

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

□ 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

