# Clustering 성능 지표

1) 정확한 값(cluster의 개수 및 소속)을 알고있는 경우.

| X\Y | $Y_1$ | $Y_2$ | $\cdots$ | $Y_s$ | Sums |
|-----|-------|-------|----------|-------|------|
| $X_1$ | $n_{11}$ | $n_{12}$ | $\cdots$ | $n_{1s}$ | $a_1$ |
| $X_2$ | $n_{21}$ | $n_{22}$ | $\cdots$ | $n_{2s}$ | $a_2$ |
| $\vdots$ | $\vdots$ | $\vdots$ | | $\vdots$ | $\vdots$ |
| $X_r$ | $n_{r1}$ | $n_{r2}$ | $\cdots$ | $n_{rs}$ | $a_r$ |
| Sums | $b_1$ | $b_2$ | $\cdots$ | $b_s$ | |

- Adjusted Rand Index

1) $n_{ij}$ : X clustering 에서는 cluster $X_i$ 있는데, Y clustering 에서는 같은 cluster $Y_j$의 수.   Data수

2) $a_i$ : X clustering 에서 같은 cluster $X_i$ 있는데, Y clustering 에서도 같은 cluster의 Data/설의 수

3) $b_j$ : X clustering 에서는 각각 cluster 있는데, Y cluster 에서도 다른 cluster $Y_j$로 표현되는 수

$$
\underset{(Adjusted Index)}{ARI} = \frac{\sum_{ij}\binom{n_{ij}}{2} - [\sum_i \binom{a_i}{2} \sum_j \binom{b_j}{2}]/\binom{n}{2}}{\underbrace{\frac{1}{2}[\sum_i \binom{a_i}{2} + \sum \binom{b_j}{2}]}_{Max~Index} - \underbrace{[\sum_i\binom{a_i}{2}\sum_j\binom{b_j}{2}]/\binom{n}{2}}_{Expected~Index}}
$$

- Adjusted Mutual Information.

$X = \{X_1, X_2, \cdots, X_r\}$

$Y = \{V_1, V_2, \cdots, V_s\}$

$P(i) = \frac{|X_i|}{N}$

$P(i,j) = \frac{|X_i \cap Y_j|}{N}$

$MI(X,Y) = \sum_{i=1}^{r}\sum_{j=1}^{s} P(i,j) \log \frac{P(i,j)}{P(i)P(j)}$