

Answer 3. (a)

Now according to the Heaps' Law:

$$v = k \cdot n^{\beta} \text{-----}(1)$$

Where $\beta = 0.5$, v is vocabulary size (number of unique words), n is the number of words in corpus,
 k, β are parameters that vary for each corpus (typical values given are $10 \leq k \leq 100$ and $\beta = 0.5$)

For full 100% of the vocabulary

So we get the equation as below for 100% vocabulary, $\beta = 0.5$, n_1 is the number of words in corpus

$$100v = k \cdot (n_1)^{0.5} \text{-----}(2)$$

Now the equation for the 90% of vocabulary, $\beta = 0.5$, n_2 is the number of words in corpus

$$90v = k \cdot (n_2)^{0.5} \text{-----}(3)$$

To find the proportion of a collection of text before 90% of its vocabulary has been encountered we have to divide equation (3) by equation (2) and calculate the n_2/n_1

After dividing equation-by-equation (2) we get the following

$$\frac{9}{10} = \left(\frac{n_2}{n_1}\right)^{0.5}$$

$$(0.9)^2 = \left(\frac{n_2}{n_1}\right)$$

Proportion of a collection of text must be read before 90% of the vocabulary has been encountered= 0.81