

**APAN PS5430**

# **Applied Text & Natural Language Analytics**

## **Week 8: Topic Modeling I**

Javid Huseynov, Ph.D.  
Thursday, March 12, 2020



- Association Rule Mining
  - Definition
  - Approaches
  - Apriori Algorithm
- Topic Modeling
- Lexical Semantics Challenges
- Latent Semantic Analysis
- Latent Dirichlet Allocation
- Class Exercises

# Association Rule Mining (ARM)

- Task of finding frequent patterns, correlations, associations, or causal structures found in databases
- Based on “market-basket” analysis
- Given a set of transactions, ARM finds the rules which predict the occurrence of some items in the transaction based on the occurrence of other items:

$\{\text{bread}\} \Rightarrow \{\text{milk}\}$

$\{\text{soda}\} \Rightarrow \{\text{chips}\}$

$\{\text{bread}\} \Rightarrow \{\text{jam}\}$

TID	Items
1	Bread, Peanuts, Milk, Fruit, Jam
2	Bread, Jam, Soda, Chips, Milk, Fruit
3	Steak, Jam, Soda, Chips, Bread
4	Jam, Soda, Peanuts, Milk, Fruit
5	Jam, Soda, Chips, Milk, Bread
6	Fruit, Soda, Chips, Milk
7	Fruit, Soda, Peanuts, Milk
8	Fruit, Peanuts, Cheese, Yogurt

## Itemset

- Collection of items (e.g. {milk, bread})
- $k$ -itemset contains  $k$  items

## Support count ( $\sigma$ )

- Frequency of occurrence of an itemset
- $\sigma(\{\text{Milk, Bread}\}) = 3$

## Support ( $s$ )

- Fraction of transactions that contain both  $X$  and  $Y$ ,

$$\text{e.g. } s = \frac{\sigma(\{\text{Milk, Bread}\})}{\# \text{ of transactions}} = 0.38$$

## Confidence ( $c$ )

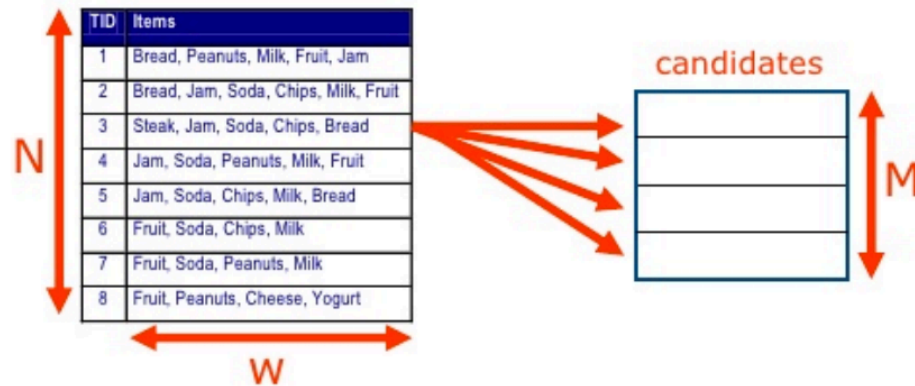
- How often items in  $Y$  appear in transactions with  $X$ ,

$$\text{e.g. } c = \frac{\sigma(\{\text{Milk, Bread}\})}{\sigma(\{\text{Bread}\})} = 0.75$$

TID	Items
1	Bread, Peanuts, Milk, Fruit, Jam
2	Bread, Jam, Soda, Chips, Milk, Fruit
3	Steak, Jam, Soda, Chips, Bread
4	Jam, Soda, Peanuts, Milk, Fruit
5	Jam, Soda, Chips, Milk, Bread
6	Fruit, Soda, Chips, Milk
7	Fruit, Soda, Peanuts, Milk
8	Fruit, Peanuts, Cheese, Yogurt

## Brute-force approach:

- ▶ Each itemset in the lattice is a **candidate** frequent itemset
- ▶ Count the support of each candidate by scanning the database



- ▶ Match each transaction against every candidate
- ▶ Complexity  $\sim O(NMw) \Rightarrow$  **Expensive** since  $M = 2^d$

Given a set of transactions  $T$ , find all rules having:

- Support  $\geq$  minimum support
- Confidence  $\geq$  minimum confidence

## Apriori Principle

- If an itemset is frequent, then its subsets must also be frequent
- Support of an itemset never exceeds the support of its subsets

$$\forall X, Y : (X \subseteq Y) \Rightarrow s(X) \geq s(Y)$$



# ARM: Apriori Algorithm



- ❑ Let  $k=1$
- ❑ Generate frequent itemsets of length 1
- ❑ Repeat until no new frequent itemsets are identified
  - ▶ Generate length  $(k+1)$  candidate itemsets from length  $k$  frequent itemsets
  - ▶ Prune candidate itemsets containing subsets of length  $k$  that are infrequent
  - ▶ Count the support of each candidate by scanning the DB
  - ▶ Eliminate candidates that are infrequent, leaving only those that are frequent

$C_k$ : Candidate itemset of size  $k$

$L_k$ : frequent itemset of size  $k$

$L_1 = \{\text{frequent items}\};$

**for** ( $k = 1; L_k \neq \emptyset; k++$ ) **do begin**

$C_{k+1}$  = candidates generated from  $L_k$ ;

**for each** transaction  $t$  in database **do**

increment the count of all candidates in  $C_{k+1}$   
that are contained in  $t$

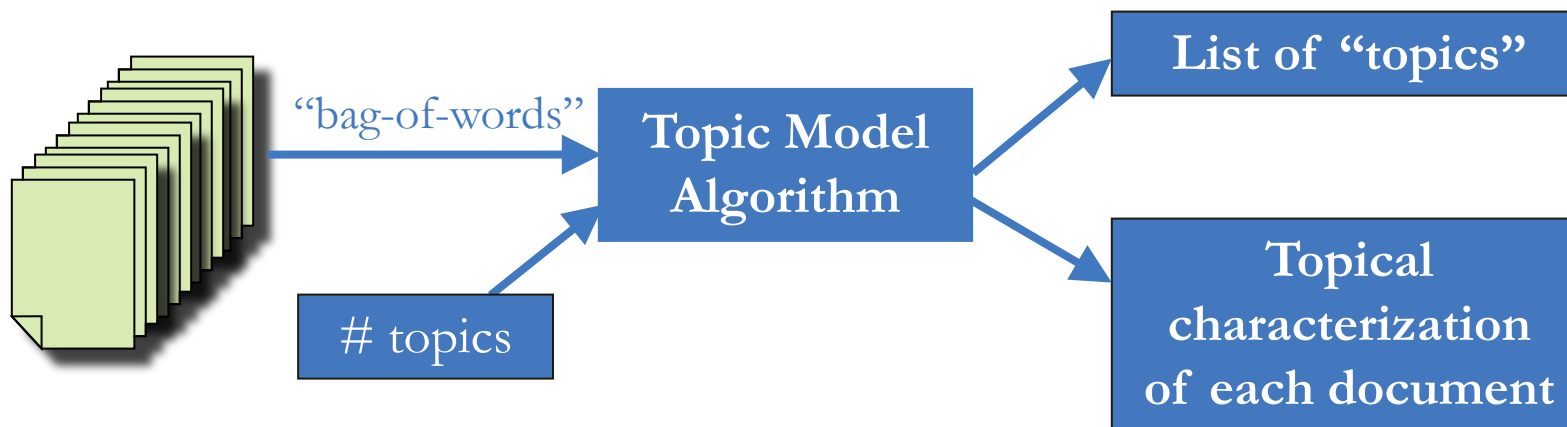
$L_{k+1}$  = candidates in  $C_{k+1}$  with min\_support

**end**

**return**  $\cup_k L_k$ ;

Unsupervised method to organize, understand, summarize a set of documents

- Discovering hidden topical patterns that are present across the collection
- Annotating documents according to these topics
- Using these annotations to organize, search and summarize texts



Useful applications:

- Search
- Automated tagging
- Summarization

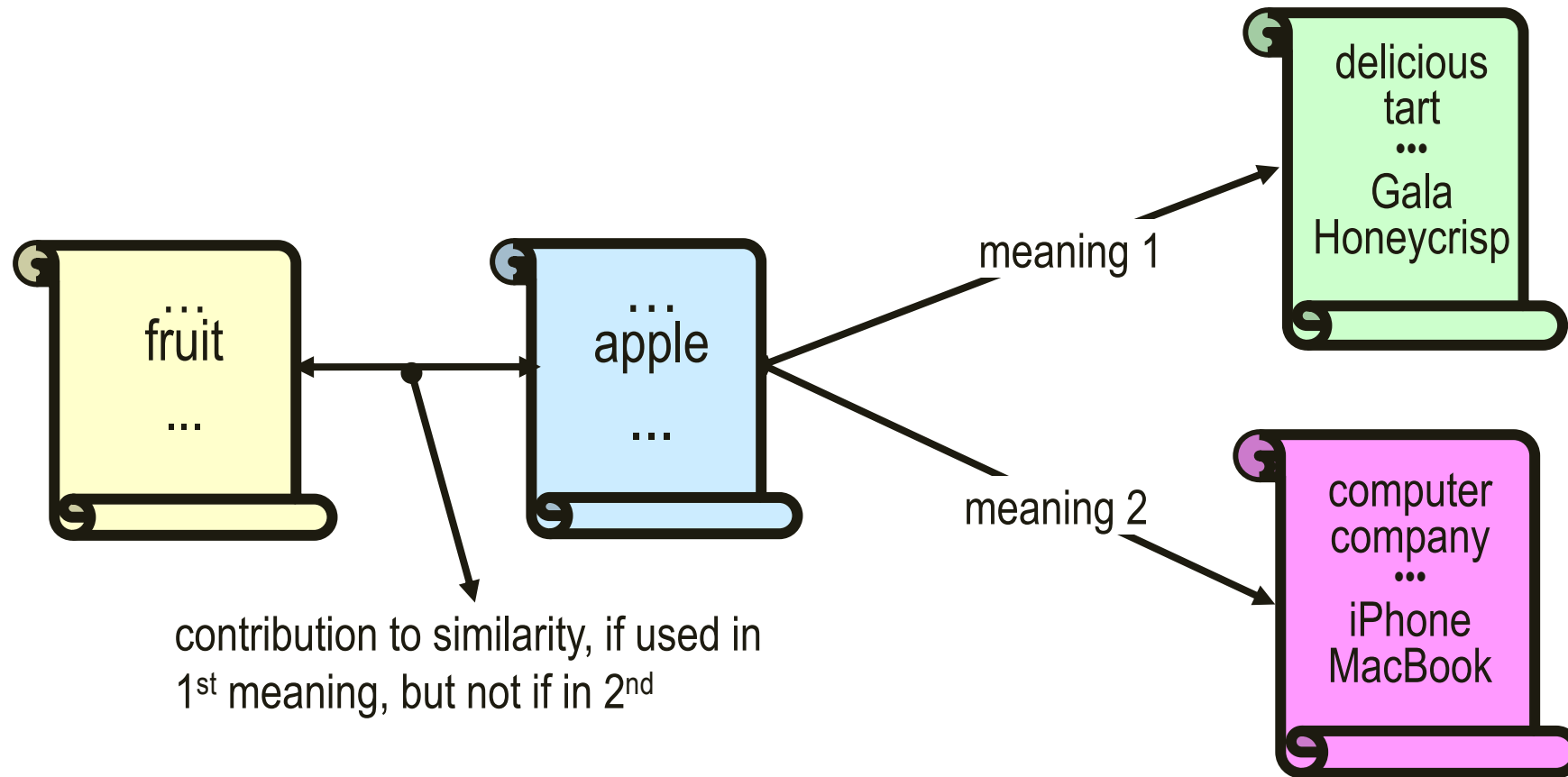
- **Polysemy (Ambiguity)**
- Same word may have different meanings when it appears in different contexts
- Vector-space model is unable to discriminate between these different meanings
- **Synonymy (Association)**
- Different words may have same meaning when they appear in a similar context
- No associations between words are captured in the vector-space representations

$$\text{sim}_{\text{true}}(d, q) < \cos(\angle(\vec{d}, \vec{q}))$$

$$\text{sim}_{\text{true}}(d, q) > \cos(\angle(\vec{d}, \vec{q}))$$

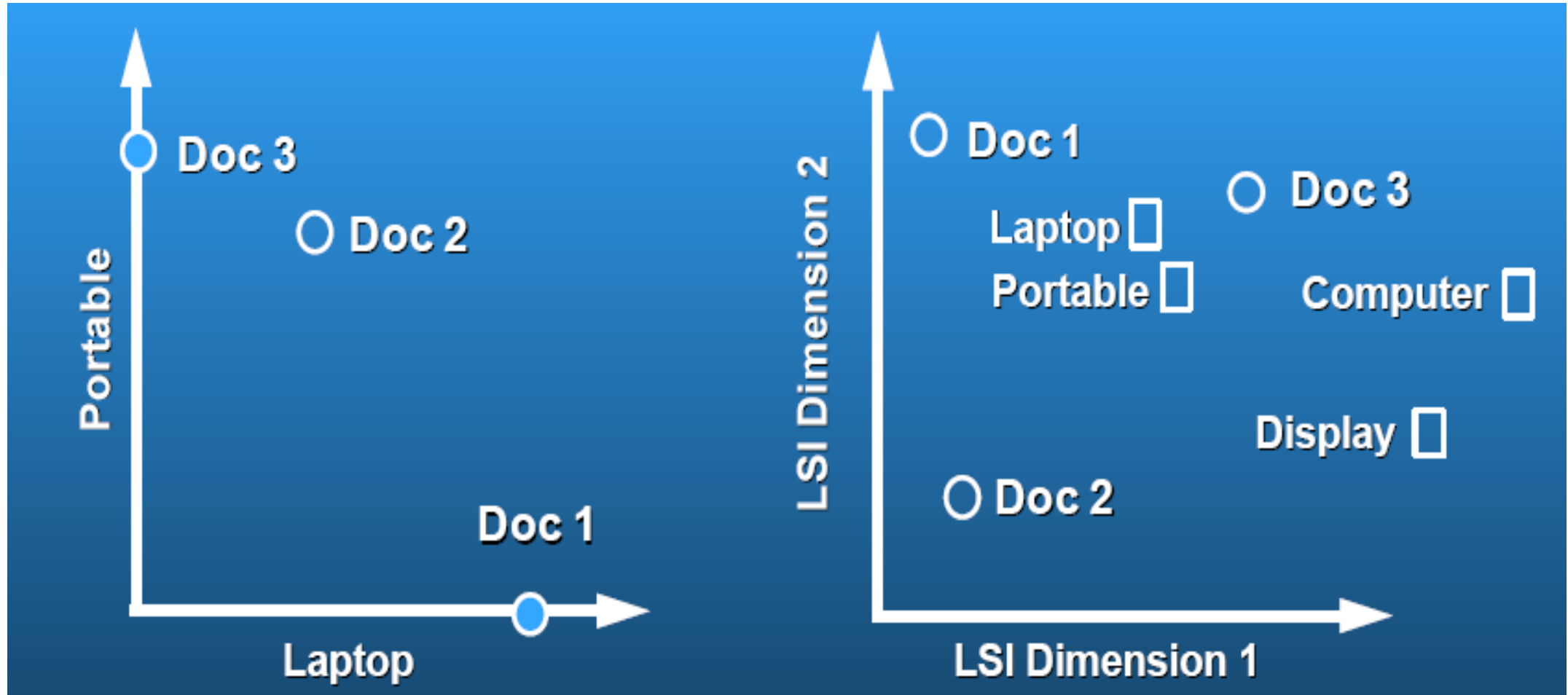


- Document similarity on a single word level: polysemy and context



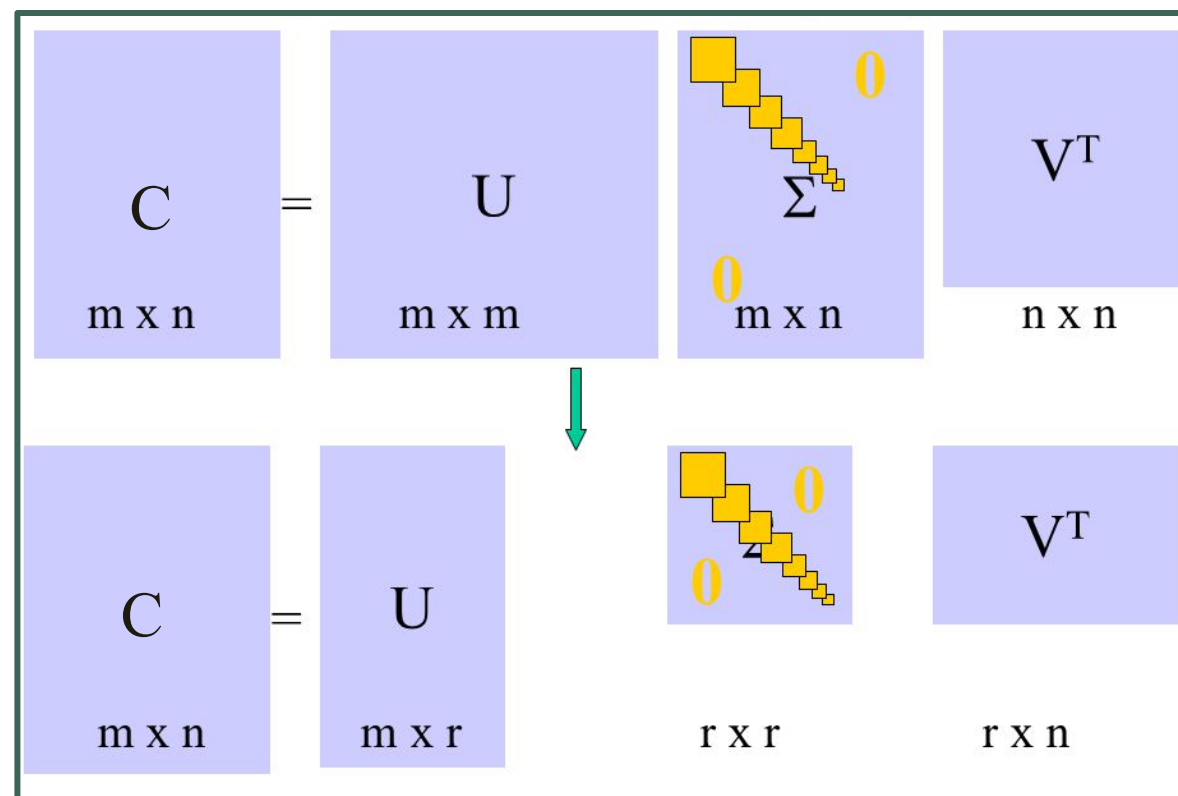
- Term-document (TFIDF) matrices can be large (sparse), while topic space is small
- LSI takes documents that are semantically similar (talk about the same topics), but are not similar in the vector space (because they use different words) and re-represents them in a reduced vector space in which they have higher similarity.
- Perform a **low-rank approximation** of term-document matrix
  - Map documents and terms to a low dimensional representation (sparsity reduction) using **Singular Value Decomposition** (SVD)
  - Design a mapping such that the low-dimensional or **latent semantic space** reflects semantic associations
  - Compute document similarity based on the inner product in this latent semantic space

# LSA: ABSTRACT Illustrative example



# Singular Value Decomposition (SVD)

- Any real  $m \times n$  matrix  $C$  can be decomposed into:
$$C = U\Sigma V^T$$
- $U$  is  $m \times m$ , column orthonormal ( $U^T U = I$ )
- $\Sigma$  is  $m \times n$  and diagonal:
  - $\Sigma = \text{diag}(\sigma_1, \sigma_2, \dots, \sigma_r)$
  - $\sigma_1$  are called *singular* values of  $C$
  - $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r \geq 0$
- $V$  is  $n \times n$  and orthonormal ( $V V^T = V^T V = I$ )



# SVD Example (LSA)

- Term-Document Matrix  $C = U\Sigma V^T$

$U$	1	2	3	4	5
ship	-0.44	-0.30	0.57	0.58	0.25
boat	-0.13	-0.33	-0.59	0.00	0.73
ocean	-0.48	-0.51	-0.37	0.00	-0.61
wood	-0.70	0.35	0.15	-0.58	0.16
tree	-0.26	0.65	-0.41	0.58	-0.09

$\Sigma$	1	2	3	4	5
1	2.16	0.00	0.00	0.00	0.00
2	0.00	1.59	0.00	0.00	0.00
3	0.00	0.00	1.28	0.00	0.00
4	0.00	0.00	0.00	1.00	0.00
5	0.00	0.00	0.00	0.00	0.39

$C$	$d_1$	$d_2$	$d_3$	$d_4$	$d_5$	$d_6$
ship	1	0	1	0	0	0
boat	0	1	0	0	0	0
ocean	1	1	0	0	0	0
wood	1	0	0	1	1	0
tree	0	0	0	1	0	1

$V^T$	$d_1$	$d_2$	$d_3$	$d_4$	$d_5$	$d_6$
1	-0.75	-0.28	-0.20	-0.45	-0.33	-0.12
2	-0.29	-0.53	-0.19	0.63	0.22	0.41
3	0.28	-0.75	0.45	-0.20	0.12	-0.33
4	0.00	0.00	0.58	0.00	-0.58	0.58
5	-0.53	0.29	0.63	0.19	0.41	-0.22

# SVD Example (LSA)



- Dimensionality reduction  $C_2 = U\Sigma_2V^T$

$U$	1	2	3	4	5
ship	-0.44	-0.30	0.00	0.00	0.00
boat	-0.13	-0.33	0.00	0.00	0.00
ocean	-0.48	-0.51	0.00	0.00	0.00
wood	-0.70	0.35	0.00	0.00	0.00
tree	-0.26	0.65	0.00	0.00	0.00

$\Sigma_2$	1	2	3	4	5
1	2.16	0.00	0.00	0.00	0.00
2	0.00	1.59	0.00	0.00	0.00
3	0.00	0.00	0.00	0.00	0.00
4	0.00	0.00	0.00	0.00	0.00
5	0.00	0.00	0.00	0.00	0.00

$V^T$	$d_1$	$d_2$	$d_3$	$d_4$	$d_5$	$d_6$
1	-0.75	-0.28	-0.20	-0.45	-0.33	-0.12
2	-0.29	-0.53	-0.19	0.63	0.22	0.41
3	0.00	0.00	0.00	0.00	0.00	0.00
4	0.00	0.00	0.00	0.00	0.00	0.00
5	0.00	0.00	0.00	0.00	0.00	0.00

$C_2$	$d_1$	$d_2$	$d_3$	$d_4$	$d_5$	$d_6$
ship	0.85	0.52	0.28	0.13	0.21	-0.08
boat	0.36	0.36	0.16	-0.20	-0.02	-0.18
ocean	1.01	0.72	0.36	-0.04	0.16	-0.21
wood	0.97	0.12	0.20	1.03	0.62	0.41
tree	0.12	-0.39	-0.08	0.90	0.41	0.49



# SVD Example (LSA)



- Similarity of  $d_2$  and  $d_3$  in the original space: 0.
- Similarity of  $d_2$  and  $d_3$  in the reduced space:  $0.52 * 0.28 + 0.36 * 0.16 + 0.72 * 0.36 + 0.12 * 0.20 + -0.39 * -0.08 \approx \mathbf{0.52}$

$C$	$d_1$	$d_2$	$d_3$	$d_4$	$d_5$	$d_6$
ship	1	0	1	0	0	0
boat	0	1	0	0	0	0
ocean	1	1	0	0	0	0
wood	1	0	0	1	1	0
tree	0	0	0	1	0	1

$C_2$	$d_1$	$d_2$	$d_3$	$d_4$	$d_5$	$d_6$
ship	0.85	0.52	0.28	0.13	0.21	-0.08
boat	0.36	0.36	0.16	-0.20	-0.02	-0.18
ocean	1.01	0.72	0.36	-0.04	0.16	-0.21
wood	0.97	0.12	0.20	1.03	0.62	0.41
tree	0.12	-0.39	-0.08	0.90	0.41	0.49

# Latent Dirichlet Allocation (LDA)



## ■ History:

- 1988: *Latent Semantic Analysis* (LSA)
  - Singular Value Decomposition (SVD) of word-document count matrix
- 1999: *Probabilistic Latent Semantic Analysis* (PLSA)
  - Non-negative matrix factorization (NMF)
- 2003: *Latent Dirichlet Allocation* (LDA)
  - Bayesian version of PLSA

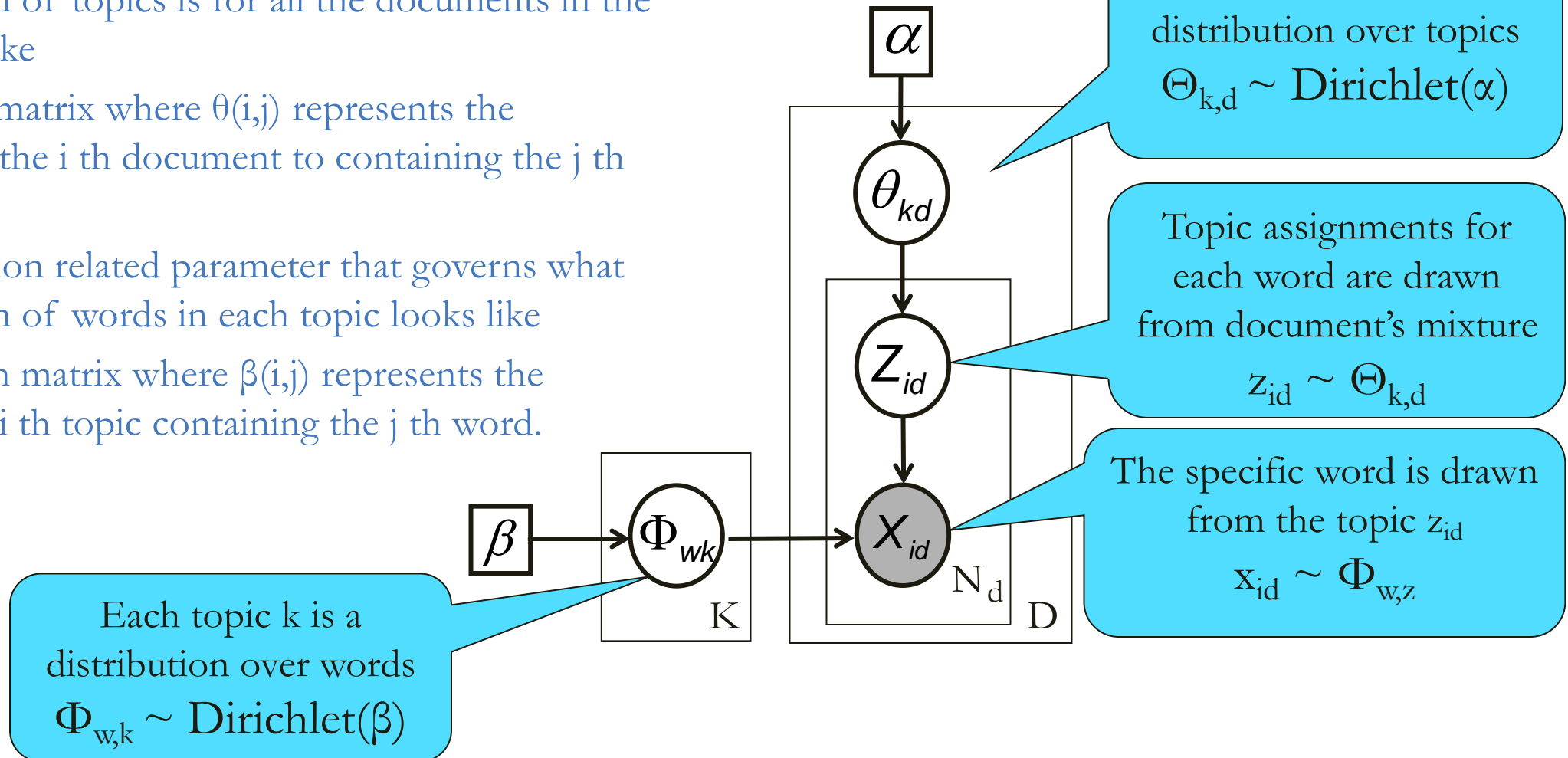
## ■ Intuition:

- Each document can be described by a distribution of topics and each topic can be described by a distribution of words

$$\begin{matrix} & D & & & K & & & D \\ \mathbb{W} & \boxed{\text{P (word | doc)}} & \approx & \mathbb{W} & \boxed{\text{P (word | topic)}} & * & K & \boxed{\text{P (topic | doc)}} \end{matrix}$$

# LDA Graphical Model

- $\alpha$  — Distribution related parameter that governs what the distribution of topics is for all the documents in the corpus looks like
- $\theta$  — Random matrix where  $\theta(i,j)$  represents the probability of the  $i$  th document to containing the  $j$  th topic
- $\eta$  — Distribution related parameter that governs what the distribution of words in each topic looks like
- $\beta$  — A random matrix where  $\beta(i,j)$  represents the probability of  $i$  th topic containing the  $j$  th word.



- Given a set of  $M$  documents, each having  $N$  words, where each word is generated by a single topic from a set of  $K$  topics:
- Find the joint posterior probability of:
  - $\theta$  — A distribution of topics, one for each document
  - $\mathbf{z}$  —  $N$  Topics for each document,
  - $\beta$  — A distribution of words, one for each topic

$$P(\theta_{1:M}, \mathbf{z}_{1:M}, \beta_{1:k} | \mathcal{D}; \alpha_{1:M}, \eta_{1:k})$$

- Assuming:
  - $\mathbf{D}$  — the corpus of documents
  - $\alpha$  — A parameter vector for each document (document — Topic distribution)
  - $\eta$  — A parameter vector for each topic (topic — word distribution)

# LDA: Another Explanation



- Go through each document and randomly assign each word to one of  $K$  topics  
[This random assignment gives topic representations of all documents and word distributions of all the topics, albeit not very good ones]
- For each document  $\mathbf{d}$ , go through each word  $\mathbf{w}$  and compute:
  - $p(\text{topic } t \mid \text{document } \mathbf{d})$ : proportion of words in document  $\mathbf{d}$  assigned to topic  $t$
  - $p(\text{word } \mathbf{w} \mid \text{topic } t)$ : proportion of assignments to topic  $t$ , over all documents  $\mathbf{d}$  that come from word  $\mathbf{w}$
- Reassign word  $\mathbf{w}$  a new topic  $t'$ , where we choose topic  $t'$  with probability
$$p(\text{topic } t' \mid \text{document } \mathbf{d}) * p(\text{word } \mathbf{w} \mid \text{topic } t')$$
- This generative model predicts the probability that topic  $t'$  generated word  $\mathbf{w}$