

# Formal Definition of Topic Mining and Analysis

- Input
  - A collection of **N** text documents **C** = {**d**<sub>1</sub>, ..., **d**<sub>N</sub>}
  - Number of topics: **k**
- Output
  - **k** topics: { **θ**<sub>1</sub>, ..., **θ**<sub>k</sub> }
  - Coverage of topics in each **d**<sub>i</sub>: { **π**<sub>i1</sub>, ..., **π**<sub>ik</sub> }
  - **π**<sub>ij</sub> = prob. of **d**<sub>i</sub> covering topic **θ**<sub>j</sub>

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

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