

# A Probabilistic Hierarchical Clustering Algorithm

- For a set of objects partitioned into  $m$  clusters  $C_1, \dots, C_m$ , the quality can be measured by,

$$Q(\{C_1, \dots, C_m\}) = \prod_{i=1}^m P(C_i)$$

where  $P()$  is the maximum likelihood

- If we merge two clusters  $C_{j1}$  and  $C_{j2}$  into a cluster  $C_{j1} \cup C_{j2}$ , the change in quality of the overall clustering is

$$\begin{aligned} & Q((\{C_1, \dots, C_m\} - \{C_{j1}, C_{j2}\}) \cup \{C_{j1} \cup C_{j2}\}) - Q(\{C_1, \dots, C_m\}) \\ = & \frac{\prod_{i=1}^m P(C_i) \cdot P(C_{j1} \cup C_{j2})}{P(C_{j1})P(C_{j2})} - \prod_{i=1}^m P(C_i) \\ = & \prod_{i=1}^m P(C_i) \left( \frac{P(C_{j1} \cup C_{j2})}{P(C_{j1})P(C_{j2})} - 1 \right) \end{aligned}$$

- Distance between clusters  $C_1$  and  $C_2$ :

$$\text{dist}(C_i, C_j) = -\log \frac{P(C_1 \cup C_2)}{P(C_1)P(C_2)}$$

- If  $\text{dist}(C_i, C_j) < 0$ , merge  $C_i$  and  $C_j$

