

# Formal Definition of Topic Mining and Analysis

- Input
  - A **collection** of **N** text documents  **$C = \{d_1, \dots, d_N\}$**
  - **Number of topics:  $k$**
- Output
  - **$k$  topics:  $\{\theta_1, \dots, \theta_k\}$**
  - **Coverage of topics in each  $d_i$ :  $\{\pi_{i1}, \dots, \pi_{ik}\}$**
  - $\pi_{ij}$  = prob. of  $d_i$  covering topic  $\theta_j$

$$\sum_{j=1}^k \pi_{ij} = 1$$

**How to define  $\theta_i$  ?**