

Formal Definition of Topic Mining and Analysis

- Input
 - A collection of **N** text documents **C** = {**d**₁, ..., **d**_N}
 - Number of topics: **k**
- Output
 - **k** topics: {**θ**₁, ..., **θ**_k}
 - Coverage of topics in each **d**_i: {**π**_{i1}, ..., **π**_{ik}}
 - **π**_{ij} = prob. of **d**_i covering topic **θ**_j

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

How to define **θ**_i ?