# A3Q4

phrase = 'to be or not to be'

alphabet = 26
space char = 1 } #|γ| = ⓔ7

① No smoothing probabilities.

| X | #X | | P(x) |
|---|---|---|---|
| t | ΙΙΙ = 3 | | 3/18 |
| o | ΙΙΙΙ = 4 | | 4/18 |
| b | ΙΙ = 2 | | 2/18 |
| e | ΙΙ = 2 | | 2/18 |
| r | Ι = 1 | | 1/18 |
| n | Ι = 1 | | 1/18 |
| ⓔ space character | ΙΙΙΙΙ = 5 | | 5/18 |
| * anything else | 0 = 0 | | 0 |

② **probabilities after Laplace smoothing**

• Laplace smoothing is to add 1 to all events in beginning.

| x | #(x) | p(x) |
|---|---|---|
| t | I I I I = 4 | 4/45 |
| o | I I I I I = 5 | 5/45 |
| b | I I I = 3 | 3/45 |
| e | I I I = 3 | 3/45 |
| r | I I = 2 | 2/45 |
| n | I I = 2 | 2/45 |
| ( ) | I I I I I I = 6 | 6/45 |
| a | I = 1 | 1/45 |
| b | I = 1 | 1/45 |
| c | I = 1 | 1/45 |
| ⋮ | | |
| ⋮ | | |
| ⋮ | | |
| z | I = 1 | 1/45 |
| | 27 + 18 = 45 | |

## (3) Witten-Bell smoothing

So, let us say we start bound this sequence and for every new character we mark it.

$ε\text{(ee)},\text{beee}\text{ecee'cn}\text{ffeeeee}$

$*t*o*\text{?}*b*c\text{!?}o*r\text{'}*n.ot\text{?}to\text{?}bc$

So whenever we see new characters we mark those.

| X | # X | | P(x) |
|---|---|---|---|
| * | ‖‖‖‖‖ = 7 | | 7/25 |
| t | ‖‖ = 3 | | 3/25 |
| o | ‖‖‖ = 4 | | 4/25 |
| b | ‖ = 2 | | 2/25 |
| c | ‖ = 2 | | 2/25 |
| r | ‖ = 1 | | 1/25 |
| n | ‖ = 1 | | 1/25 |
| ( ) | ‖‖‖ = 5 | | 5/25 |
| | 25 | | |

$$P(*) = 7/25 \Big\}$$ So, this is probability of all unseen events

$$P(*) \begin{cases} a \\ b \\ c \\ \vdots \\ z \end{cases} = 7/25$$

So, probability of each unseen event will be

$$\frac{7/25}{(27-7)} = \frac{7/25}{20} = \frac{7}{25 \times 20}$$

$$\boxed{P(\text{unseen}) = \frac{7}{500}}$$

probability of single unseen event.