# Cluster Stability

# Cluster Stability



- ❑ Clusterings obtained from several datasets sampled from the same underlying distribution as $D$ should be similar or "stable"
- ❑ Typical approach:
  - ❑ Find good parameter values for a given clustering algorithm
- ❑ Example: Find a good value of $k$, the correct number of clusters
- ❑ A **bootstrapping approach** to find the best value of $k$ (judged on stability)
  - ❑ Generate $t$ samples of size $n$ by sampling from $D$ with replacement
  - ❑ For each sample $D_i$, run the same clustering algorithm with $k$ values from 2 to $k_{max}$
  - ❑ Compare the distance between all pairs of clusterings $C_k(D_i)$ and $C_k(D_j)$ via some distance function
    - ❑ Compute the expected pairwise distance for each value of $k$
  - ❑ The value $k^*$ that exhibits the least deviation between the clusterings obtained from the resampled datasets is the best choice for $k$ since it exhibits the most stability
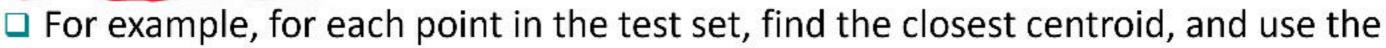
# Other Methods for Finding K, the Number of Clusters

- ❑ **Empirical method**
  - ❑ # of clusters: $k \approx \sqrt{n/2}$ for a dataset of n points (e.g., $n = 200$, $k = 10$)
- ❑ **Elbow method**: Use the turning point in the curve of the sum of within cluster variance with respect to the # of clusters



Elbow for KMeans clustering

- ❑ **Cross validation method**
  - ❑ Divide a given data set into $m$ parts
  - ❑ Use $m - 1$ parts to obtain a clustering model
  - ❑ Use the remaining part to test the quality of the clustering
    - ❑ For example, for each point in the test set, find the closest centroid, and use the sum of squared distance between all points in the test set and the closest centroids to measure how well the model fits the test set
  - ❑ For any $k > 0$, repeat it $m$ times, compare the overall quality measure w.r.t. different $k$'s, and find # of clusters that fits the data the best

3