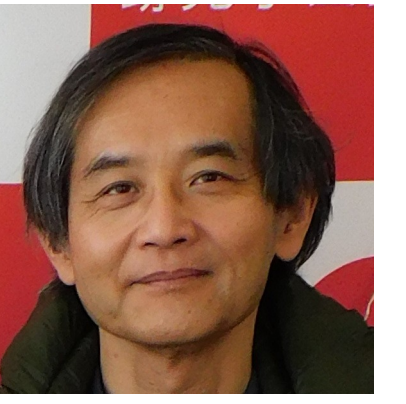
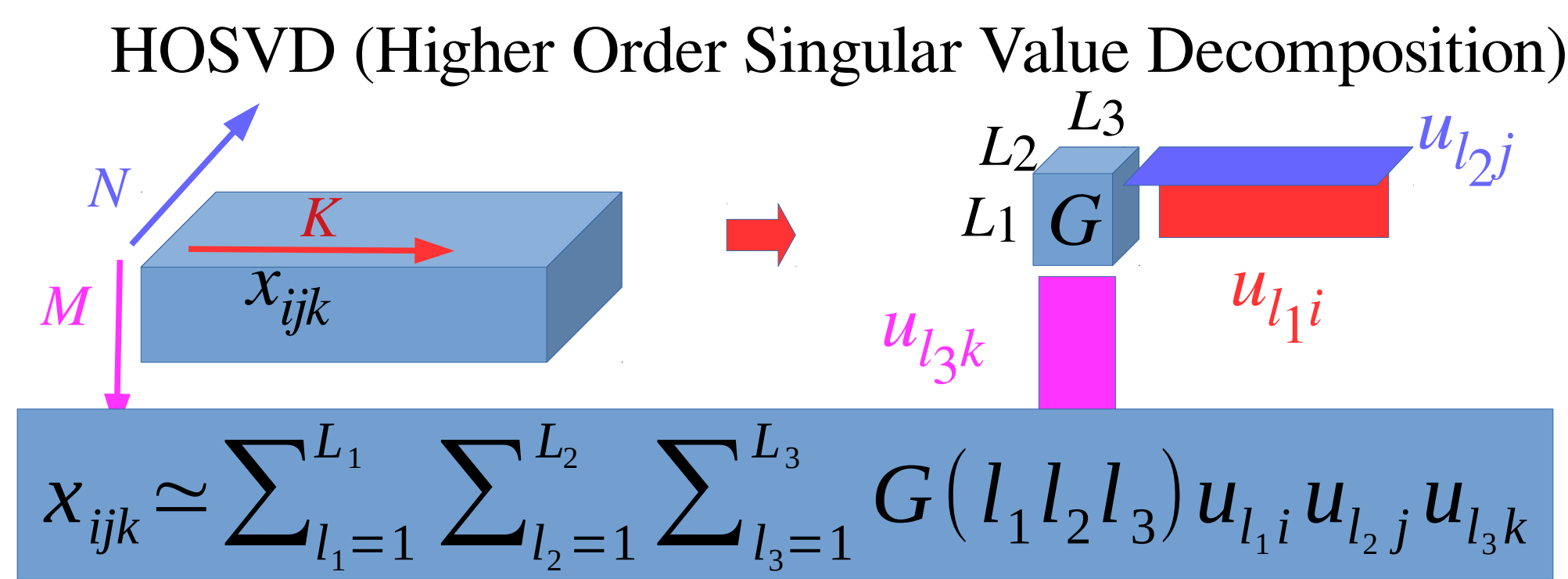


# テンソル分解を用いた教師無し学習による変数選択法の のバイオインフォマティクスへの応用

中央大学理工学部物理学科 田口善弘

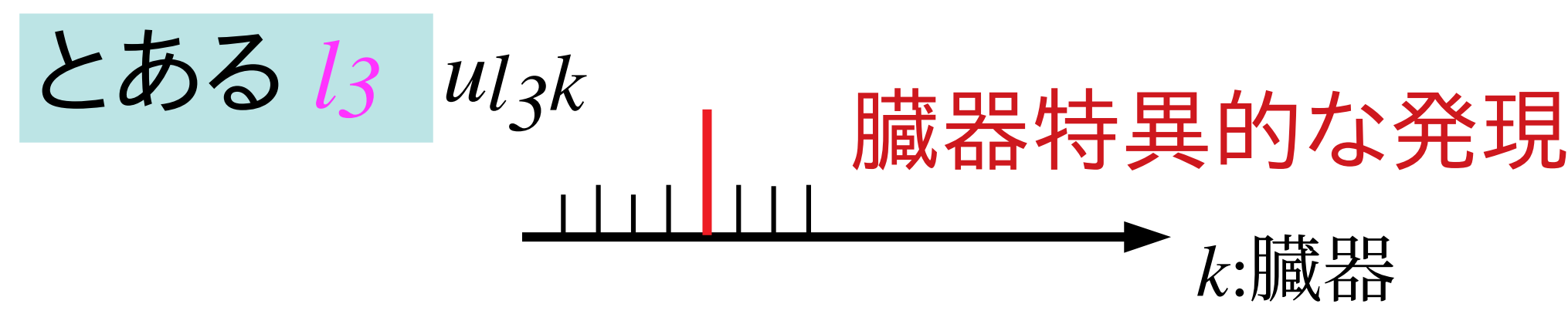
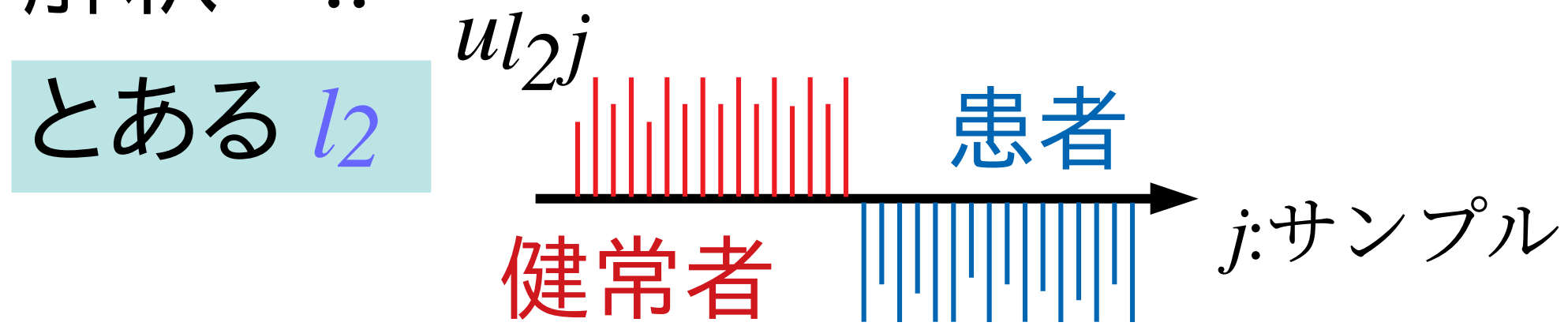


## ※テンソル分解とは？

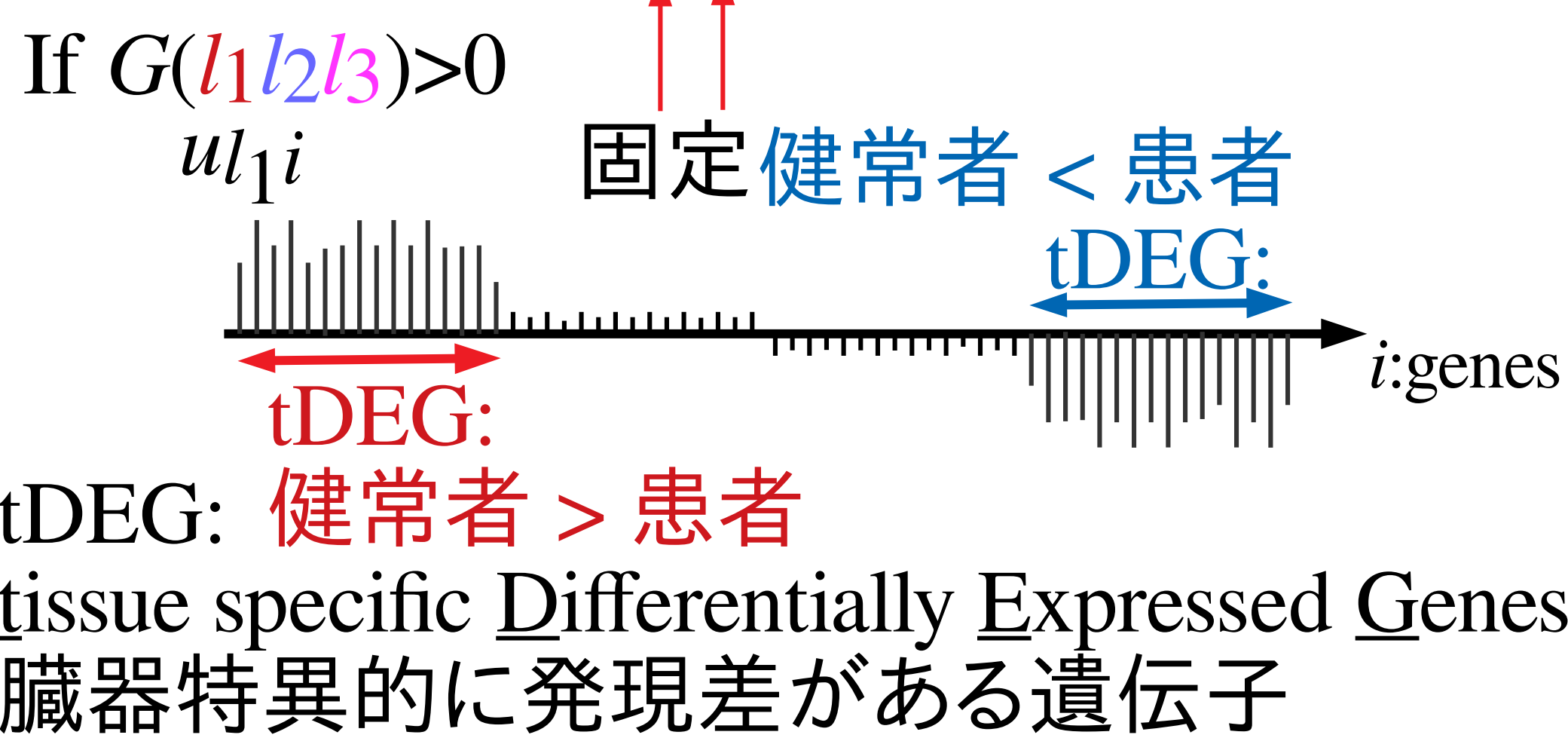


例  $x_{ijk}$ : 遺伝子  
 $N$ : 遺伝子数 ( $I$ ),  $M$ : サンプル数 ( $J$ ),  
 $K$ : 臓器数 ( $K$ )

解釈……



とある  $l_1$  が  $|G(l_1 l_2 l_3)|$  最大とする



遺伝子の選択基準:

$u_{1i}$  がガウス分布であると仮定 (帰無仮説)

→ 累積  $\chi^2$  分布で遺伝子に P 値を付与

$$P_i = P_{\chi^2} \left[ > \left( \frac{u_{1i}}{\sigma} \right)^2 \right] \quad (1)$$

→ P 値を多重比較補正 (Benjamini-Hochberg 基準)

→ 補正 P 値がしきい値 (通常 0.01) 以下の遺伝子を選択。

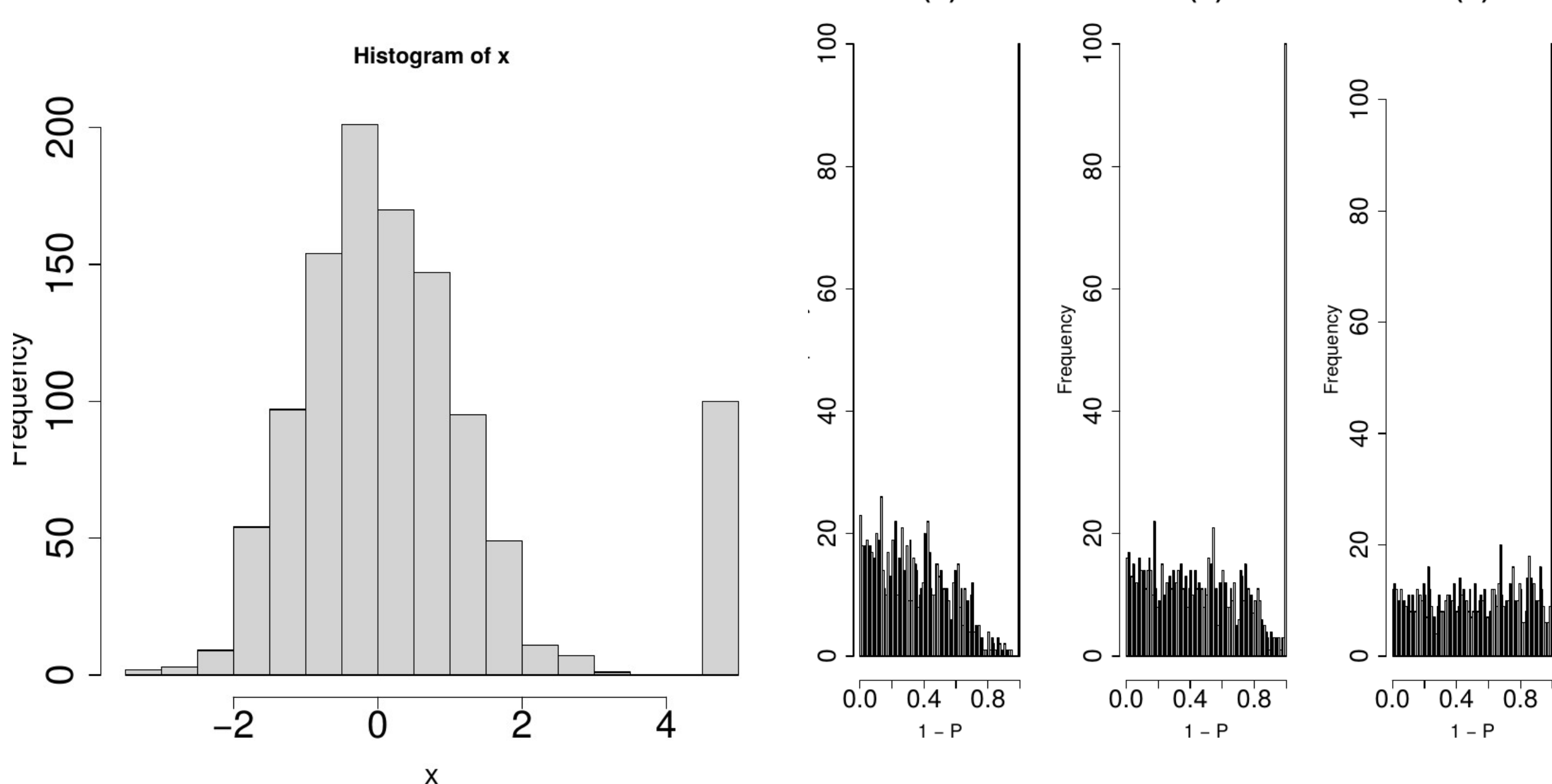
最近、 $\sigma_{1i}$  の最適化をすると性能が大きく向上することを発見

1) Tensor decomposition- and principal component analysis-based unsupervised feature extraction to select more reasonable differentially expressed genes: Optimization of standard deviation versus state-of-art methods, Y-h. Taguchi, Turki Turki bioRxiv 2022.02.18.481115; doi: <https://doi.org/10.1101/2022.02.18.481115>

2) Principal component analysis- and tensor decomposition-based unsupervised feature extraction to select more reasonable differentially methylated cytosines: Optimization of standard deviation versus state-of-the-art methods, Y-H. Taguchi, Turki Turki bioRxiv 2022.04.02.486807; doi: <https://doi.org/10.1101/2022.04.02.486807>

3) Tensor decomposition and principal component analysis-based unsupervised feature extraction outperforms state-of-the-art methods when applied to histone modification profiles Sanjiban Sekhar Roy, Y-h. Taguchi bioRxiv 2022.04.29.490081; doi: <https://doi.org/10.1101/2022.04.29.490081>

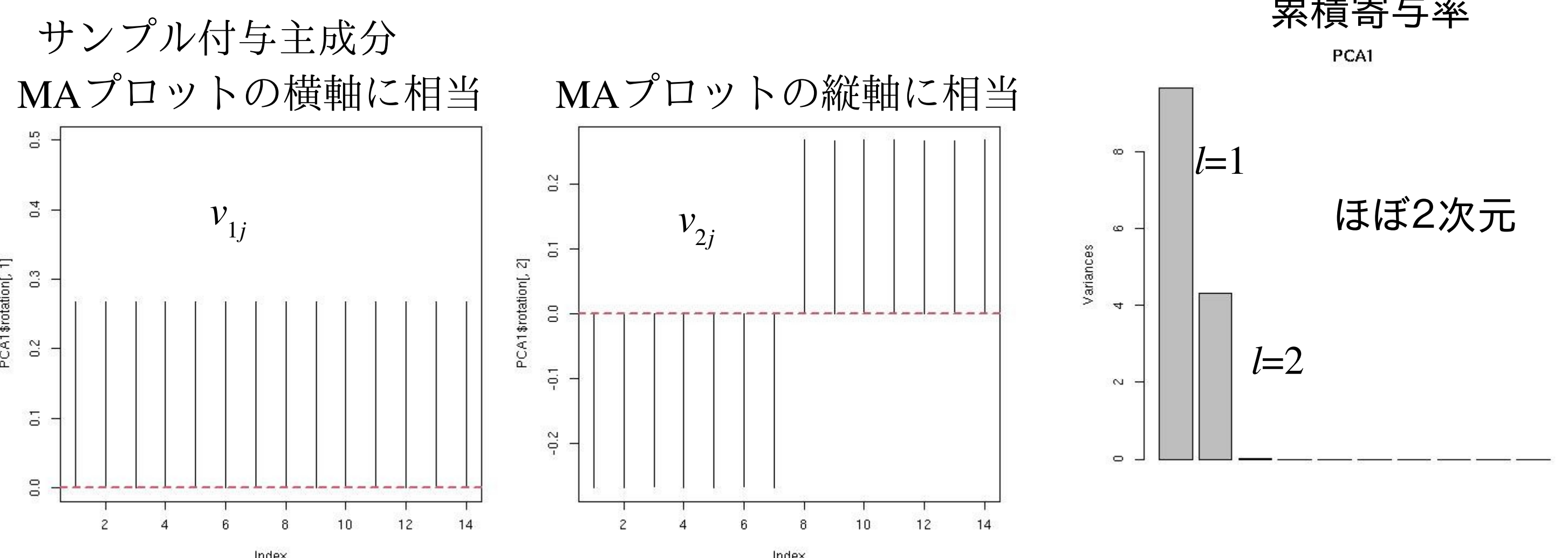
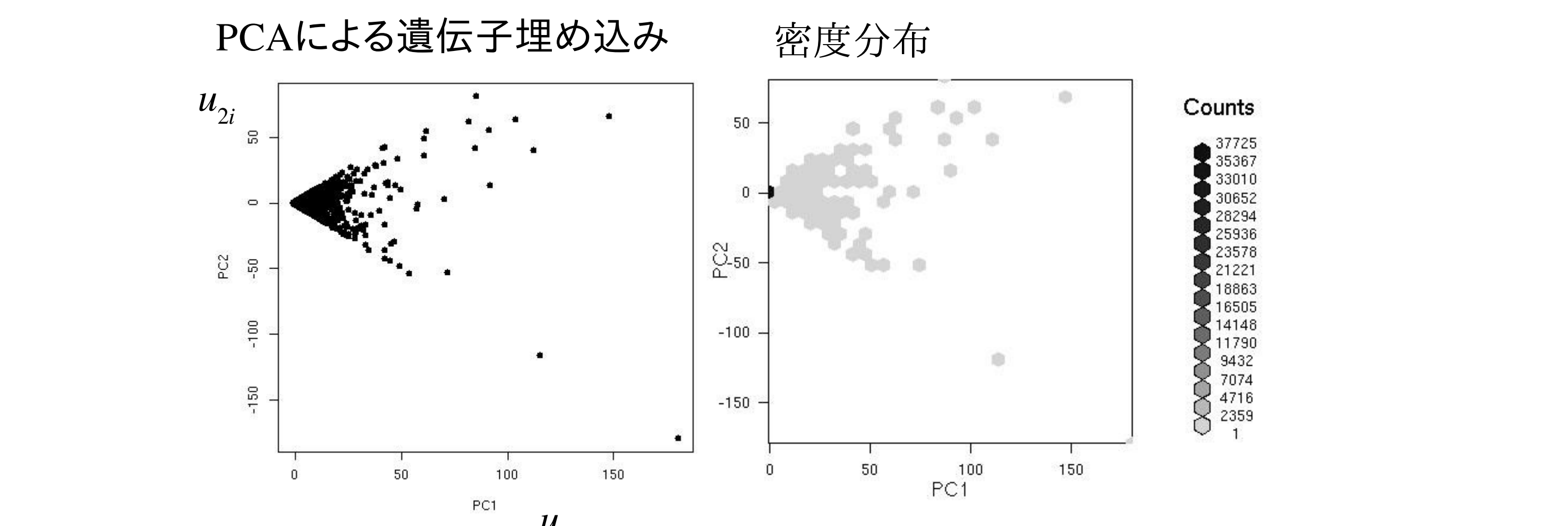
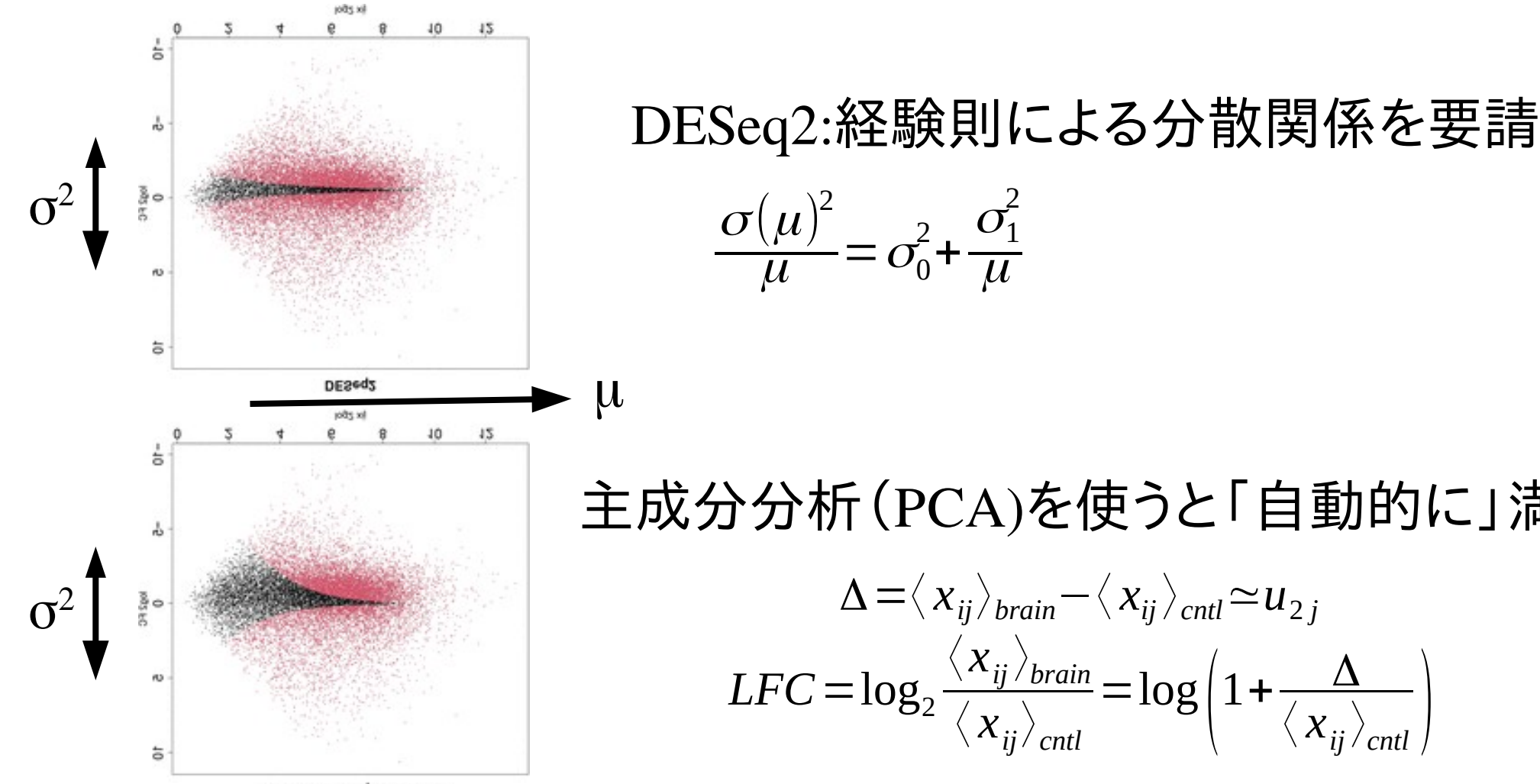
○例: 外れ値があるガウス分布



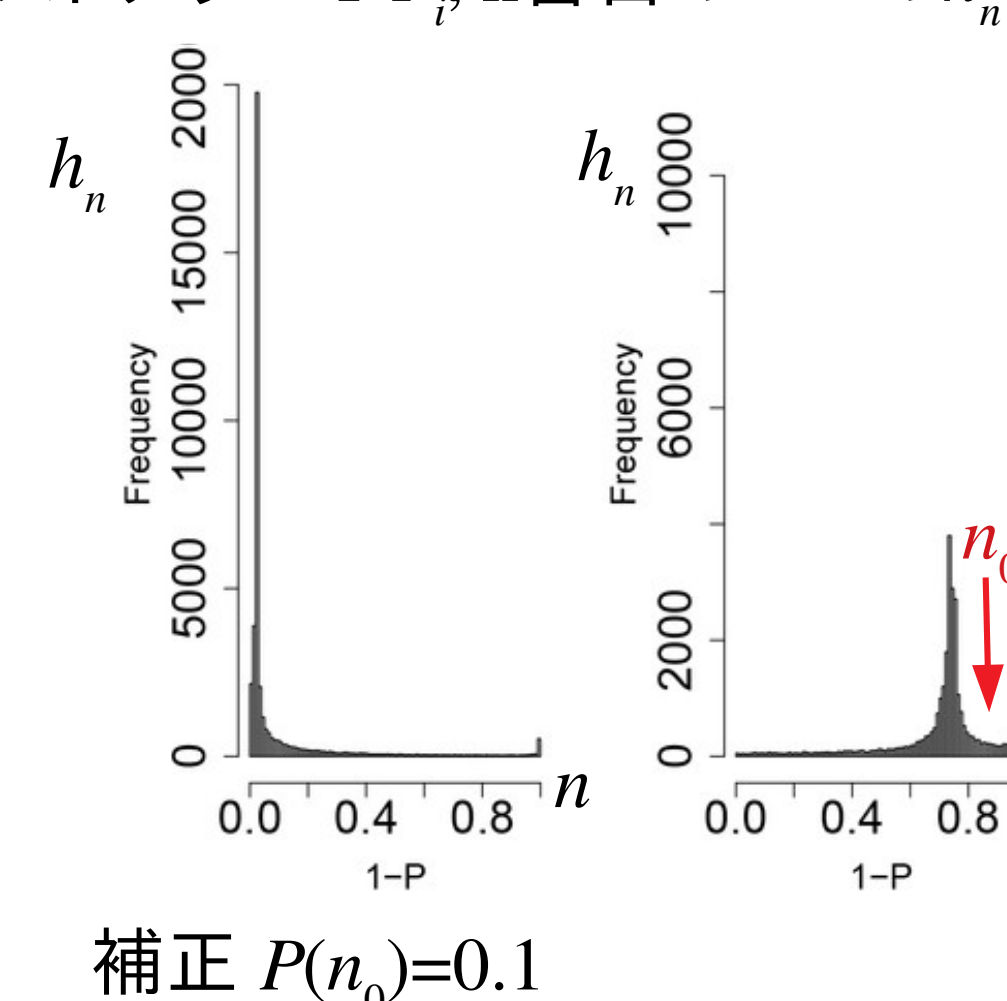
(A) 左の分布から  $\sigma$  を計算してその値から (1) 式で P 値を出した場合の  $1-P$  のヒストグラム  
 (B)  $\sigma$  を最適化した場合の P 値を用いた場合の  $1-P$  のヒストグラム  
 (C) ガウス分布を生成した  $\sigma$  を用いた場合の P 値の  $1-P$  のヒストグラム

応用例1 遺伝子発現プロファイルへの応用 (文献1)

「発現量の多い遺伝子を選ばれやすいべき」という要請が「自動的に」満たされる。



ヒストグラム  $1-P$ ,  $n$  番目の bin の  $h_n$



帰無仮説

$u_{2i}$  は Gaussian 分布 →

$P_i = P_{\chi^2} \left[ > \left( \frac{u_{2i}}{\sigma} \right)^2 \right]$  累積  $\chi^2$  分布

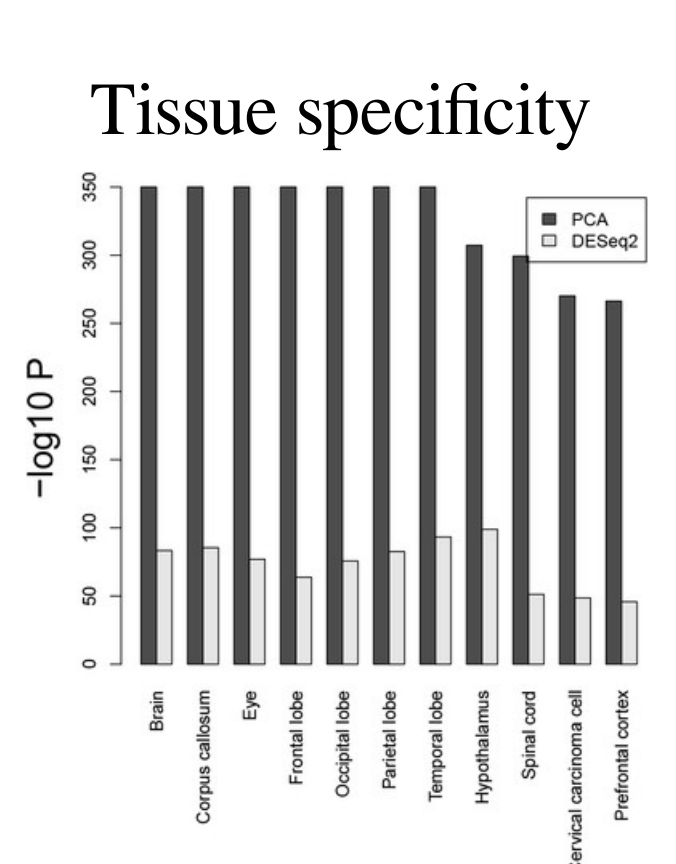
左:  $\sigma = \sqrt{\frac{\sum_i (u_{2i} - \langle u_{2i} \rangle)^2}{N}}$

右: 最適  $\sigma$  は  $\sigma_n$  を最小化

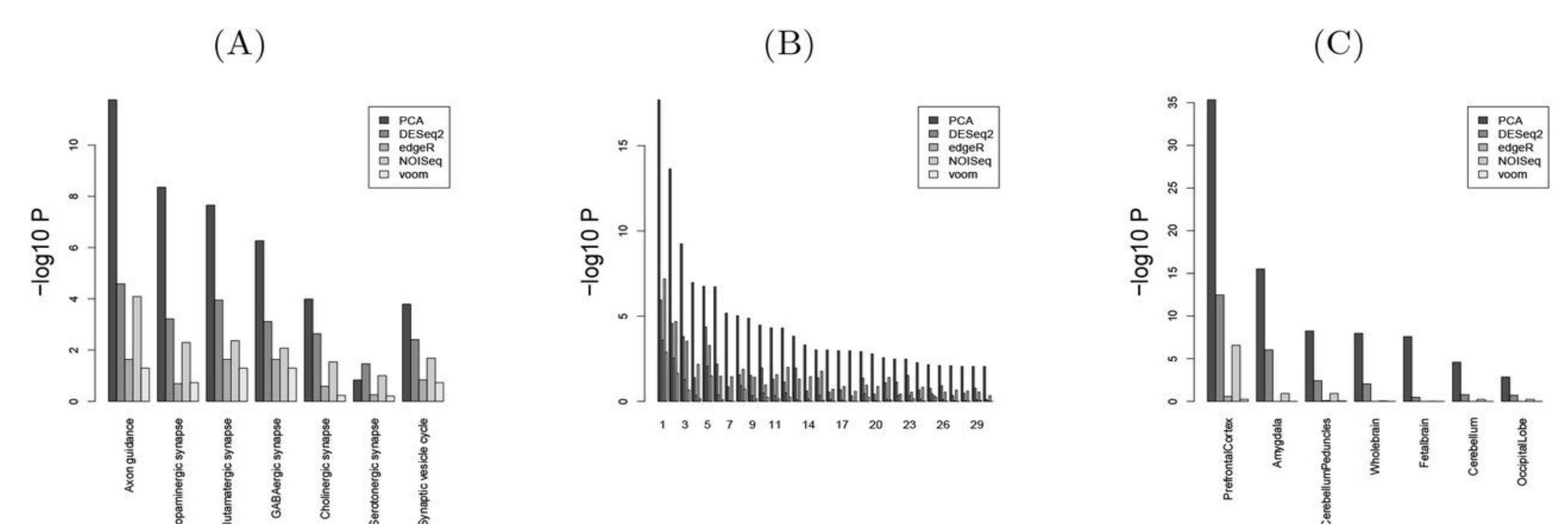
$$\sigma_n = \sqrt{\frac{\sum_{n < n_0} (h_n - \langle h_n \rangle)^2}{N(n < n_0)}}$$

補正 P 値  $P_i < 0.1$  の遺伝子を選択

生物学的な評価: PCA 対 DESeq2



(A) KEGG (B) GO BP (C) Human gene atlas  
 PCA vs DESeq2 vs edgeR vs NOISeq vs voom



手法を全く変えずに DNA のメチル化やヒストン修飾の  
 差分解析にも使えることが解った (文献2, 3).  
 DEG やメチル化、ヒストン修飾解析のデファクトスタン  
 ダードになるように努力したい。

演題発表内容に関連し、発表者に開示すべき  
 COI 関係にある企業などはありません。