# Entropy-Based Measures (I): Conditional Entropy

- **Entropy of clustering** $C$: $\quad H(\mathcal{C}) = -\sum_{i=1}^{r} p_{C_i} \log p_{C_i}$ $\quad p_{C_i} = \dfrac{n_i}{n}$ (i.e., the probability of cluster $C_i$)
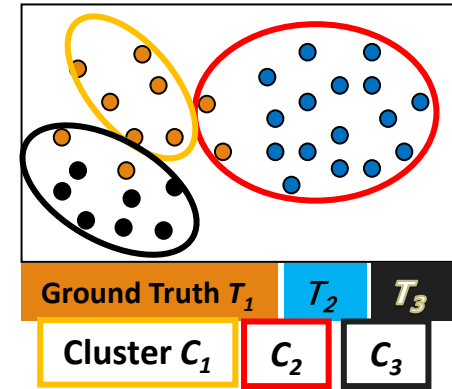
- **Entropy of partitioning** $T$: $\quad H(\mathcal{T}) = -\sum_{j=1}^{k} p_{T_i} \log p_{T_j}$

- **Entropy of $T$ with respect to cluster** $C_i$: $\quad H(\mathcal{T}|C_i) = -\sum_{j=1}^{k} (\dfrac{n_{ij}}{n_i}) \log(\dfrac{n_{ij}}{n_i})$

- **Conditional entropy of $T$ with respect to clustering** $C$: $\quad H(\mathcal{T}|\mathcal{C}) = -\sum_{i=1}^{r} (\dfrac{n_i}{n}) H(\mathcal{T}|C_i) = -\sum_{i=1}^{r}\sum_{j=1}^{k} p_{ij} \log(\dfrac{p_{ij}}{p_{C_i}})$



Ground Truth $T_1$ $T_2$ $T_3$

Cluster $C_1$ $C_2$ $C_3$

  - The more a cluster's members are split into different partitions, the higher the conditional entropy

  - For a perfect clustering, the conditional entropy value is 0, where the worst possible conditional entropy value is *log k*

$$H(\mathcal{T}|\mathcal{C}) = -\sum_{i=1}^{r}\sum_{j=1}^{k} p_{ij}(\log p_{ij} - \log p_{C_i}) = -\sum_{i=1}^{r}\sum_{j=1}^{k} p_{ij} \log p_{ij} + \sum_{i=1}^{r}(\log p_{C_i} \sum_{j=1}^{k} p_{ij})$$

$$= -\sum_{i=1}^{r}\sum_{j=1}^{k} p_{ij} \log p_{ij} + \sum_{i=1}^{r}(p_{C_i} \log p_{C_i}) = H(\mathcal{C}, \mathcal{T}) - H(\mathcal{C})$$