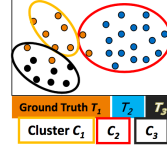


Entropy-Based Measures (II): Normalized Mutual Information (NMI)

□ Mutual information:

- Quantifies the amount of shared info between the clustering C and partitioning T

$$I(C, T) = \sum_{i=1}^r \sum_{j=1}^k p_{ij} \log \left(\frac{p_{ij}}{p_{C_i} \cdot p_{T_j}} \right)$$
- Measures the dependency between the observed joint probability p_{ij} of C and T , and the expected joint probability $p_{C_i} \cdot p_{T_j}$ under the independence assumption
- When C and T are independent, $p_{ij} = p_{C_i} \cdot p_{T_j}$, $I(C, T) = 0$. However, there is no upper bound on the mutual information



□ Normalized mutual information (NMI)

$$NMI(C, T) = \sqrt{\frac{I(C, T)}{H(C)} \cdot \frac{I(C, T)}{H(T)}} = \frac{I(C, T)}{\sqrt{H(C) \cdot H(T)}}$$

- Value range of NMI: $[0, 1]$. Value close to 1 indicates a good clustering

Entropy of clustering C :
$$H(C) = - \sum_{i=1}^r p_{C_i} \log p_{C_i}$$

Entropy of partitioning T :
$$H(T) = - \sum_{j=1}^k p_{T_j} \log p_{T_j}$$

Pairwise Measures: Four Possibilities for Truth Assignment

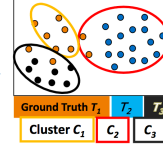
□ Four possibilities based on the agreement between cluster label and partition label

- TP : true positive—Two points x_i and x_j belong to the same partition T , and they also in the same cluster C

$$TP = |\{(x_i, x_j) : y_i = y_j \text{ and } \hat{y}_i = \hat{y}_j\}|$$

where y_i : the true partition label, and \hat{y}_i : the cluster label for point x_i

- FN : false negative: $FN = |\{(x_i, x_j) : y_i = y_j \text{ and } \hat{y}_i \neq \hat{y}_j\}|$
- FP : false positive: $FP = |\{(x_i, x_j) : y_i \neq y_j \text{ and } \hat{y}_i = \hat{y}_j\}|$
- TN : true negative: $TN = |\{(x_i, x_j) : y_i \neq y_j \text{ and } \hat{y}_i \neq \hat{y}_j\}|$



□ Calculate the four measures:

$$N = \binom{n}{2} \quad \text{Total \# of pairs of points}$$

$$TP = \sum_{i=1}^r \sum_{j=1}^k \binom{n_{ij}}{2} = \frac{1}{2} \left(\sum_{i=1}^r \sum_{j=1}^k n_{ij}^2 \right) - n$$

$$FN = \sum_{j=1}^k \binom{m_j}{2} - TP$$

$$FP = \sum_{i=1}^r \binom{n_i}{2} - TP \quad TN = N - (TP + FN + FP) = \frac{1}{2} \left(n^2 - \sum_{i=1}^r n_i^2 - \sum_{j=1}^k m_j^2 + \sum_{i=1}^r \sum_{j=1}^k n_{ij}^2 \right)$$

Pairwise Measures: Jaccard Coefficient and Rand Statistic

- Jaccard coefficient:** Fraction of true positive point pairs, but after ignoring the true negatives (thus asymmetric)
- $Jaccard = TP / (TP + FN + FP)$ [i.e., denominator ignores TN]
- Perfect clustering: $Jaccard = 1$

□ Rand Statistic:

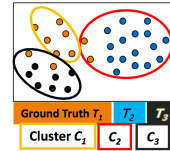
- $Rand = (TP + TN) / N$
- Symmetric; perfect clustering: $Rand = 1$

□ Fowlkes-Mallow Measure:

- Geometric mean of precision and recall

$$FM = \sqrt{prec \times recall} = \frac{TP}{\sqrt{(TP + FN)(TP + FP)}}$$

- Using the above formulas, one can calculate all the measures for the green table (leave as an exercise)



$C \backslash T$	T_1	T_2	T_3	Sum
C_1	0	20	30	50
C_2	0	20	5	25
C_3	25	0	0	25
m_j	25	40	35	100