

Document Clustering Revisited

Which cluster does d belong to? \rightarrow Which θ_i was used to generate d ?

$d = x_1 x_2 \dots x_L$ where $x_i \in V$ 

$$\begin{aligned}\text{cluster}(d) &= \arg \max_i p(\theta_i | d) \\ &= \arg \max_i p(d | \theta_i) p(\theta_i) \\ &= \arg \max_i \prod_{j=1}^L p(x_j | \theta_i) p(\theta_i) \\ &= \arg \max_i \prod_{w \in V} p(w | \theta_i)^{c(w, d)} p(\theta_i)\end{aligned}$$

$$\begin{aligned}p(\theta_i | d) &= \frac{p(d | \theta_i) p(\theta_i)}{p(d)} \\ &= \frac{p(d | \theta_i) p(\theta_i)}{\sum_{j=1}^k p(d | \theta_j) p(\theta_j)}\end{aligned}$$

