

Lab Meeting 9/16/13

Nature Methods UCLA

digitally scanned light sheet

(original sheet done w/ elliptical lens
→ susceptible to higher order aberrations)

Rand derivation notes

Contingency table for clusterings

$$X = \{\{x_1, x_2\}, \{x_3, x_4, x_5\}, \dots, \{x_i\}\}$$

$$Y = \{\{x_1\}, \{x_2, x_3\}, \{x_4, x_5, x_6, x_7\}, \dots, \{x_i\}\}$$

X/Y	Y ₁	Y ₂	...	Y _j
x ₁	n ₁₁	n ₁₂	...	n _{1j}
x ₂	n ₂₁	n ₂₂		
...				
x _i	n _{i1}			n _{ij}

n_{ij} = # of elements
in common btw
cluster X_i and Y_j

Important identity: $\binom{N}{2} = \frac{N(N-1)}{2}$

~~See Fowlkes & Mallows:~~

~~$T_X \neq$~~

See Meilă:

N # point PAIRS in same cluster in $X + Y$
 N_{01} " " " in same cluster in X diff clusters in Y
 N_{10} " " " " " " in Y " " in X
 N_{00} " " " in diff clusters in X and Y

~~$N_{01} + N_{10} + N_{00} + N_{11} = \binom{N}{2}$~~
 where N is # of points (total # possible pairs)

$$\text{Rand}(X, Y) = \frac{\# \text{ pairs together in } X/Y \text{ or apart in } X/Y}{\# \text{ possible pairs}}$$

$$= \frac{N_{11} + N_{00}}{\binom{N}{2}}$$

See Fowlkes + Mallows:

$$T_K = 2 N_{11} = 2 \sum_i \sum_j \binom{n_{ij}}{2} \rightarrow \# \text{ of pairs together in } X \text{ and } Y$$

$$= 2 \sum_i \sum_j \frac{n_{ij}(n_{ij}-1)}{2} = \frac{1}{2} T_K$$

$$= \sum_i \sum_j n_{ij}^2 + \underbrace{\sum_i \sum_j n_{ij}}_{= N \text{ (# pts)}}$$

$$P_K = \cancel{P_K} = 2 \sum_i \binom{\sum_j n_{ij}}{2} \rightarrow \# \text{ of pairs in } Y$$

$$= \frac{1}{2} P_K$$

$$= \sum_i \left(\sum_j n_{ij} \right)^2 \neq N$$

$$Q_K = \cancel{Q_K} = 2 \sum_j \binom{\sum_i n_{ij}}{2} \neq N \quad \# \text{ of pairs in } X$$

$$= \frac{1}{2} Q_K$$

$$\# \text{ pairs in } Y + \# \text{ pairs in } X = 2 N_{11} + N_{01} + N_{10}$$

2 b/c Common pairs counted twice in sum

$$\Rightarrow \frac{1}{2} P_K + \frac{1}{2} Q_K = 2 N_{11} + N_{01} + N_{10} \quad 2 N_{11} = T_K$$

$$\textcircled{1} N_{01} + N_{10} = \frac{1}{2} [P_K + Q_K] - T_K$$

$$N_{11} + N_{01} + N_{10} + N_{00} = \binom{N}{2}$$

$$N_{11} + N_{00} = \binom{N}{2} - [N_{01} + N_{10}] \quad (2)$$

$$(1) + (2) \Rightarrow N_{11} + N_{00} = \binom{N}{2} - \left[\frac{1}{2}(P_K + Q_K) - T_K \right]$$

N 's cancel out

$$\Rightarrow \frac{N_{11} + N_{00}}{\binom{N}{2}} = \text{Rand}(X, Y)$$

$$= \frac{\binom{N}{2} - \frac{1}{2} \left[\sum_j \left(\sum_i n_{ij} \right)^2 + \sum_i \left(\sum_j n_{ij} \right)^2 \right] + \sum_i \sum_j n_{ij}}{\binom{N}{2}}$$

Ranges $\phi \rightarrow 1$, Adjusted Rand Index ^(normalized) _(after subtract expected value /)

$$= \frac{\text{Rand}(X, Y) - E[R]}{1 - E[R]}$$

Meila

\Rightarrow subtract expected value for randomized clusters and then normalize by range

So new range is still $\phi \rightarrow 1$ although some really ~~different~~ different clusterings can produce negative vals