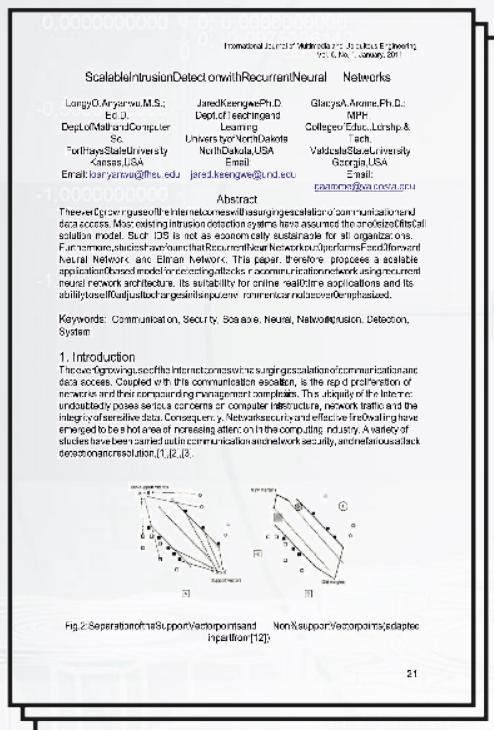


# Topic modeling: a way to navigate through text collections

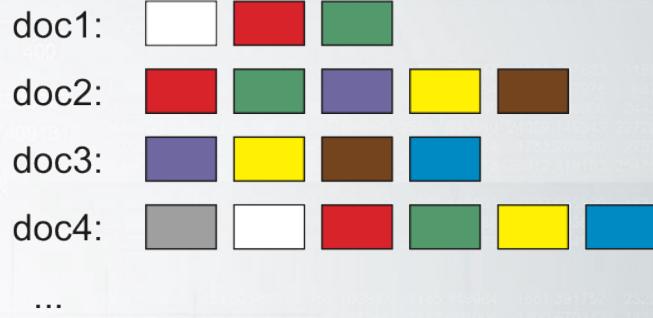


# From texts to topics



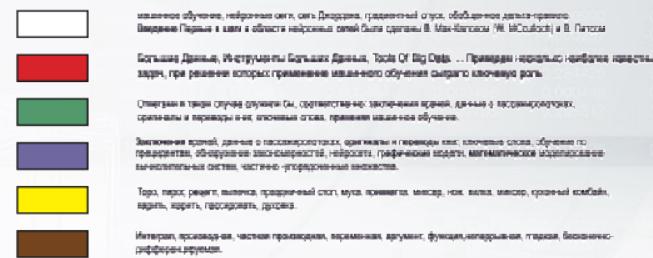
## Topics of documents

Documents



## Words and keyphrases of topics

Topics



# The formal task

**Given:**

- Collection of texts as bags-of-words:  
 $n_{wd}$  is a count of the word  $w$  in the document  $d$

**Find:**

- Probabilities of word in topics:  
 $\phi_{wt} = p(w|t)$
- Probabilities of topics in documents:  
 $\theta_{td} = p(t|d)$

# The formal task

Given:

- Collection of texts as bags-of-words:

$n_{wd}$  is a count of the word  $w$  in the document  $d$

Find:

- Probabilities of word in topics:

$$\phi_{wt} = p(w|t) \quad \text{Definition of a topic!}$$

- Probabilities of topics in documents:

$$\theta_{td} = p(t|d)$$

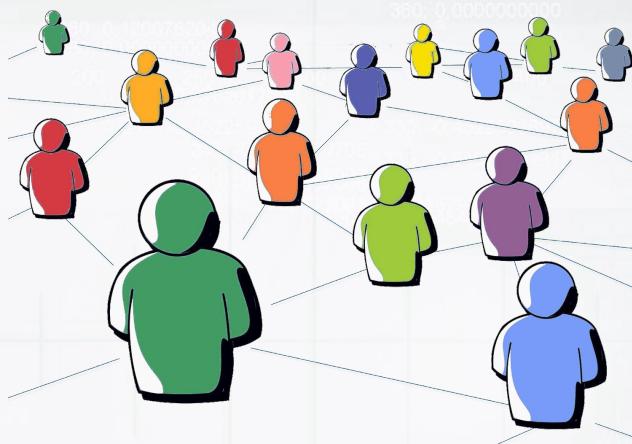
# Where do we need that?

Exploration and navigation through large text collections

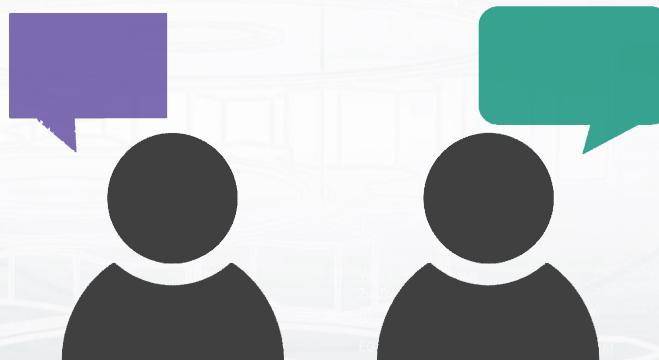


# Where do we need that?

- Social network analysis



- Dialogue manager in chat-bots



# Why do we need it?

Topic models provide hidden semantic representation of texts.

## Many more applications:

- Categorization and classification of texts
- Document segmentation and summarization
- News flows aggregation and analysis
- Recommender systems
- Image captioning
- Bioinformatics (genome annotation)
- Exploratory search
- ...

# Generative model of texts

## Probabilistic Latent Semantic Analysis (PLSA):

$$p(w|d) = \sum_{t \in T} p(w|t, d) p(t|d) = \sum_{t \in T} p(w|t) p(t|d)$$

### Notation:

- $w$  – word
- $d$  – document
- $t$  – topic

# Generative model of texts

## Probabilistic Latent Semantic Analysis (PLSA):

$$p(w|d) = \sum_{t \in T} p(w|t, d) p(t|d) = \sum_{t \in T} p(w|t) p(t|d)$$

*Law of total probability*

$$p(w) = \sum_{t \in T} p(w|t) p(t)$$

### Notation:

- $w$  – word
- $d$  – document
- $t$  – topic

# Generative model of texts

## Probabilistic Latent Semantic Analysis (PLSA):

$$p(w|d) = \sum_{t \in T} p(w|t, d) p(t|d) = \sum_{t \in T} p(w|t) p(t|d)$$

*Law of total probability*

$$p(w) = \sum_{t \in T} p(w|t) p(t)$$

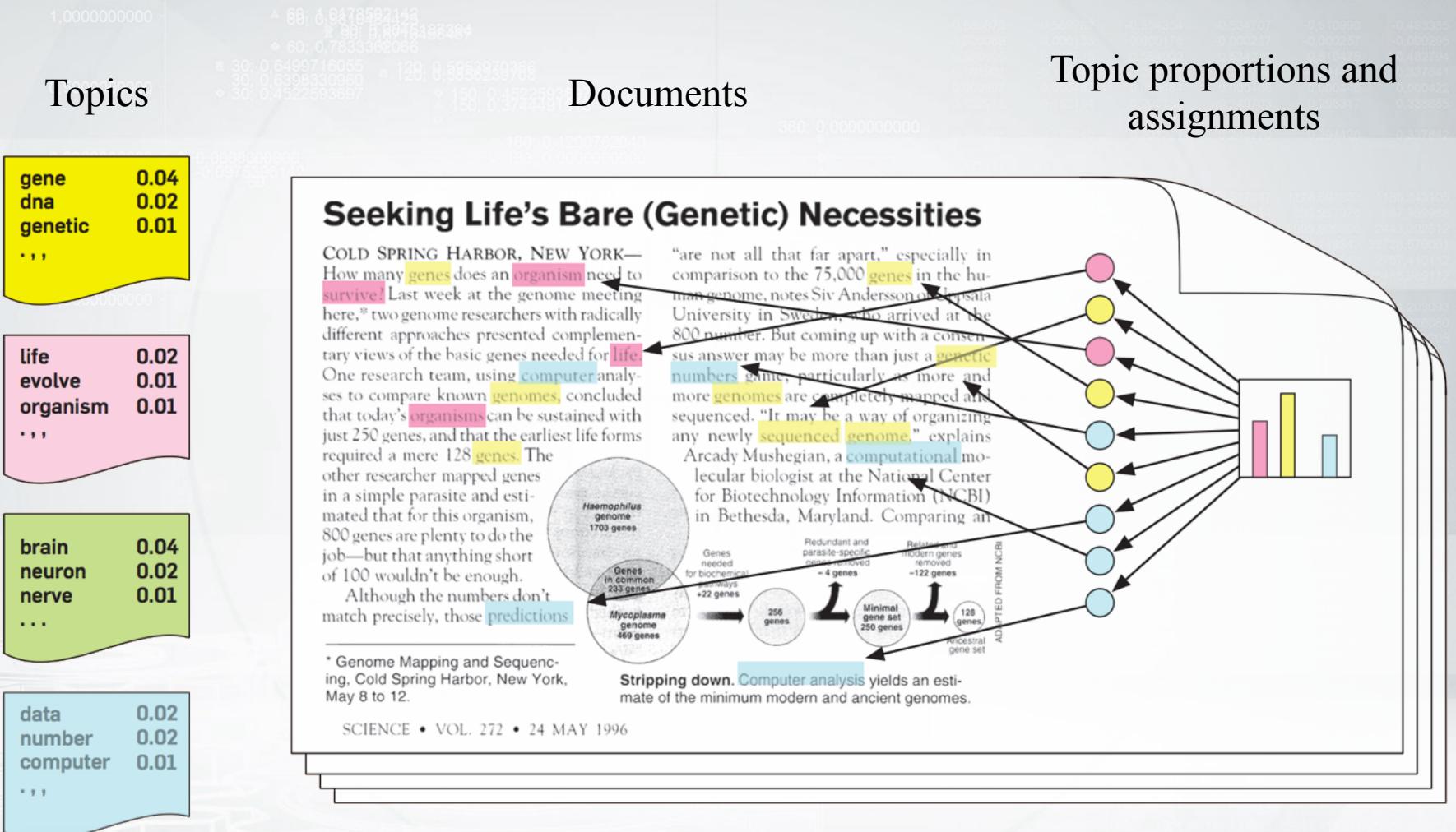
*Assumption of conditional independence*

$$p(w|t, d) = p(w|t)$$

### Notation:

- $w$  – word
- $d$  – document
- $t$  – topic

# Generative model of texts



David Blei, Probabilistic topic models, 2012.

# Matrix way of thinking

## Probabilistic Latent Semantic Analysis:

$$p(w|d) = \sum_{t \in T} p(w|t) p(t|d) = \sum_{t \in T} \phi_{wt} \theta_{td}$$

