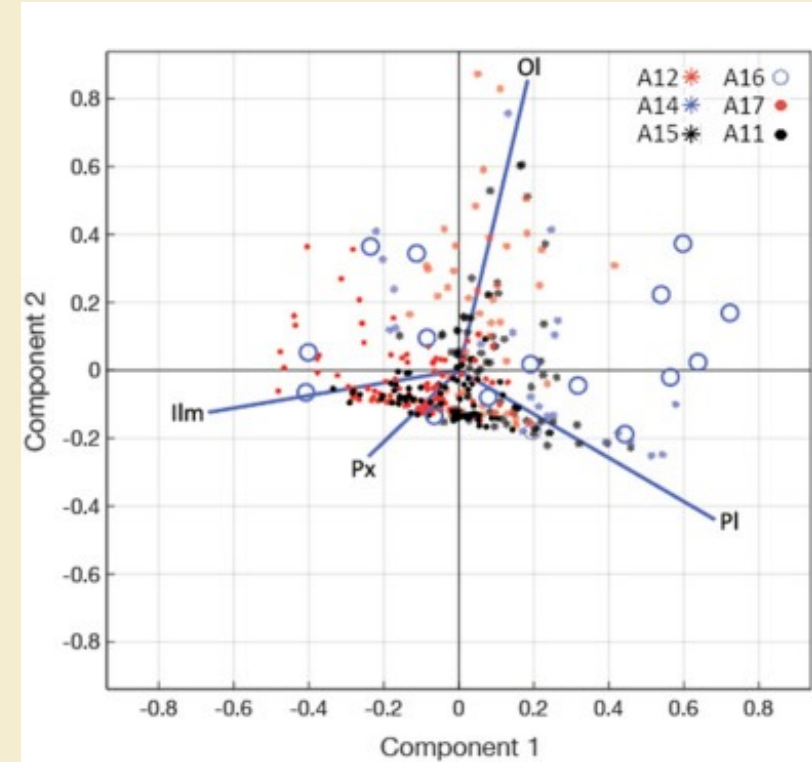
# 地球化学への応用: PCA

Jones + (2020) *Lithos*



東オーストラリアの新生代火山の同位体, 微量元素などのデータセットに主成分分析

Cone + (2020) *Icarus*



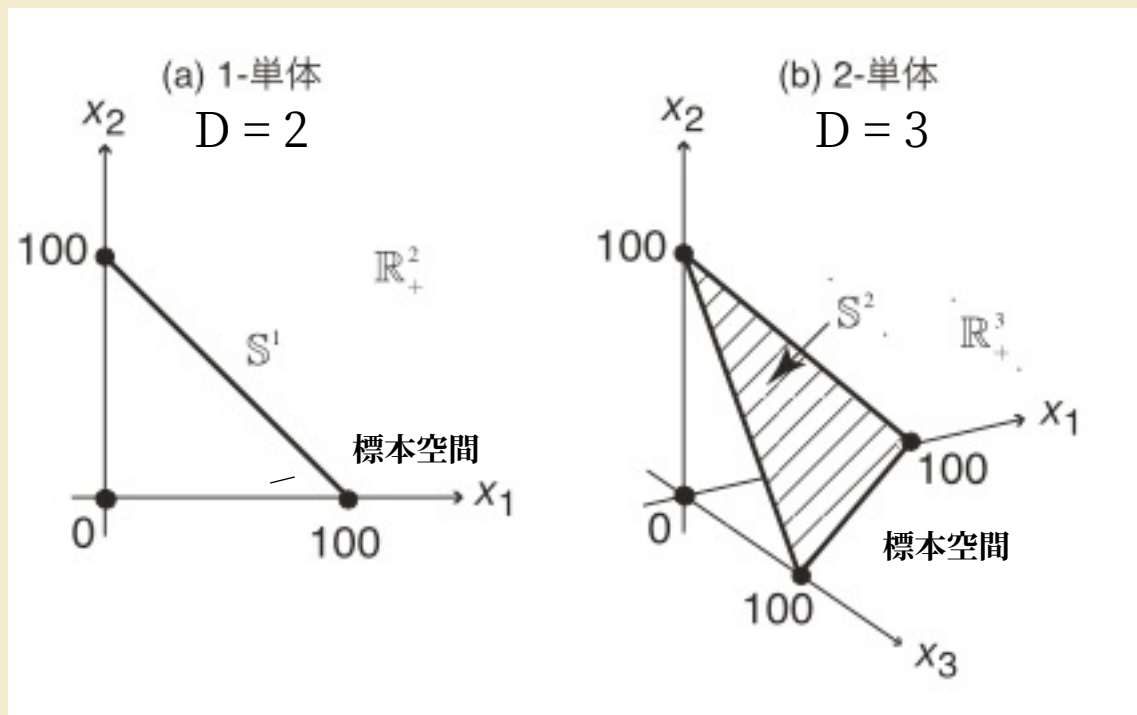
Apollo sampleの玄武岩データベースの化学組成や鉱物モードに適用



# 組成データの注意点

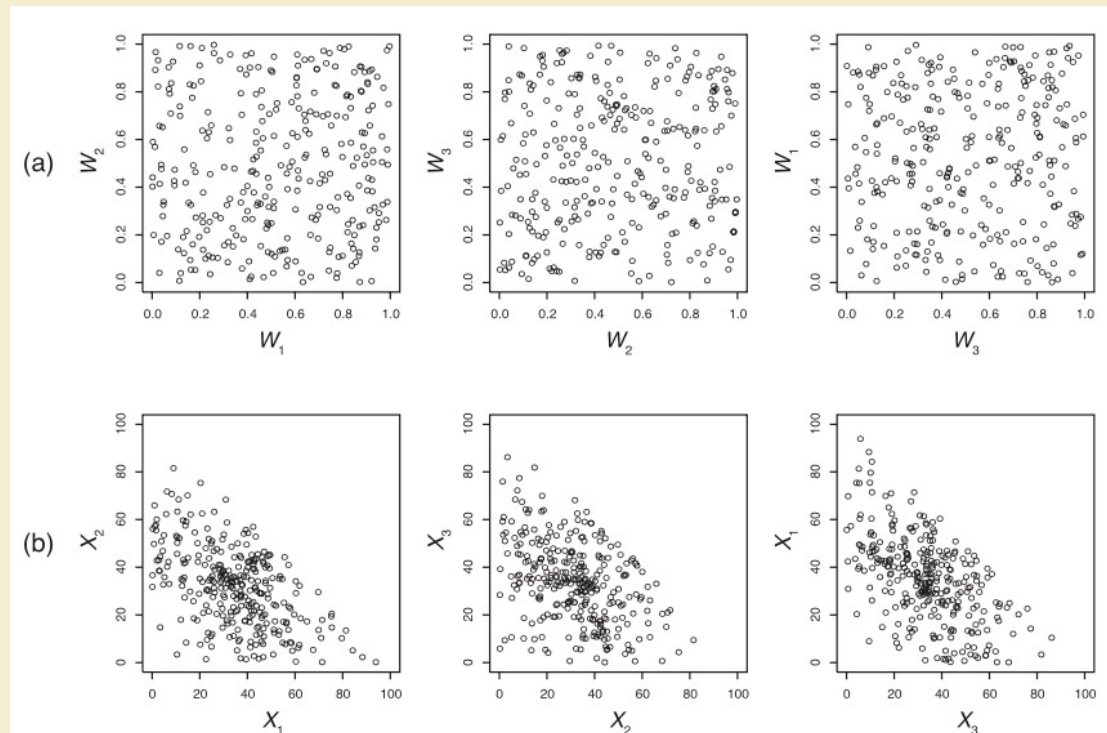
Aitchison (1982, 1986)

- 組成データ: 変数が全て正, 変数の総和が定数
- 定数和制約: 変数の独立性は保たれない



Butler(1978, 1979a)

統計学的な性質に注意が必要



位相幾何学的には  $D-1$  次元空間を単体空間( $S^{D-1}$ )

$$S^{D-1} = \left\{ (x_1, \dots, x_D) \mid x_i > 0 \ (i=1, \dots, D), \sum_{i=1}^D x_i = 100 \right\}$$

3次元の組成データは3次元空間に属する 2-単体空間に配列する

- データはもともとの3次元的な広がりを持たず, ある平面の領域に集約
- 無理に直行座標系にプロットすると, 実在しない配列や散布度が出現

# 対数比変換・対数比解析

組成データを  $S^{D-1}$  から  $R^D$  に写像する



写像後の対数比データに対して、  
実数と同等に演算と統計解析が実行可能

= 定数和制約からデータを開放

相加対数比変換

(alr: additive log-ratio transformation)

$$alr(x) = \left( \ln \frac{x_1}{x_D}, \ln \frac{x_2}{x_D}, \dots, \ln \frac{x_{D-1}}{x_D} \right)^t = y(y_1, y_2, \dots, y_{D-1})^t$$

有心対数比変換

(clr: centered log-ratio transformation)

$$clr(x) = \left( \ln \frac{x_1}{g(x)}, \ln \frac{x_2}{g(x)}, \dots, \ln \frac{x_D}{g(x)} \right)^t = z(z_1, z_2, \dots, z_D)^t$$

\* 変換結果として使用できる変数は  
元の変数から一つ減ってしまう

$$g(x) = \left( \prod_{i=1}^D x_i \right)^{1/D}$$