A reliable taxonomy is crucial for the assessment of biodiversity and for the comparison of habitats based on their species composition. Determining taxon boundaries is challenging in the case of organisms for which often only molecular data are available, such as bacteria, fungi, and many unicellular eukaryotes. Even in the case of organisms with well-established microscopical characteristics, molecular taxonomy is necessary to determine misidentified and mislabeled GenBank sequences, to identify incompletely known specimens and cryptic species, and last but not least to analyze sequences directly sampled from the environment as in metagenomics studies. For molecular taxonomy, researchers mostly use a predefined threshold for pair-wise genetic distances in clustering algorithms to assign sequences to molecular operational taxonomic units. However, thresholds applied differ in literature, even if applied to the same organisms and molecular markers,, and are often based on subjective criteria or just on tradition. Furthermore, the clustering algorithm applied also has a profound impact on the clustering outcome, but it is seldom addressed which algorithm is most appropriate for molecular taxonomy. Finally, the calculation of the distance matrices may also cause considerable methodological problems because of alignment ambiguity, rate heterogeneity between sites, and other potential sources of biases. To address these issues, we have designed and implemented a simple yet effective and flexible clustering optimization method. Using biologically sensible reference partitions, it automatically distinguishes between within-taxon and between-taxon sequence heterogeneity in the course of identifying optimal thresholds, clustering algorithms, and distance methods. Usage examples for clustering optimization with alternative types of biological data are provided, and it is discussed as a general method for improving molecular taxonomy.