# Basic Concepts: Measuring Similarity between Objects

# What Is Good Clustering?

❑ A good clustering method will produce high quality clusters which should have

   ❑ **High intra-class similarity:** Cohesive within clusters

   ❑ **Low inter-class similarity:** Distinctive between clusters

❑ **Quality function**

   ❑ There is usually a separate "quality" function that measures the "goodness" of a cluster

   ❑ It is hard to define "similar enough" or "good enough"

      ❑ The answer is typically highly subjective

❑ There exist many similarity measures and/or functions for different applications

❑ Similarity measure is critical for cluster analysis

# Similarity, Dissimilarity, and Proximity

❑ **Similarity measure** or **similarity function**

  ❑ A real-valued function that quantifies the similarity between two objects

  ❑ Measure how two data objects are alike: The higher value, the more alike

  ❑ Often falls in the range [0,1]:  0: no similarity; 1: completely similar

❑ **Dissimilarity** (or **distance**) measure

  ❑ Numerical measure of how different two data objects are

  ❑ In some sense, the inverse of similarity:  The lower, the more alike

  ❑ Minimum dissimilarity is often 0 (i.e., completely similar)

  ❑ Range [0, 1] or [0, ∞) , depending on the definition

❑ **Proximity** usually refers to either similarity or dissimilarity