1. Based on your experiments in Problems 2 and 3, what is the relation between the number of clusters and the quality of a clustering? What would a "good" number of clusters be for this particular data set?

Based on our experiments in Problems 2 and 3, we have seen that the quality of clustering improves as the number of clusters (k) increases.

For this particular data set, a "good" number of clusters is 7 that maximizes the Rand Index.

2. Why is it important to run an LDA model for multiple passes, and not just one? Why is it important to monitor an LDA model for convergence (like you did in Problem 5) and not simply run it for, say, 1000 passes?

Running an LDA model for multiple passes helps it learn better about topics in the data. It iteratively improves its understanding. Monitoring convergence is like making sure the learning process is stable and won't get much better with more tries. Doing too many passes, like 1000, might make the model too focused on the training data and not work well with new data (overfit). Checking convergence saves time and resources, ensuring the model are just right.

3. What are the differences between k-means and LDA? When would you use one, when the other?

K-means and LDA are distinct clustering algorithms with different approaches. K-means is an example of hard clustering which assigns each observation particularly to a single cluster. In contrast, LDA does soft clustering, assigning probabilities to indicate the likelihood of an observation belonging to various clusters. K-means is like sorting the observations into clusters, where each observation goes into one specific cluster, and the clusters don't overlap. It's great when you have clear group of observations. On the other hand, LDA is like recognizing that the observations can belong to more than one cluster, like a review being both about music and sound-bar. This is useful when topics can overlap, like in documents talking about different things.