

RAND INDEX

①

Given a set of n objects $S = \{o_1, o_2, o_3, \dots, o_n\}$

Suppose $U = \{u_1, u_2, u_3, \dots, u_R\}$ and $V = \{v_1, v_2, v_3, \dots, v_S\}$
 represents two different partitions of the objects
 such that $\bigcup_{i=1}^R u_i = S = \bigcup_{j=1}^C v_j$ and $U_i \cap U_{i'} = \emptyset$
 $= V_j \cap V_{j'} = \emptyset$

Suppose U is actual criteria.
 and V is clustering result.

Let

- a \Rightarrow The number of pairs of objects in the same class in U and in the same cluster in V .
- b \Rightarrow The number of pairs of objects in the same class in U but not in the same cluster V .
- c \Rightarrow The number of pairs of objects in the different class in U , but in same cluster in V .
- d \Rightarrow The number of pairs of objects in the different class in U , and in the different cluster in V .

U/V	Same cluster V	Different cluster V
Same class	A (TP)	C (FN)
Different class	B (FP)	D (TN)

- TP \Rightarrow Two similar documents (objects) assigned to same cluster.
 TN \Rightarrow Two dissimilar documents (objects) assigned to different clusters.
 FP \Rightarrow Two dissimilar documents (objects) assigned to same cluster.
 FN \Rightarrow Two similar documents (objects) assigned to different clusters.

Example:

class/cluster	v_1	v_2	...	v_c	SUM
u_1	n_{11}	n_{12}	...	n_{1c}	$n_{1..}$
u_2	n_{21}	n_{22}	...	n_{2c}	$n_{2..}$
\vdots	\vdots	\vdots			\vdots
u_R	n_{R1}	n_{R2}	...	n_{Rc}	$n_{R..}$
SUM	$n_{..1}$	$n_{..2}$...	$n_{..c}$	$n_{...} = n$

class/cluster	v_1	v_2	v_3	SUMS
u_1	1	1	0	2
u_2	1	2	1	4
u_3	0	0	4	4
SUMS	2	3	5	<u>$n=10$</u>

For the given example a is ②

$$a = \sum_{i,j} \binom{n_{ij}}{2} = \binom{2}{2} + \binom{4}{2} = \underline{\underline{7}}$$

$$b = \sum_i \binom{m_{i\cdot}}{2} - \sum_{i,j} \binom{n_{ij}}{2} = \binom{2}{2} + \binom{4}{2} + \binom{4}{2} - 7 = \underline{\underline{6}}$$

$$c = \sum_j \binom{n_{\cdot j}}{2} - \sum_{i,j} \binom{n_{ij}}{2} = \binom{2}{2} + \binom{3}{2} + \binom{5}{2} - 7 = \underline{\underline{7}}$$

$$d = \binom{n}{2} - (a+b+c)$$

$$= \binom{10}{2} - 7 - 6 - 7 = \underline{\underline{25}}$$

class/cluster	same	different
Same	a (TP) 7	b (FN) 7
different	c (FP) 6	d (TN) 25

$$\text{Rand Index} = \frac{a+d}{a+b+c+d}$$

$$P = \frac{A}{A+B}, \quad R = \frac{A}{A+C}$$

$$F-\text{ME} = \frac{2PR}{P+R}$$

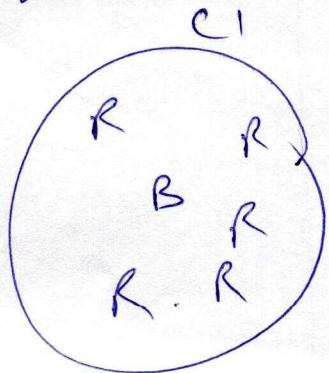
Purity \rightarrow maximum count of correctly assigned documents divided by total no. of documents.

$$\text{Purity}(A, C) = \frac{1}{N} \sum_k \max_j |w_k \cap c_j|$$

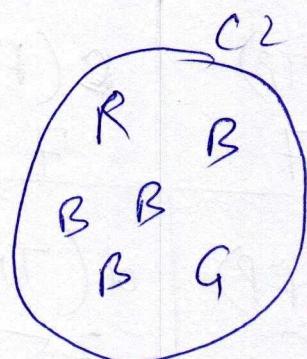
$$A = \{w_1, w_2, w_3, \dots, w_K\}$$

$$C = \{c_1, c_2, c_3, \dots, c_J\}.$$

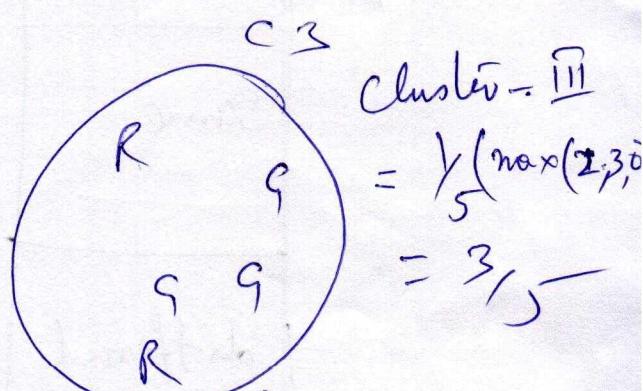
Value is in between 0 and 1.



Cluster - I
cluster - II



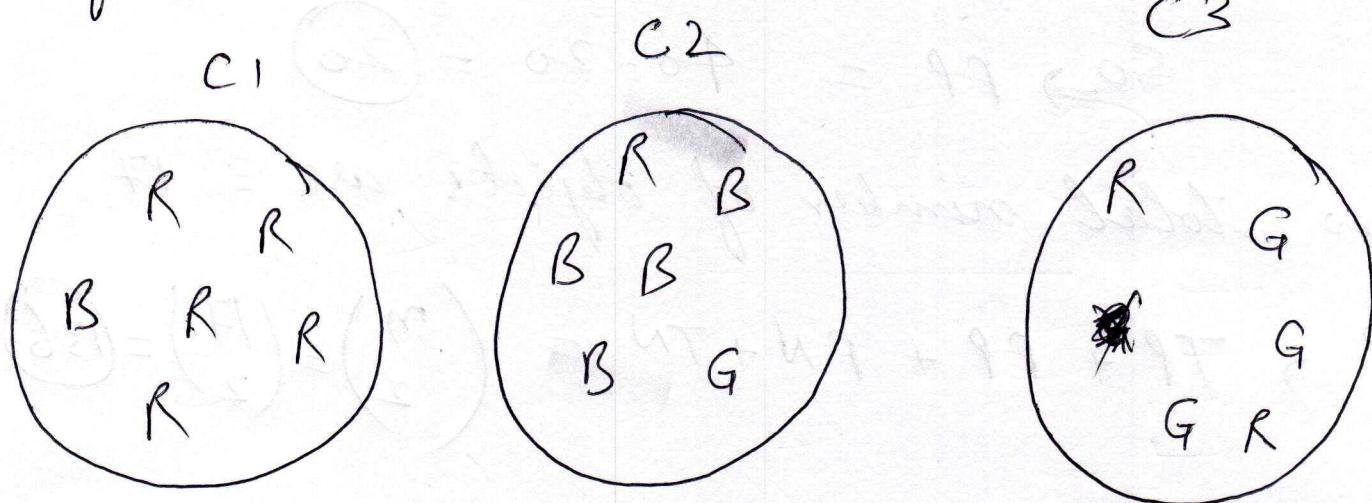
$$\text{Purity} = \frac{\max(5, 0, 1)}{6} = \frac{5}{6}$$



$$\begin{aligned} \text{Cluster - III} &= \frac{\max(2, 3, 0)}{5} \\ &= \frac{3}{5} \end{aligned}$$

(3)

Rand Index for the same example can be calculated as follows:



As per the definition of Rand index, we count the number of pairs of objects belong to same class and whether they are in the same cluster as well.

$\frac{TP + FP}{TP + FP}$. The three cluster contains 6, 6 and 5 objects so total number of pairs of object that are in the same cluster are:

$$TP + FP = \binom{6}{2} + \binom{6}{2} + \binom{5}{2} = 40$$

out of these, the R pairs in cluster (C1), the B pairs in cluster (C2) and R and G pairs in cluster (C3) are true positive.

$$\begin{aligned} TP &= \frac{\binom{5}{2}}{(C1)} + \frac{\binom{4}{2}}{(C2)} + \frac{\binom{3}{2}}{(C3)} + \frac{\binom{2}{2}}{(C3)} \\ &= 10 + 6 + 3 + 1 = 20 \end{aligned}$$

$$\text{Thus } TP + FP = 40$$

$$\Rightarrow TP = 20$$

$$\text{So } PP = 40 - 20 = 20$$

Now total number of objects are = 17

$$\text{So } \Rightarrow TP + PP + FN + TN = \binom{n}{2} = \binom{17}{2} = 136$$

FN \Rightarrow (False negative) is when the number of pair of objects in the different class but in same cluster.

$TN \Rightarrow$ (True negative) is when the number of pair of objects in the different class and in different clusters as well.

$$\text{So } TN = \binom{7}{2} + \binom{10}{2} + \binom{4}{2}$$

$$\binom{7}{2} \Rightarrow 4 'B' \text{ in } C_2 \text{ and } 3 'G' \text{ in } C_3 \\ \Rightarrow 4 + 3 = 7 \text{ (cluster 2)}$$

$$\binom{10}{2} \Rightarrow 5 'R' \text{ in } C_1 \text{ and } (2 'R' \text{ and } 3 'G' \text{ in } C_3) \\ \Rightarrow 5 + 5 = 10 \text{ (cluster 2)}$$

$$\binom{4}{2} \Rightarrow \text{Cluster 2 contains } 4 'B' \\ = 4 \text{ (cluster 3)}$$

$$TN = \binom{7}{2} + \binom{10}{2} + \binom{4}{2}$$

(4)

$$= \frac{7 \times 6}{2} + \frac{10 \times 9}{2} + \frac{4 \times 3}{2}$$

$$= 21 + 45 + 6 = 72$$

$$\boxed{TN = 72}$$

$$\text{Now } TP + FP + FN + TN = 136$$

$$FN = 136 - (TP + FP) - TN$$

$$= 136 - 40 - 72$$

$$\boxed{FN = 24}$$

So the confusion matrix is also known as Contingency table :

	TP (a) 20	FN (c) 24
	FP (b) 20	TN (d) 72

$$RA = \frac{a+d}{a+b+c+d} = \frac{92}{136} = 0.6764$$

$$P = \frac{a}{a+b}, R = \frac{a}{a+c}, F = \frac{2PR}{P+R}$$

So F-measure for the given example is

$$P = \frac{20}{40} = 0.5, R = \frac{20}{44} = 0.45$$

$$F = \frac{2 \times 0.5 \times 0.45}{0.5 + 0.45}$$

$$\boxed{F = 0.47}$$