# K-Modes: Clustering Categorical Data

❑ *K-Means* cannot handle non-numerical (categorical) data

  ❑ Mapping categorical value to 1/0 cannot generate quality clusters for high-dimensional data

❑ **K-Modes**: An extension to *K-Means* by replacing means of clusters with **modes**

❑ Dissimilarity measure between object X and the center of a cluster Z

  ❑ $\Phi(x_j, z_j) = 1 - n_j^r/n_l$ when $x_j = z_j$ ; 1 when $x_j \ddagger z_j$

    ❑ where $z_j$ is the categorical value of attribute j in $Z_l$, $n_l$ is the number of objects in cluster *l*, and $n_j^r$ is the number of objects whose attribute value is r

❑ This dissimilarity measure (distance function) is **frequency-based**

❑ Algorithm is still based on iterative *object cluster assignment* and *centroid update*

❑ A ***fuzzy K-Modes*** method is proposed to calculate a ***fuzzy cluster membership value*** for each object to each cluster

❑ A mixture of categorical and numerical data: Using a **K-Prototype** method