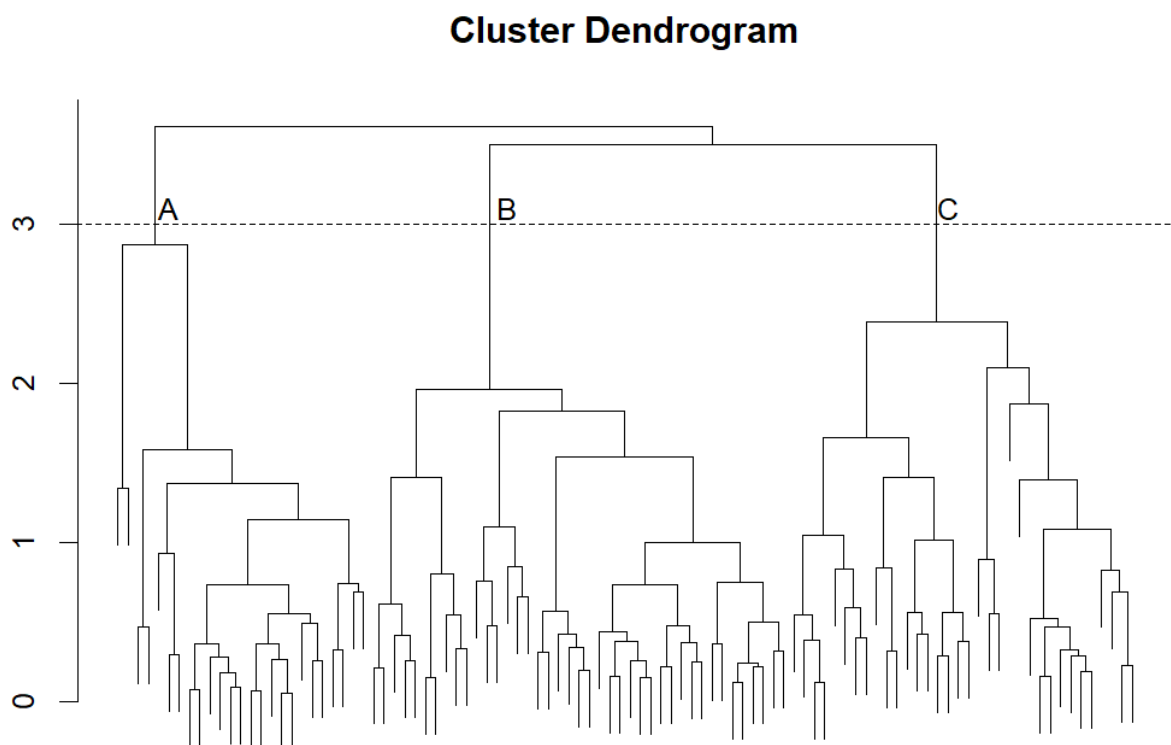Explain in your own words how much of the variance is explained by a PCA retaining p components in a dataset containing p variables?

You have performed hierarchical clustering on a dataset with Euclidean distance and the average linkage method. You obtain the following dendrogram, which you decide to cut at height = 3 (horizontal dashed line).

**Cluster Dendrogram**



You label the clustered observations belonging to the leftmost branch A, the middle branch B and the rightmost branch C. Which clusters are, on average, most similar, as measured by the Euclidian distance?

You have a corpus of 20 000 tweets about a governmental election in the Netherlands, with a total number of 50 000 unique words. Your goal is to classify sentiment orientation (positive, negative) of each tweet. The first step is to create a vector representation of your text data using one of the two methods:

- Bag-of-Words (e.g., tf-idf)
- Word Embedding (e.g., word2vec)

Which method would be better in terms of memory allocation and capturing words relations? Explain your reasoning.

Text preprocessing is an important step for text mining. In your own words explain the purpose of text preprocessing.