

①

A probability distribution should satisfy these conditions below:

① The produce probability must be between $[0, 1]$:

Ans: $P_{\text{Laplace}}(w) = \frac{\text{Count}_w + 1}{\sum_w \text{Count}_w + |V|} \rightarrow \text{size of vocabulary}$

i) we know $\text{Count}_w + 1 \geq 1$, so $P_{\text{Laplace}}(w) \geq 0$

in case of $\text{Count}_w = 0 \rightarrow P(w) = \frac{0+1}{\sum_w \text{Count}_w + |V|} > 0$

ii) we know $\text{Count}_w \leq \sum_w \text{Count}_w \Rightarrow P(w) \leq 1$

in case of $\text{Count}_w = \sum_w \text{Count}_w$ if $|V| = 1 \rightarrow P(w) = 1$
else $P(w) < 1$

② Sum of probabilities must be 1:

$$\sum_w P(w) = \frac{\sum_w (\text{Count}_w + 1)}{\sum_w \text{Count}_w + |V|} = \frac{\sum_w \text{Count}_w + \sum_w 1}{\sum_w \text{Count}_w + |V|} \xrightarrow{\text{equal to number of vocab}} \frac{\sum_w \text{Count}_w + |V|}{\sum_w \text{Count}_w + |V|} = 1$$

②

likelihood of a corpus :

$$P(w_1) \times P(w_2|w_1) \times P(w_3|w_1, w_2) \times \dots \times P(w_n|w_1, \dots, w_{n-1})$$

we can replace $P(w_k|w_1, \dots, w_{k-1})$ with bigram probability $P(w_k|w_{k-1})$ then:

$$L(w_1, w_2, w_3, \dots, w_n | \theta) = P(w_1) \times P(w_2|w_1) \times P(w_3|w_2) \times \dots \times P(w_n|w_{n-1})$$

$$\theta = \{P(w_i), P(w_i|w_j) : w_i \in V, i=1, 2, 3, \dots, n, i \neq j\}$$

* $P(w_1)$ can be probability of $\langle \text{start} \rangle$ ($P(\text{start}) \leq 1$) or $P(w_1 | \langle \text{start} \rangle)$.
in the first case last probability would be $P(\langle \text{Eos} \rangle | w_n)$. In second case we finish with $P(w_n | w_{n-1})$.

$$** P(w_i | w_{i-1}) = \frac{\text{Count}_{w_i, w_{i-1}} + 1}{\text{Count}_{w_{i-1}} + |V|}$$

③

$$f(x, y) = x^2 - y^2$$

$$\nabla_{(x,y)} f(x, y) = \begin{bmatrix} \frac{\partial f(x,y)}{\partial x} \\ \frac{\partial f(x,y)}{\partial y} \end{bmatrix} = \begin{bmatrix} 2x \\ -2y \end{bmatrix}$$

④ $\begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} 2 \\ 1 \end{bmatrix}, \eta = 0.1$

$$x_{\text{new}} = x_{\text{old}} - \eta \nabla_x f(x, y) = 2 - 0.1(2x) \Big|_{x=2} = 2 - 0.1 \times 4 = 1.6$$

$$y_{\text{new}} = y_{\text{old}} - \eta \nabla_y f(x, y) = 1 - 0.1 \times (-2y) \Big|_{y=1} = 1 + 0.2 = 1.2$$

$$f_{\text{old}}(x, y) = x^2 - y^2 \Big|_{\substack{x=2 \\ y=1}} = 4 - 1 = 3$$

$$f_{\text{new}}(x, y) = x^2 - y^2 \Big|_{\substack{x=1.6 \\ y=1.2}} = (1.6)^2 - (1.2)^2 = 2.56 - 1.44 = 1.12$$

$f_{\text{old}} > f_{\text{new}} \rightarrow$ we know that $\begin{matrix} x^* \\ y^* \end{matrix} = \begin{bmatrix} 0 \\ \max y \end{bmatrix}$ and gradient descent will merge to those values eventually by decreasing x , and $f(x, y)$ and also increasing y