

# Topic Modeling

CPE 393: Text Analytics

***Dr. Sansiri Tarnpradab***

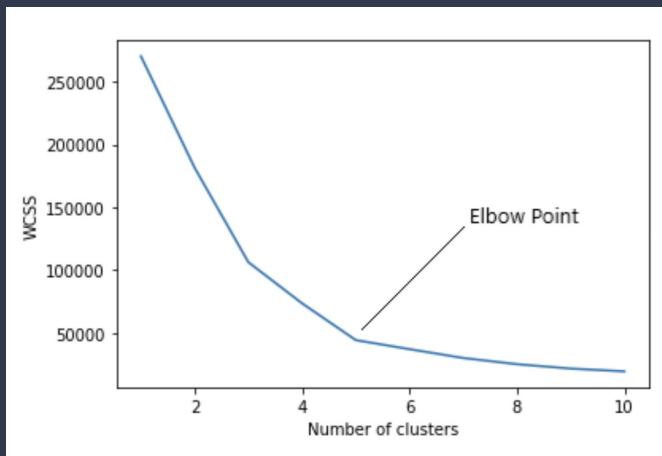
*Department of Computer Engineering  
King Mongkut's University of Technology Thonburi*

# Review

## *Text Clustering*

- **Text Clustering**
  - Similarity
  - Distance functions
  - Quality (inter-cluster, intra-cluster)
  - K-means clustering
- Some **drawbacks** for K-means clustering
  - Dependent on K value
  - Trial-and-error
- **Elbow method**

# Elbow Method



- A technique to determine an optimal number of clusters
  - **Plotting** output from different K values
  - **Identifying** the elbow point
- Computation
  - **WCSS (Within-Cluster Sum of Square)**
  - Squared average distance of all the points within a cluster to the cluster centroid
  - As the number of clusters increases, the WCSS value decreases
- Weaknesses
  - May not hold for complex datasets with irregularly shaped or differently sized clusters
  - Sensitive to initial cluster centroids
  - Inefficient for large datasets
  - Only works for K-means clustering

That was *Hard* Clustering...

*Intro*

*Pattern  
Matching*

*Text  
Visualization*

*Web Scraping*

*Text  
Preparation*

*Text Feature  
Representation*

*Text  
Classification*

*Text  
Summarization*

*Text  
Clustering*

*Topic  
Modeling*

*TBA*

*Presentation*



# Outline

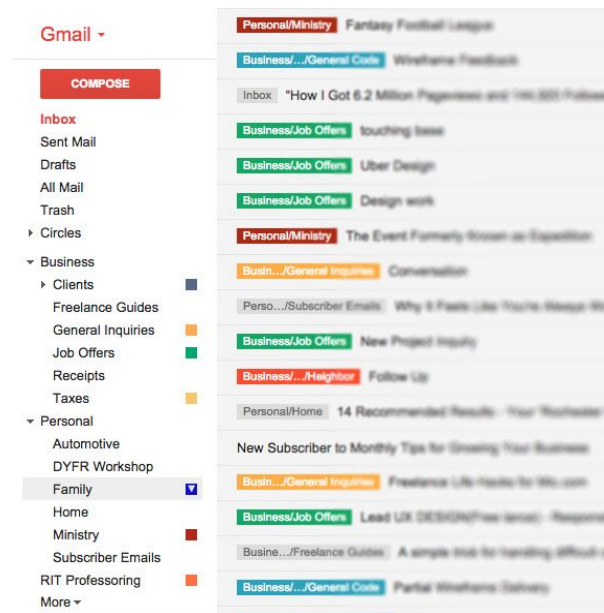
- Intro to Topic Modeling
- **LSA**: Latent Semantic Analysis
- **LDA**: Latent Dirichlet Allocation
- **BERTopic**

# Topic Modeling

Why?

Imagine you have **1M** emails to sort through..

- School
- Meeting
- Finance
- Internship



Ref: <https://mattolpinski.com/articles/gmail-label-system/>

# Topic Modeling

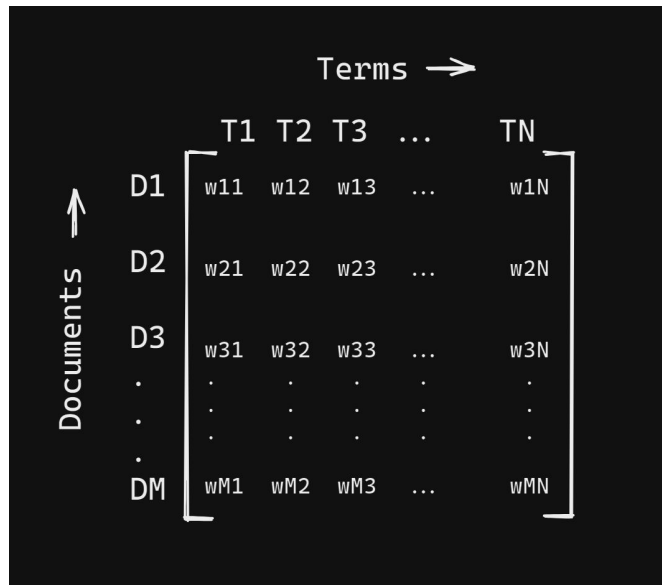
## *In a Nutshell*

- **Unsupervised Learning**
- A **technique** for discovering abstract topics in a collection of documents.
- **Goal:**
  - To discover latent topics within a corpus
  - Latent = Hidden
- **Key Elements:**
  - Every document is a mix of topics
  - Every topic is a mix of words
- **Input**
  - Topics & words
  - Document-term matrix
- **Output**
  - Various topics



# Revisit

## Document-Term Matrix



The diagram illustrates a Document-Term Matrix. It features a grid of weights  $w_{ij}$  where  $i$  represents the document index and  $j$  represents the term index. The rows are labeled D1, D2, D3, ..., DM on the left, with an upward arrow and the word 'Documents' indicating the direction of increasing document index. The columns are labeled T1, T2, T3, ..., TN at the top, with a rightward arrow and the word 'Terms' indicating the direction of increasing term index. The matrix is enclosed in large square brackets.

	Terms →					
	T1	T2	T3	...	TN	
D1	$w_{11}$	$w_{12}$	$w_{13}$	...	$w_{1N}$	
D2	$w_{21}$	$w_{22}$	$w_{23}$	...	$w_{2N}$	
D3	$w_{31}$	$w_{32}$	$w_{33}$	...	$w_{3N}$	
⋮	⋮	⋮	⋮	⋮	⋮	
⋮	⋮	⋮	⋮	⋮	⋮	
⋮	⋮	⋮	⋮	⋮	⋮	
DM	$w_{M1}$	$w_{M2}$	$w_{M3}$	...	$w_{MN}$	

# Latent Semantic Analysis (LSA)

# Latent Semantic Analysis: *LSA*

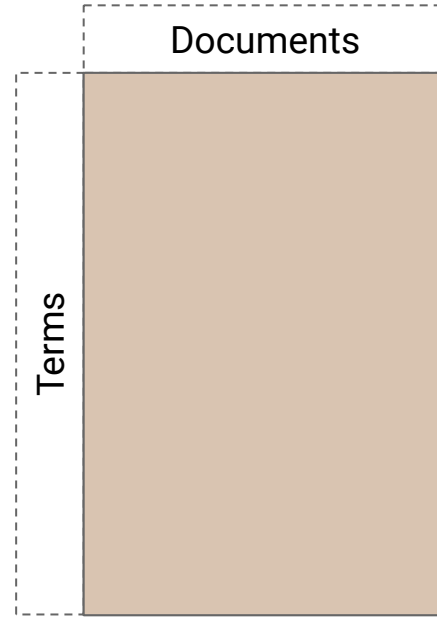
- Introduced in the late 1980s
- Statistical Method
- To analyze terms & documents relationships
- Text data → **Term-Document Matrix**

		Terms →				
		T1	T2	T3	...	TN
Documents ↑	D1	w11	w12	w13	...	w1N
	D2	w21	w22	w23	...	w2N
	D3	w31	w32	w33	...	w3N
	.	.	.	.	.	.
	.	.	.	.	.	.
	.	.	.	.	.	.
	DM	wM1	wM2	wM3	...	wMN

Documents  $\rightarrow$

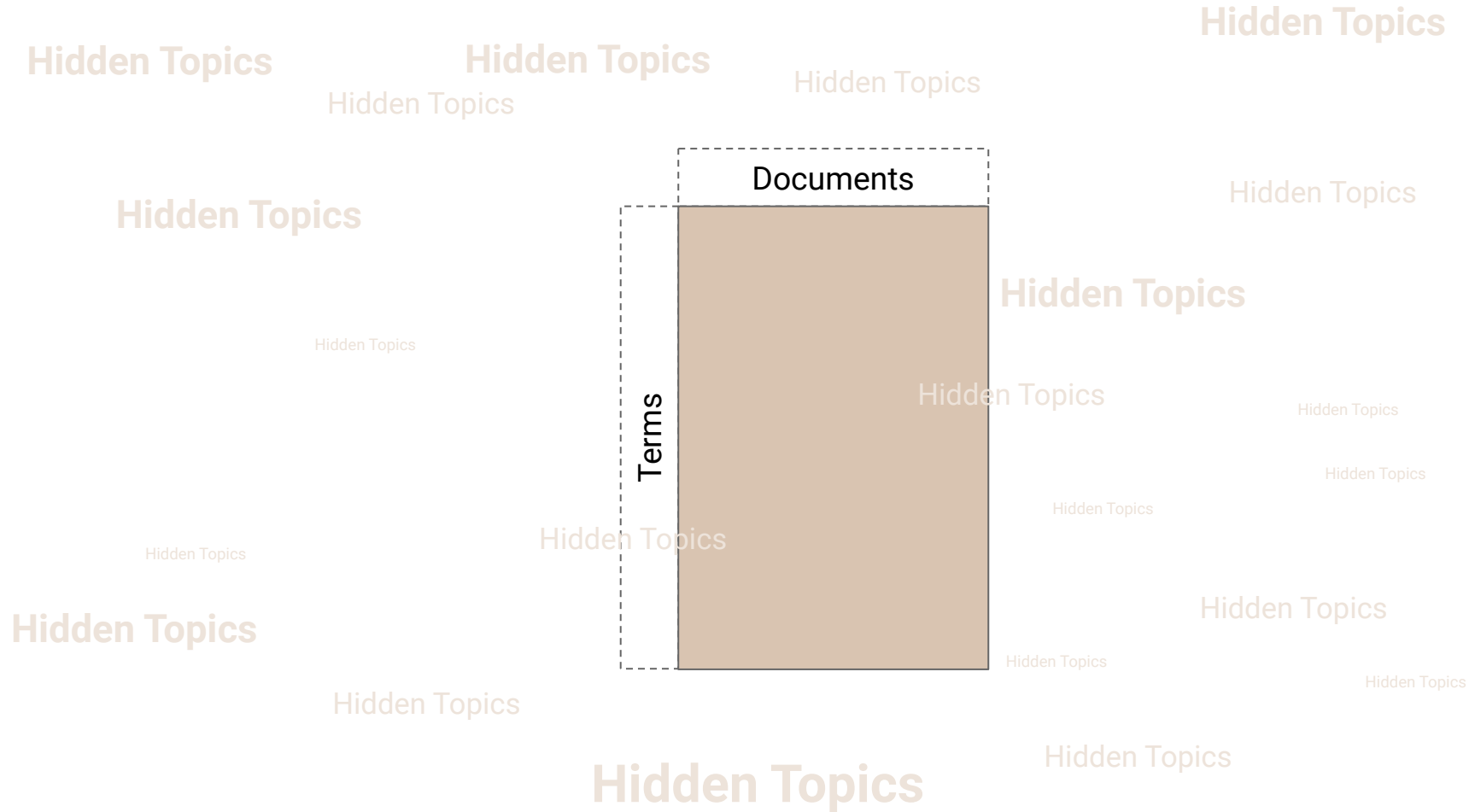
	Terms $\rightarrow$				
	T1	T2	T3	...	TN
D1	$w_{11}$	$w_{12}$	$w_{13}$	...	$w_{1N}$
D2	$w_{21}$	$w_{22}$	$w_{23}$	...	$w_{2N}$
D3	$w_{31}$	$w_{32}$	$w_{33}$	...	$w_{3N}$
.	.	.	.	.	.
.	.	.	.	.	.
DM	$w_{M1}$	$w_{M2}$	$w_{M3}$	...	$w_{MN}$

		Terms →				
Documents →		T1	T2	T3	...	TN
D1		$w_{11}$	$w_{12}$	$w_{13}$	...	$w_{1N}$
D2		$w_{21}$	$w_{22}$	$w_{23}$	...	$w_{2N}$
D3		$w_{31}$	$w_{32}$	$w_{33}$	...	$w_{3N}$
⋮		⋮	⋮	⋮	⋮	⋮
DM		$w_{M1}$	$w_{M2}$	$w_{M3}$	...	$w_{MN}$



Let's say..

This matrix presents a total of  $n$  documents and  $m$  terms.  
What's the dimension of this matrix?



# Matrix:

## *Dimension Property*

In order for matrix multiplication to be defined, the number of columns in the first matrix must be equal to the number of rows in the second matrix.

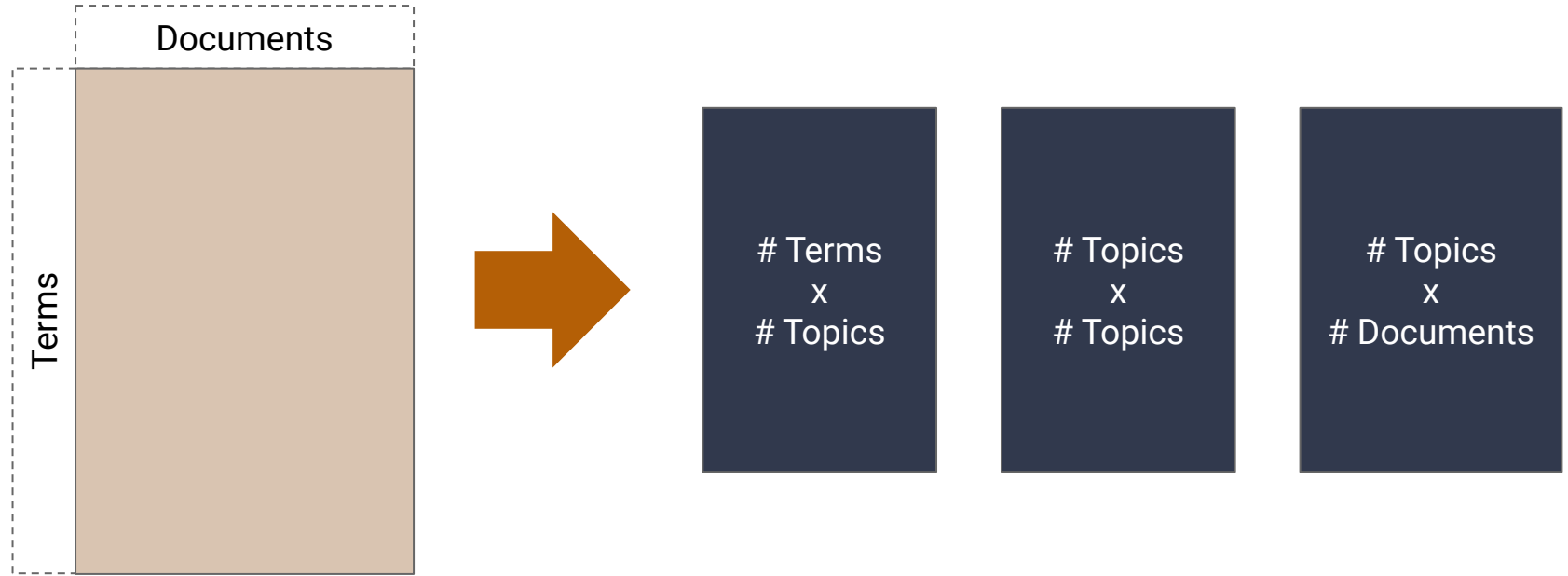
1	3
2	4
2	5

**A**

1	3	2	2
2	4	5	1

**B**

What's the dimension of AB?



# Matrix Factorization

$$A = U\Sigma V^T$$

## 1. **A**: Data Matrix

- $[m \times n] \rightarrow m$  terms,  $n$  documents
- Aka Term-Document Matrix

Singular Value Decomposition (SVD) decomposes the term-document matrix  $A$  (1) into three matrices (2)-(4):

## 2. **U**: Left singular vectors

- $[m \times k] \rightarrow m$  terms,  $k$  concepts
- Word Assignment to Topics

## 3. **$\Sigma$** : Diagonal matrix of singular values

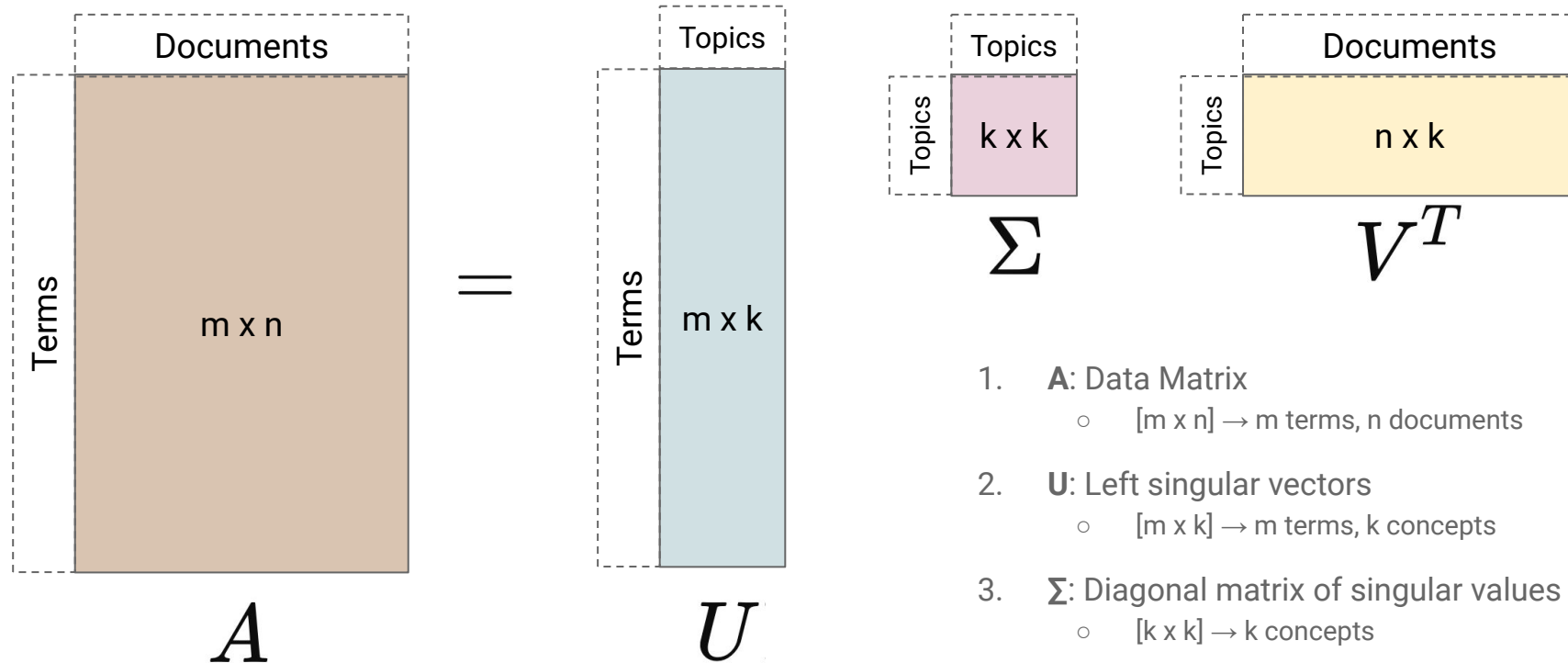
- $[k \times k] \rightarrow k$  concepts
- Topic Importance

## 4. **$V^T$** : Right singular vectors

- $[n \times k] \rightarrow n$  documents,  $k$  concepts
- Topic Distribution Across Documents

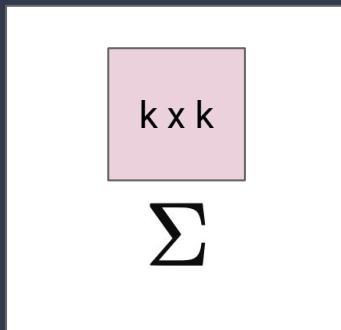


$$A = U\Sigma V^T$$



1. **A**: Data Matrix
  - $[m \times n] \rightarrow m$  terms,  $n$  documents
2. **U**: Left singular vectors
  - $[m \times k] \rightarrow m$  terms,  $k$  concepts
3.  **$\Sigma$** : Diagonal matrix of singular values
  - $[k \times k] \rightarrow k$  concepts
4.  **$V^T$** : Right singular vectors
  - $[n \times k] \rightarrow n$  documents,  $k$  concepts

# Diagonal Matrix of Singular Values



- Capture the significance of each latent dimension in the data
- Values arranged in descending order along the diagonal

12.5	0	0
0	9.0	0
0	0	1.4

- Each **singular value** indicates the importance of the concept in the reduced space
- Larger value:
  - Capture more of the variance in the data
  - Represent the most significant underlying pattern
- Hence, Singular Value Decomposition

# Singular Value Decomposition:

*SVD*

## Capture Hidden Patterns

→ Topics

## Dimensionality Reduction

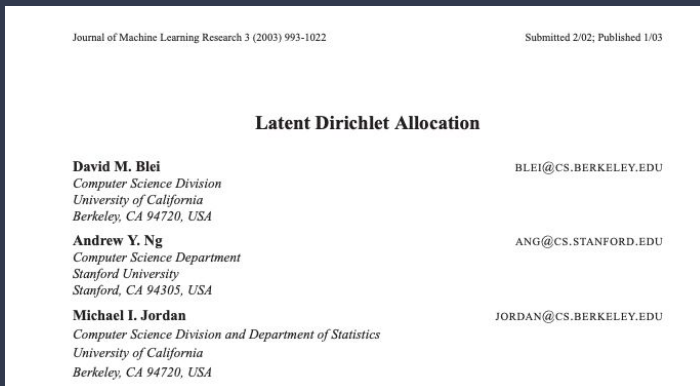
→ Retaining only the top k singular values

## Noise Reduction

- SVD in LSA helps filter out noise from the data
- Focusing on the most significant singular values and vectors

# Latent Dirichlet Allocation (LDA)

# Latent Dirichlet Allocation: *LDA*



[Read the paper](#)

- Introduced in 2003
- **Probabilistic** model
- **Latent** → Hidden
- **Dirichlet** → Type of probability distribution
- **Allocation** → Allocation
- Key Elements:
  - Every document is a mix of topics  
➔ **Every document is a distribution of topics**
  - Every topic is a mix of words  
**Every topic is a distribution of words**

# Example

## 5 Documents

The apple was crisp and sweet, bursting with flavor.

The playful puppy chased its tail in circles around the yard.

The lion roared loudly in the jungle, asserting its dominance.

The ripe banana was yellow and fragrant, ready to be enjoyed as a healthy snack.

The curious monkey reached for the juicy mango hanging from the tree.

The apple was crisp and sweet, bursting with flavor.	<b>TOPIC A</b>
The playful puppy chased its tail in circles around the yard.	<b>TOPIC B</b>
The lion roared loudly in the jungle, asserting its dominance.	<b>TOPIC B</b>
The ripe banana was yellow and fragrant, ready to be enjoyed as a healthy snack.	<b>TOPIC A</b>
The curious monkey reached for the juicy mango hanging from the tree.	<b>TOPIC A &amp; B</b>

The apple was crisp and sweet, bursting with flavor.	<b>TOPIC A</b>
The playful puppy chased its tail in circles around the yard.	<b>TOPIC B</b>
The lion roared loudly in the jungle, asserting its dominance.	<b>TOPIC B</b>
The ripe banana was yellow and fragrant, ready to be enjoyed as a healthy snack.	<b>TOPIC A</b>
The curious monkey reached for the juicy mango hanging from the tree.	<b>TOPIC A &amp; B</b>

**TOPIC A:** apple, banana, mango, crisp, sweet, flavor, juicy, ...

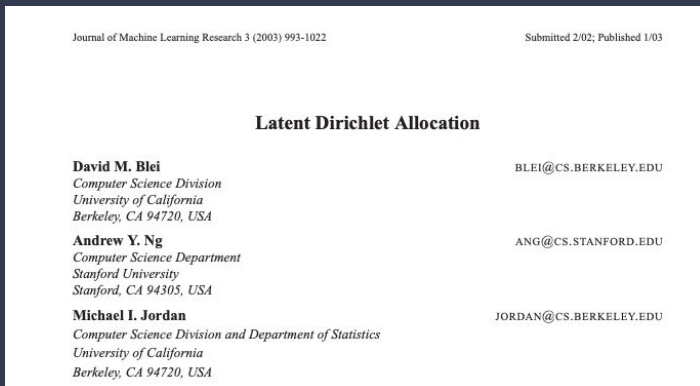
**TOPIC B:** puppy, tail, lion, monkey, roared, dominance, curious, ...

(each has percentage of distribution)



The apple was crisp and sweet, bursting with flavor.	<b>FRUIT</b>
The playful puppy chased its tail in circles around the yard.	<b>ANIMAL</b>
The lion roared loudly in the jungle, asserting its dominance.	<b>ANIMAL</b>
The ripe banana was yellow and fragrant, ready to be enjoyed as a healthy snack.	<b>FRUIT</b>
The curious monkey reached for the juicy mango hanging from the tree.	<b>FRUIT &amp; ANIMAL</b>

# Latent Dirichlet Allocation: *LDA*



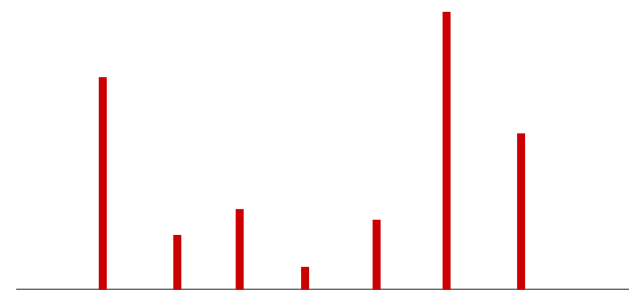
[Read the paper](#)

- 2003
- **Probabilistic** model
- **Latent** → Hidden
- **Dirichlet** → Type of probability distribution
- **Allocation** → Allocation
- Key Elements:
  - Every document is a mix of topics  
Every document is a distribution of topics
  - Every topic is a mix of words  
Every topic is a distribution of words

The apple was crisp and sweet, bursting with flavor.	<b>TOPIC A</b>
The playful puppy chased its tail in circles around the yard.	<b>TOPIC B</b>
The lion roared loudly in the jungle, asserting its dominance.	<b>TOPIC B</b>
The ripe banana was yellow and fragrant, ready to be enjoyed as a healthy snack.	<b>TOPIC A</b>
The curious monkey reached for the juicy mango hanging from the tree.	<b>TOPIC A &amp; B</b>

**TOPIC A:** apple, banana, mango, crisp, sweet, flavor, juicy, ...

**TOPIC B:** puppy, tail, lion, monkey, roared, dominance, curious, ...



Word distribution of TOPIC A

# How *LDA* works

FRUIT

ANIMAL

The playful puppy chased its tail in circles around the yard.

- Choose the number of topics beforehand (K)
- Randomly assign each word in each document to one of the K topics
- Go through every word and its assignment
  - How often the topic occurs in the document?
  - How often the word occurs in the topic overall?
- Assign the word to a new topic accordingly
- Perform this for multiple iterations

# Walk-Through

## Step 1

The playful puppy chased its tail in circles around the yard.

Choose the number of topics beforehand (K)

Let the number of topics be two:

1. Fruit
2. Animal

# Walk-Through

## Step 2

The playful puppy chased its tail in circles around the yard.

- Randomly assign each word in each document to one of the K topics
- Aka **Topic-word Distribution**

**Topic 1** (Fruit) might start with words like "chased" and "yard."

**Topic 2** (Animal) might start with words like "puppy" and "tail."

chased  
circles  
yard

Fruit

playful  
puppy  
tail

Animal

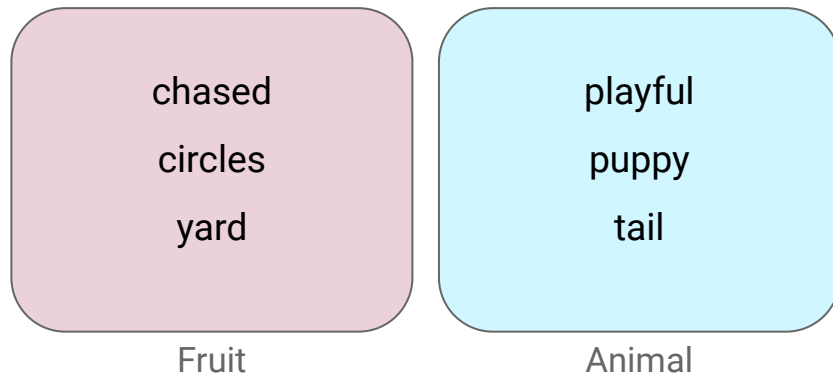
# Walk-Through

## Step 3 (Iterative)

The playful puppy chased its tail in circles around the yard.

Go through every word and its assignment

- How often the topic occurs in the document?
  - Document-Topic Distribution
- How often the word occurs in the topic overall?
  - Topic-Word Distribution



- Assign the word to a new topic accordingly

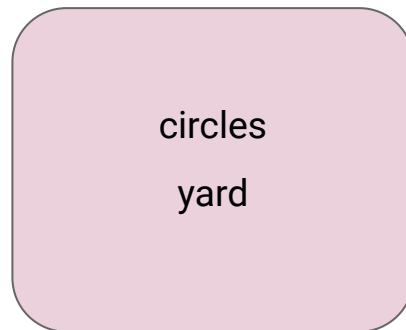
# Walk-Through

## Step ...

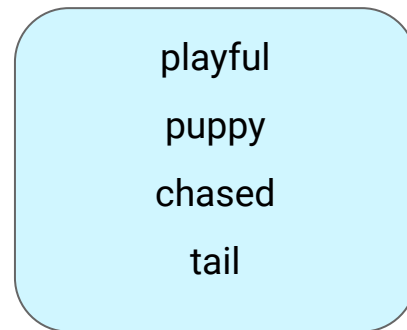
The playful puppy chased its tail in circles around the yard.

After several iterations, we may observe the following refinements:

- Words like "playful," "puppy," and "tail" increasingly align with Animal.
- Words like "yard" and "circles" might stay with Fruit if they lack strong relevance to Animal.



Fruit



Animal



# Walk-Through

## Final Step

The playful puppy chased its tail in circles around the yard.

Let's talk about the final output..

### Document-Topic Distribution:

- LDA will estimate that this document is primarily about the Animal topic (e.g., 80% Animal, 20% Fruit).

### Topic-Word Distribution:

- **Topic 1** (Fruit): Words like "circles" and "yard" may have a small association.
- **Topic 2** (Animal): Strongly associated with "playful," "puppy," "chased," and "tail."

# Multiple Documents

"I love playing football."

"The match was exciting."

"The player scored a goal."

## Topics

**Sports** and **Emotions**

# Adjustment of distribution through iterations..

I love playing football.

The match was exciting.

The player scored a goal.

With each iteration, LDA updates:

- Document-Topic Distribution
- Topic-Word Distribution

Let's talk about the final output..

Topic-Word Distribution

Document-Topic Distribution

# In Practice

```
class gensim.models.ldamodel.LdaModel(corpus=None, num_topics=100, id2word=None,
distributed=False, chunksize=2000, passes=1, update_every=1, alpha='symmetric', eta=None,
decay=0.5, offset=1.0, eval_every=10, iterations=50, gamma_threshold=0.001,
minimum_probability=0.01, random_state=None, ns_conf=None, minimum_phi_value=0.01,
per_word_topics=False, callbacks=None, dtype=<class 'numpy.float32'>)
```

## Input

- Document-Term Matrix
- Number of topics (K)
- Number of iterations

## Gensim

- Go through every word
- Find the best word & topic distribution
- Assignment

## Output

- Top words in each topic
- Adjust parameter as needed

# LDA VS LSA

	LSA	LDA
<b>Type</b>	Statistical method	Probabilistic method
<b>Technique</b>	Singular Value Decomposition (SVD)	Bayesian inference with Dirichlet distributions
<b>Document Representation</b>	Concepts are derived from word co-occurrence patterns	Topics are probabilistic distributions over words
<b>Scalability</b>	Effective with smaller datasets; performance decreases with size	Scales well with large datasets, widely used in big data applications
<b>Interpretability</b>	Concepts are less interpretable; dimensions represent abstract concepts	Topics are explicit with assigned probabilities
<b>Applications</b>	Better for semantic search, synonym discovery	Better for text generation, topic categorization

# Applications

- Text summarization
- Information retrieval
- Recommendation systems
- Sentiment analysis
- And more...

# Conclusion

- Intro to Topic Modeling
- **LSA**: Latent Semantic Analysis
- **LDA**: Latent Dirichlet Allocation
- **BERTopic**

# Q & A