

- Гипотеза условной независимости. Появление слов в документе d , относящихся к теме t , описывается общим для всей коллекции распределением $p(w | t)$ и не зависит от документа.
- Вероятностная модель. Согласно формуле полной вероятности и гипотезе условной независимости

$$p(w | d) = \sum p(t | d) p(w | t)$$

- Алгоритм 1:
- Вход: распределения $p(t | d)$, $p(w | t)$.
 - 1 для всех d
 - здать длину n документа d ;
 - 2 для всех $i=1 \dots n$
 - выбрать случайную тему t из распределения $p(t | d)$,
 - выбрать случайный термин w из распределения $p(w | t)$.
 - Добавить в выборку пару (d, w) . Тема забывается.
- Выход: выборка пар (d_i, w_i) , где $i=1, \dots, n$.