

x_{ij} expression, gene i , sample j $i: 1 \sim p$ $j: 1 \sim n$.

class: $1 \sim K$.

C_k k -th class, which has n_k samples. C_k 是类 k .

$\bar{x}_{ik} = \sum_{j \in C_k} \frac{x_{ij}}{n_k}$ gene i of the centroid for class k .

$\bar{x}_i = \sum_{j=1}^n \frac{x_{ij}}{n}$ overall centroid, gene i .

shrink class centroid $\bar{x}_{ik} \rightarrow \bar{x}_i$

Standardization for gene i .

$d_{ik} = \frac{\bar{x}_{ik} - \bar{x}_i}{m_k (s_i + s_0)}$ $\bar{x}_{ik} - \bar{x}_i$: between class.

$s_i = \frac{1}{n-k} \sum_k \sum_{j \in C_k} (x_{ij} - \bar{x}_{ik})^2$: within class (pooled: sum of all class)
within class
pooled.

$m_k = \sqrt{\frac{1}{n_k} + \frac{1}{n}}$ $m_k \cdot s_i \sim (\bar{x}_{ik} - \bar{x}_i)$'s standard error.

So positive constant (for all genes)

防止 d_{ik} 过大 (防止低表达量 level's genes)

$s_0 = \text{median}(s_i)$

rewrite: $\bar{x}_{ik} = \bar{x}_i + m_k (s_i + s_0) d_{ik}$.

Solution: shrink each d_{ik} toward 0 (why?).

$\bar{x}'_{ik} = \bar{x}_i + m_k (s_i + s_0) d'_{ik}$

We call this soft thresholding. $d'_{ik} = \text{sign}(d_{ik}) (|d_{ik}| - \Delta)_+$ $\left. \begin{array}{l} t_+ = t \quad t > 0 \\ t_+ = 0 \quad t \leq 0 \end{array} \right\}$
when $\Delta \uparrow$, more gene could be eliminated from the class.

for a gene i , $d_{ik} \xrightarrow{\Delta \uparrow} 0$, the centroid of gene i is \bar{x}_i , that means gene i does not contribute to the nearest-centroid computation.

1st. Get gene expression profile of each test sample

2nd. 计算 Δ 到各个类重心的平方距离.

$S_i | z_i = k$ for class k .

$$J_k(x^*) = \sum_{i=1}^I \frac{(x_i^* - \bar{x}_{ik}')^2}{(s_i + s_0)^2} - 2 \log \pi_k.$$

$$\begin{cases} \bar{x}_{ik}' = \bar{x}_i + m_k (s_i + s_0) d_{ik} \\ d_{ik} = \text{sign}(d_{ik}) (|d_{ik}| - \Delta)_+ \end{cases}$$

π_k : class prior prob probability.

$$\sum_{k=1}^K \pi_k = 1. \quad \Leftrightarrow \quad \hat{\pi}_k = \frac{n_k}{n}. \quad \text{or} \quad \hat{\pi}_k = \frac{1}{K}$$

~~Then $C(x^*) = \arg \min_k J_k(x^*)$~~

$$C(x^*) = l \quad \text{where} \quad J_l(x^*) = \min_k J_k(x^*). \quad l \in 1 \sim K.$$

* * * estimates.

$$\hat{P}_k(x^*) = \frac{e^{-\frac{1}{2} J_k(x^*)}}{\sum_{l=1}^K e^{-\frac{1}{2} J_l(x^*)}}$$

Appendix.

~~$$J_k^{LDA}(x^*) = (x^* - \bar{x}_k) W^T (x^* - \bar{x}_k) - 2 \log \pi_k$$~~

$$J_k^{LDA}(x^*) = (x^* - \bar{x}_k) W^{*-1} (x^* - \bar{x}_k) - 2 \log \pi_k.$$