# Use PLSA for Text Mining

- PLSA would be able to generate
  - Topic coverage in each document: $p(\pi_d = j)$
  - Word distribution for each topic: $p(w|\theta_j)$
  - Topic assignment at the word level for each document
  - The number of topics must be given in advance
- These probabilities can be used in many different ways
  - $\theta_j$ naturally serves as a word cluster
  - $\pi_{d,j}$ can be used for document clustering $\quad j^* = \arg\max_j \pi_{d,j}$
  - Contextual text mining: Make these parameters conditioned on context, e.g.,
    - $p(\theta_j |time)$, from which we can compute/plot $p(time| \theta_j )$
    - $p(\theta_j |location)$, from which we can compute/plot $p(loc| \theta_j )$