

# 多変量高次元データ解析の理論と応用

研究代表者 杉山高一 研究員

## 多変量高次元相関行列の各固有値の同等性検定 ～パーミュテーションテストを用いて～

多変量高次元ベクトル変量  $X^{(i)}$ ,  $i = 1, 2$

平均ベクトル  $\mu^{(i)}$   
相関行列  $P^{(i)}$

$$H_0 : \xi_j^{(1)} = \xi_j^{(2)}$$

$$H_1 : \xi_j^{(1)} \neq \xi_j^{(2)} \text{ for } j = 1, \dots, p$$

のような仮説を考える。ただし、 $\xi_j^{(i)}$  を第  $i$  母集団の  $j$  番目固有値とする

$N_i$  個の無作為標本  $x_1^{(i)}, \dots, x_{N_i}^{(i)}$  より共分散行列  $S^{(i)} = (s_{jj}^{(i)})$  を作成  
固有ベクトル  $r_1^{(i)}, \dots, r_p^{(i)}$ , ( $r_1^{(i)} \geq \dots \geq r_p^{(i)}$ )  
固有値  $b_1^{(i)}, \dots, b_p^{(i)}$

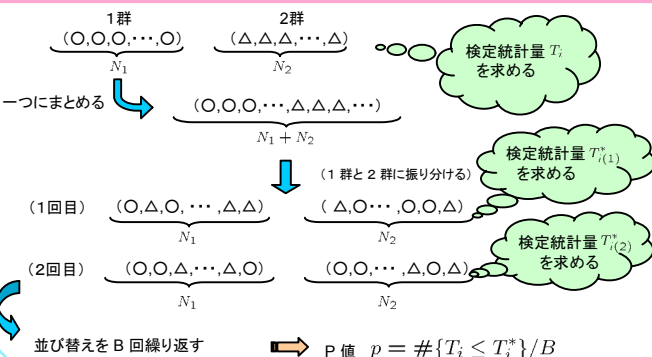
### 主成分スコア

$$z_{j\alpha}^{(i)} = b_j^{(i)} (x_{\alpha}^{(i)} - \bar{x}^{(i)}) / \sqrt{s_{jj}^{(i)}} \quad i = 1, 2, \alpha = 1, \dots, N_i$$

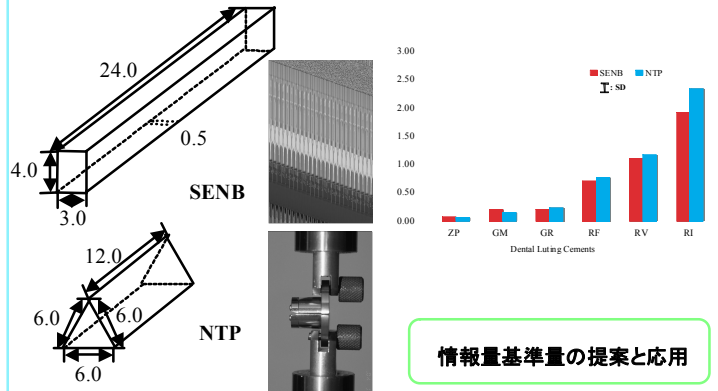
### 検定統計量

$$T_1 = \sqrt{N} \log(r_j^{(2)}/r_j^{(1)}) \quad T_2 = \sqrt{N}(r_j^{(2)} - r_j^{(1)})$$

$$T_3 = \frac{A_N - E[A_N]}{\sqrt{\text{var}[A_N]}} \left( A_N = \frac{N_1(N+1)}{2} - \sum_{i=1}^N |i - \frac{N+1}{2}| V_i \right), T_4 = \frac{\bar{M} - E[\bar{M}]}{\sqrt{\text{var}[\bar{M}]}} \left( \bar{M} = \sum_{i=1}^N \left( i - \frac{N+1}{2} \right)^2 V_i \right)$$



## 情報量基準による関数関係モデルの選択問題 歯科用合着用セメントの2種類の強度測定法 SENB; NTP の同等性研究



$$M1 : \mu_{2j} = c\mu_{1j}, (c \neq 0) \quad (j = 1, \dots, m)$$

$$M2 : \mu_{2j} = \alpha + \beta\mu_{1j}, (\alpha, \beta \neq 0) \quad (j = 1, \dots, m) \quad (m \geq 2)$$

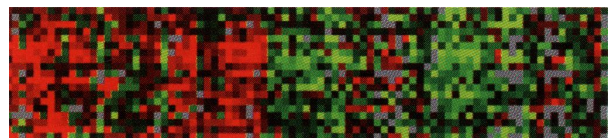
$$M3 : \mu_{2j} \neq \mu_{1j} \quad (j = 1, \dots, m)$$

M1 のもとでの推定方程式

$$\sum_{j=1}^m \frac{(c\bar{x}_{1j} - \bar{x}_{2j})(\bar{x}_{1j} + c\bar{x}_{2j})}{(1+c^2)(s_{1j}^2 + s_{2j}^2) + (c\bar{x}_{1j} - \bar{x}_{2j})^2} = 0$$

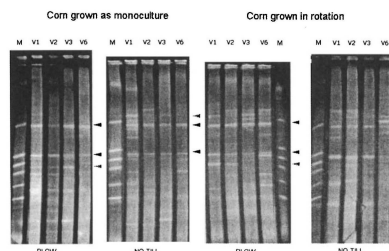
## 高次元多変量2値データの解析 ～どの遺伝子で判別できるか～

### フィンガープリントデータ



Terry, Speed: *Interdisciplinary Statistics Statistical Analysis of Gene Expression Microarray Data*. Chapman & Hall/CRC(2003)

### 4種類のとうもろこし: 84個の遺伝子(変数)



どの遺伝子で判別できる?

### 変数選択法の開発

Wilbur, J.D., Ghosh, J.K., Nakatsu, C.H., Broader, S.M., and Doerge, R.W.: Variable selection in high-dimensional multivariate binary data with application to the analysis of microbial community DNA fingerprints. *Biometrics*, 58, 378-386. (2002)

$$AIC = -n \log \frac{|W_{22.1}|}{|T_{22.1}|} + np(1 + \log 2\pi) + 2\{qk + p - k + \frac{1}{2}p(p+1)\}$$

## 高次元漸近理論の開発 ～大標本漸近理論から高次元漸近理論へ～

$$\lim_{n \rightarrow \infty} \longrightarrow \lim_{\substack{n \rightarrow \infty, p \rightarrow \infty \\ n/p \rightarrow c}}$$

- ・主成分寄与率
  - ・小さい主成分分散の同等性検定
  - ・準相関係数
- 等の高次元漸近分布を導出

$$\sigma_{m,n}^{-1} \left\{ \left( \sum_{j=q+1}^p \log l_j - m \log \frac{1}{m} \sum_{j=q+1}^p l_j \right) - \mu_{m,n} \right\} \xrightarrow{d} N(0, 1)$$

$l_j$  : 標本共分散行列の  $j$  番目に大きい固有値

$$\mu_{m,n} = m \log m - m\psi\left(\frac{mn}{2}\right) + \psi_m\left(\frac{n}{2}\right)$$

$$\sigma_{m,n}^2 = \psi'_m\left(\frac{n}{2}\right) - m^2\psi'\left(\frac{mn}{2}\right), \psi_m a = \sum_{j=1}^m \psi\left(a - \frac{1}{2}(j-1)\right)$$