We have $V = \{v_k\}_k$ as the vocabulary and the unigram model $p_v = \{p_k\}_{k=0}^{|V|-1}$

Prove that: $p_k = \dfrac{n_k}{\sum\limits_{k=0}^{|V|-1} n_k}$ is optimal, where $n_k$: number of occurrence of $v_k$ observed.

Proof :

For any model $p = \{p_k\}_{k=0}^{|V|-1}$ used to model the vocabulary $V$

We know that $\sum\limits_{k=0}^{|V|-1} p_k = 1$ ——— ①

must be satisfied.

We also have that the probability of the observations is:

$$p_{obs} = \prod_{k=0}^{|V|-1} (p_k)^{n_k} \quad\quad — ②$$

The most optimal model is one that maximizes the probability of observed data

That is, the model $p = \{p_k\}_{k=0}^{|V|-1}$ where the assignment of probabilities $p_k$ satisfies eq ① and maximizes eq ②.

$$\max_{P} \prod_{k=0}^{|V|-1} (p_k)^{n_k} \quad\quad s.t \quad p = \{p_k\}_{k=0}^{|V|-1} \text{ and } \sum_{k=0}^{|V|-1} p_k = 1$$

This is equivalent to maximizing the log of $p_{obs}$, since log is a monotonically increasing function

$$\approx \max_{P} \log\left[ \prod_{k=0}^{|V|-1} (p_k)^{n_k} \right] \quad\quad s.t \quad p = \{p_k\}_{k=0}^{|V|-1} \text{ and } \sum_{k=0}^{|V|-1} p_k = 1$$

$$= \max_{p} \sum_{k=0}^{|V|-1} n_k \log \beta_k \qquad \text{s.t} \qquad \sum_{k=0}^{|V|-1} \beta_k = 1$$

Using the lagrangian method to solve the constrained optimization problem

$$L(\beta_0, \beta_1 \ldots, \beta_{|V|-1}, \lambda) = \sum_{k=0}^{|V|-1} n_k \log \beta_k - \lambda \left[ \left( \sum_{k=0}^{|V|-1} \beta_k \right) - 1 \right]$$

$$\frac{\partial L(\beta_0, \beta_1, \ldots, \beta_{|V|-1}, \lambda)}{\partial p_k} = 0$$

$$\Rightarrow \quad \frac{n_k}{\beta_k} - \lambda = 0$$

$$\Rightarrow \quad \boxed{\beta_k = \frac{n_k}{\lambda}} \quad \text{———} \quad ③$$

From ①

$$\sum_{k=0}^{|V|-1} \beta_k = 1$$

$$\Rightarrow \quad \sum_{k=0}^{|V|-1} \frac{n_k}{\lambda} = 1$$

$$\Rightarrow \quad \boxed{\lambda = \sum_{k=0}^{|V|-1} n_k} \quad \text{———} \quad ④$$

Therefore, $\beta_k = \dfrac{n_k}{\sum_{k=0}^{|V|-1} n_k}$ maximizes $\beta_{obs}$ while satisfying $\sum_{k=0}^{|V|-1} \beta_k = 1$