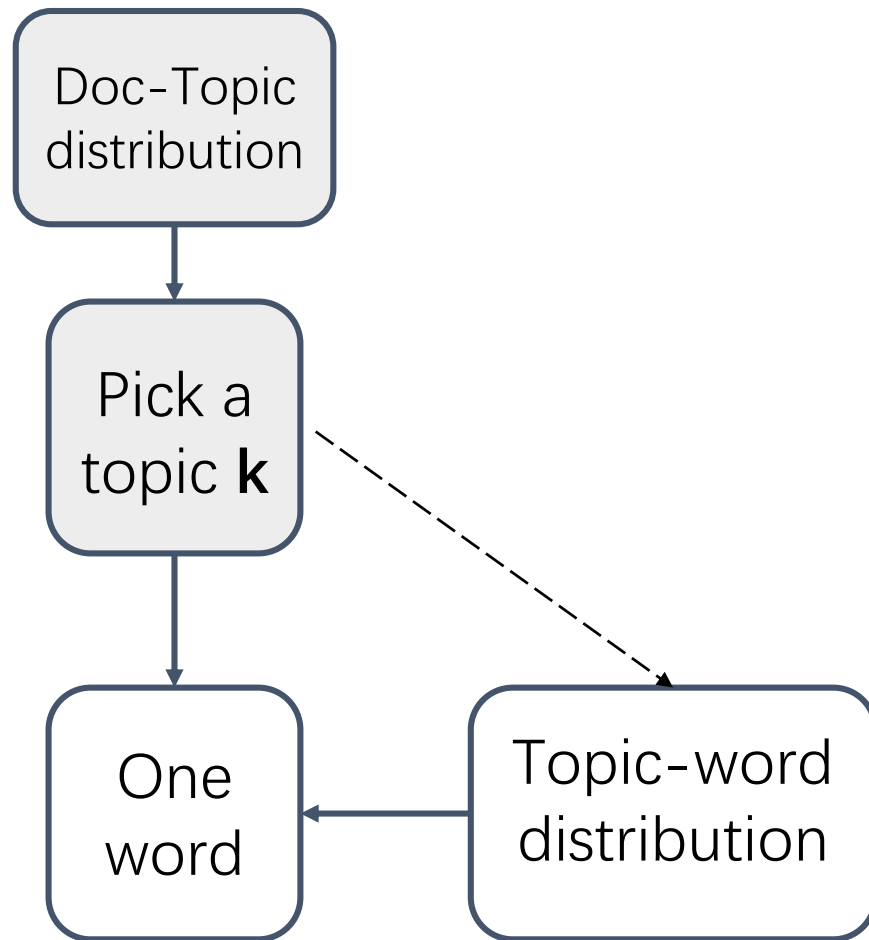# Document Classification

## based on

# LDA

Latent Dirichlet Allocation
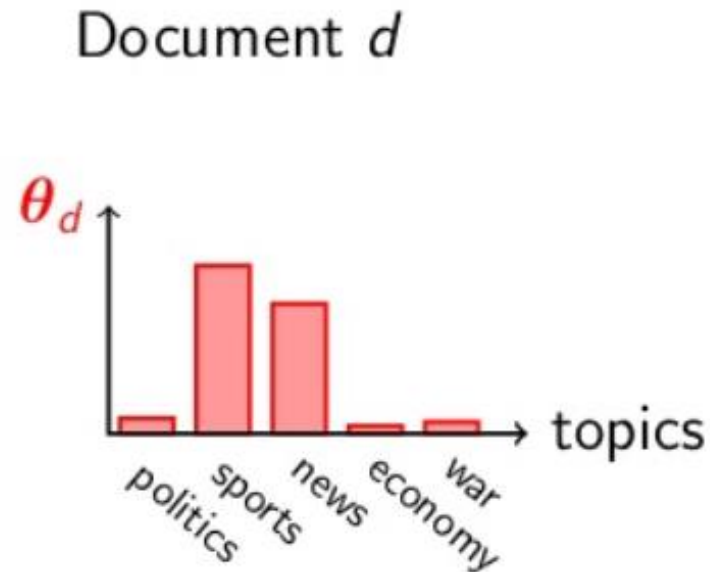
Chunyuan Li
Jiarui Xing

# Contents

- Introduction to LDA

- Experiment Design
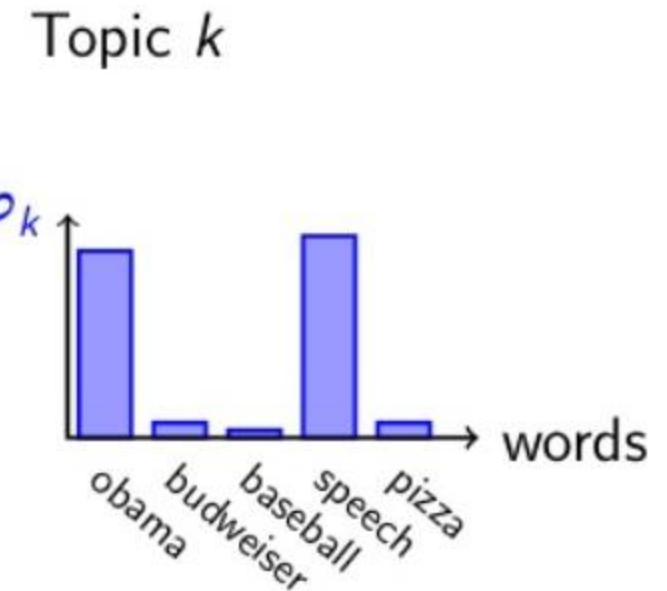
- Experiment Result

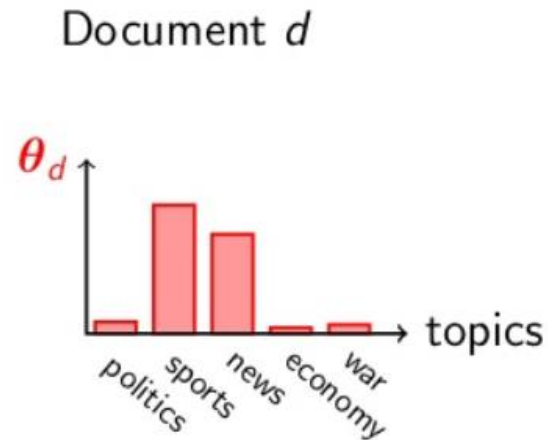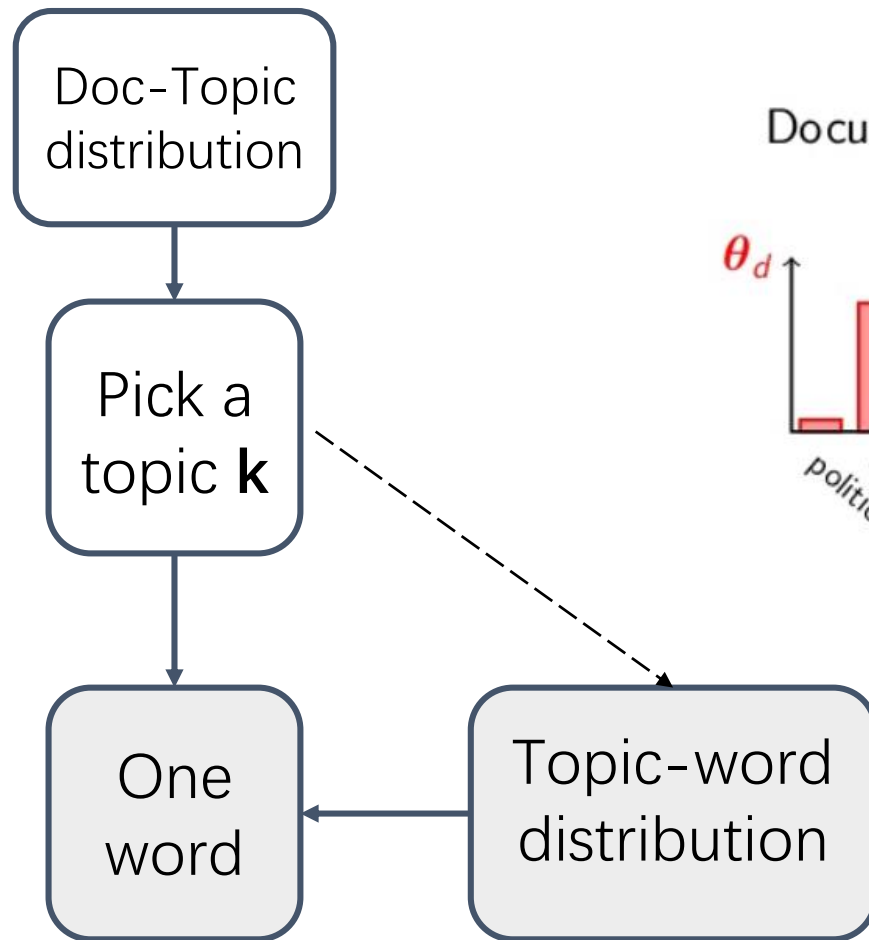## LDA Document Generation Model (Simplified)

Doc-Topic distribution

↓

Pick a topic **k**

↓

One word ← Topic-word distribution

First Step:
**choose a topic from document-topic distribution**



Document $d$

$\theta_d$ → topics

politics  sports  news  economy  war

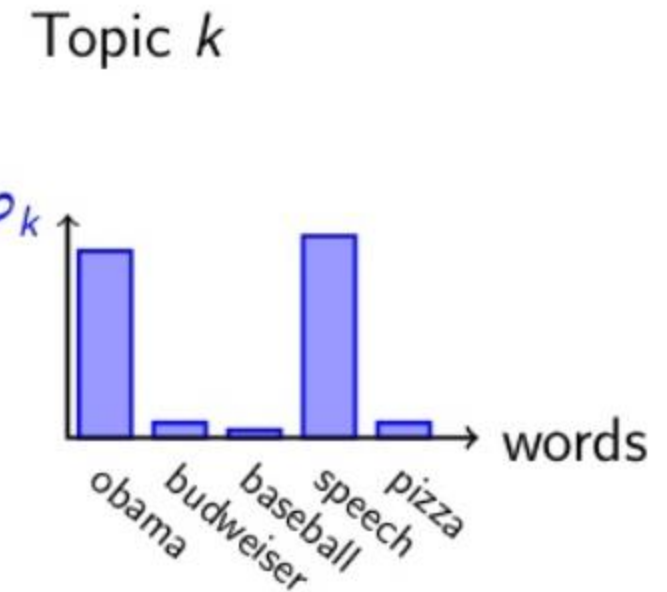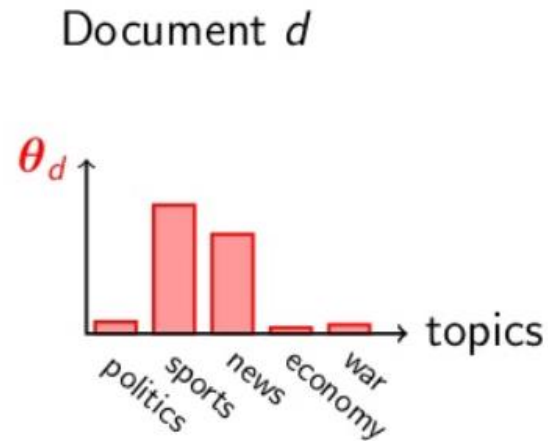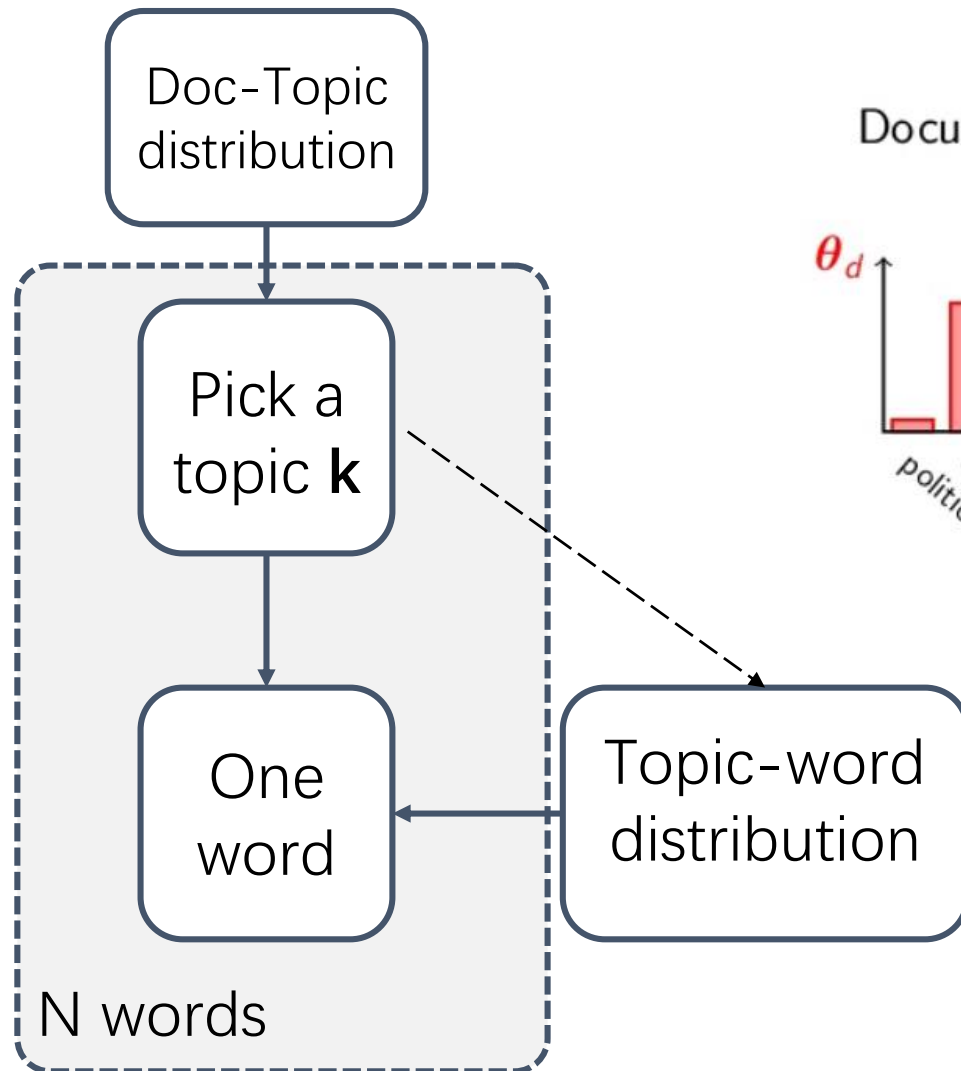Figure: datasciencecentral.com/profiles/blogs/a-tale-about-lda2vec-when-lda-meets-word2vec

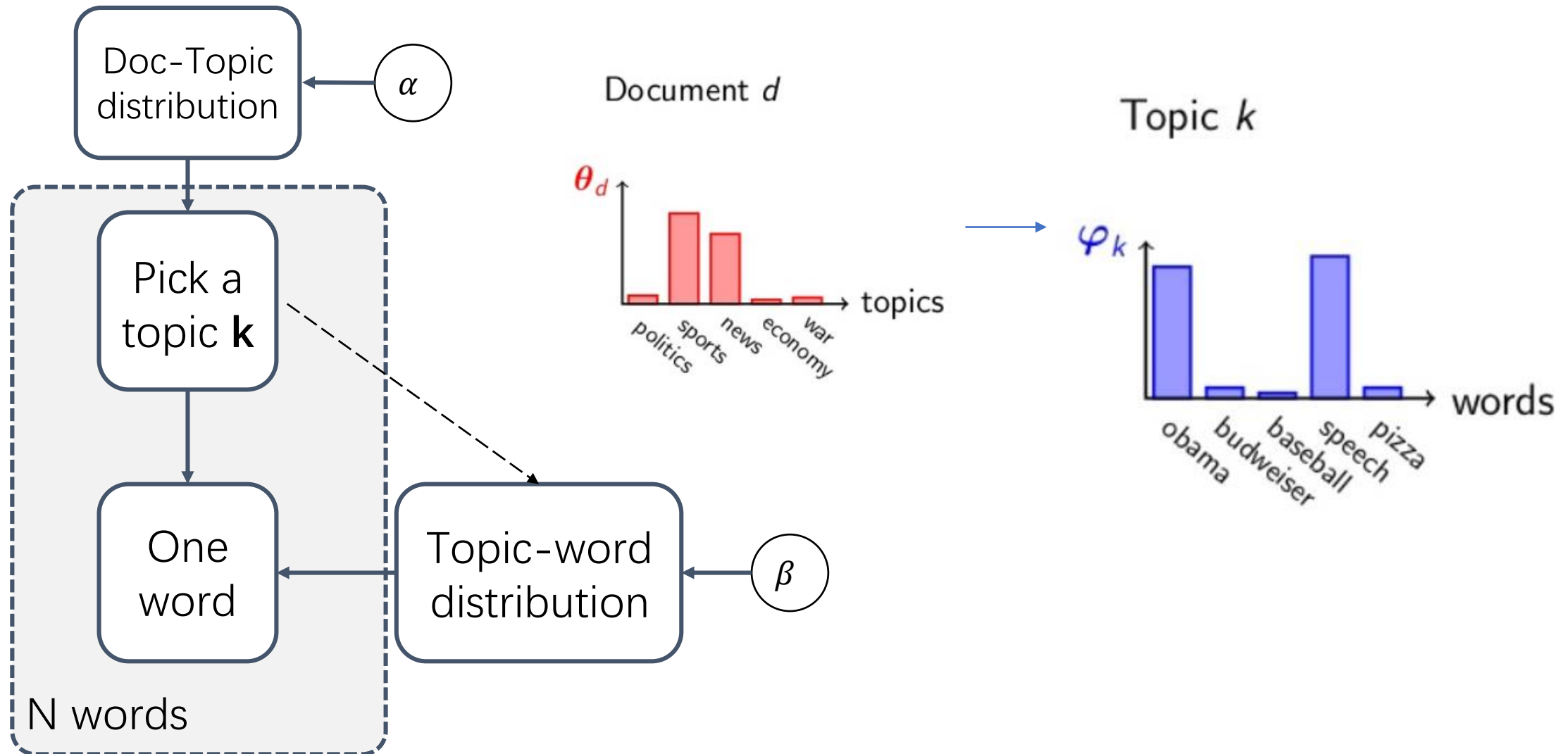## LDA Document Generation Model (Simplified)



Second Step:
**choose a word from topic-word distribution**

## LDA Document Generation Model (Simplified)



Do it N times to generate a document with N words

## LDA Document Generation Model (Simplified)

- **Basic Idea**

As different documents have difference in the

usage of words (e.g. term frequency/tf-idf), they

should also **differ in topics**!



words → game, video, Nintendo, Japan, …

topics → {Nintendo, game, video,…}
{media, franchise, company,…}

Wikipedia:
Pokémon



words → statistical, math, inference, …

topics → {Bayesian, statistical, inference}
{prior, posterior, use,…}

Wikipedia:
Bayesian Inference

- ## Basic Idea

  As different documents have difference in the usage

  of words (e.g. term frequency or tf-idf), they should also

  differ in topics!

- ## Implementation: LDA as Vectorization!

  Transform a document into **distribution of topics**

Document

| 0.12 | 0.08 | 0.01 | 0.23 | ... | 0.01 | 0.01 | 0.02 | 0.23 | 0.13 |

Topic Distribution

- **When should LDA work better (than other document vectorization methods)?**

    - With lower vector dimension?

        Each topic referring to several words, which may represent a document in lower dimensions (than tf and tf-idf).

# of words

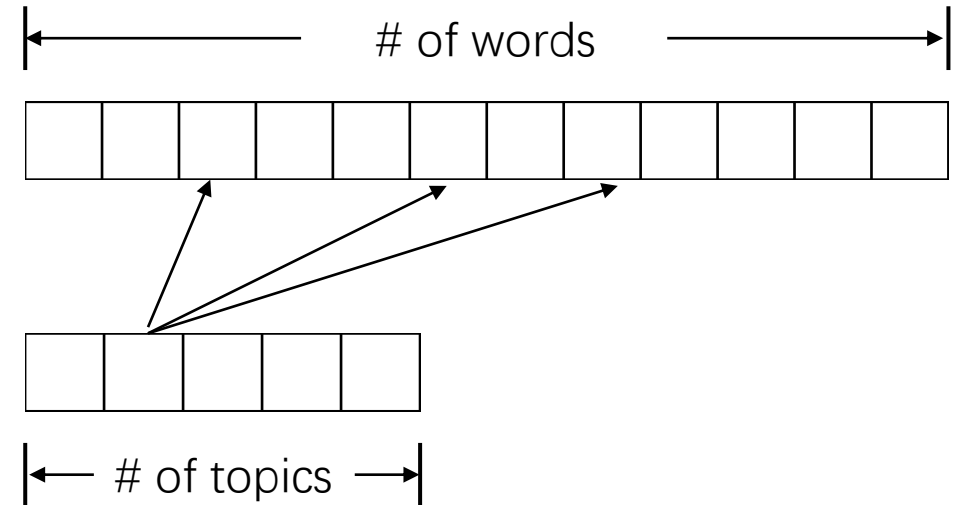# of topics

- **When should LDA work better (than other document vectorization methods)?**

  - With lower vector dimension?

    Each topic referring to several words, which may represent a document in lower dimensions (than tf and tf-idf).

  - With small dataset?

    Prior helps us avoid overfitting (think of pseudo-count in coin flipping case)
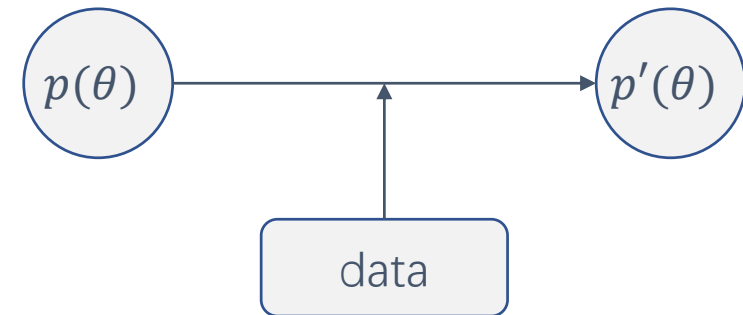
- **When should LDA work better (than other document vectorization methods)?**

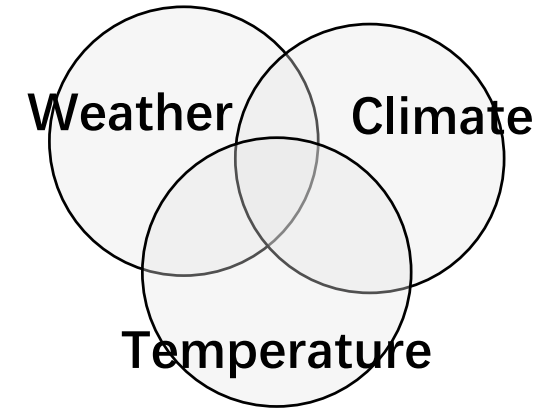  - With lower vector dimension?

    Each topic referring to several words, which may represent a document in lower dimensions (than tf and tf-idf).
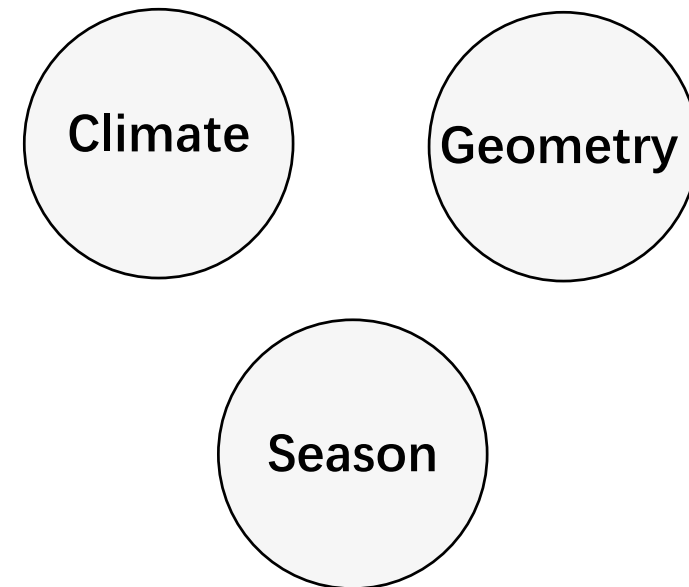
  - With small dataset?

    Prior helps us avoid overfitting(think of pseudo-count in coin flipping case)

  - With smaller document-topic-prior ($\alpha$)?

    "Sparse" topics work better
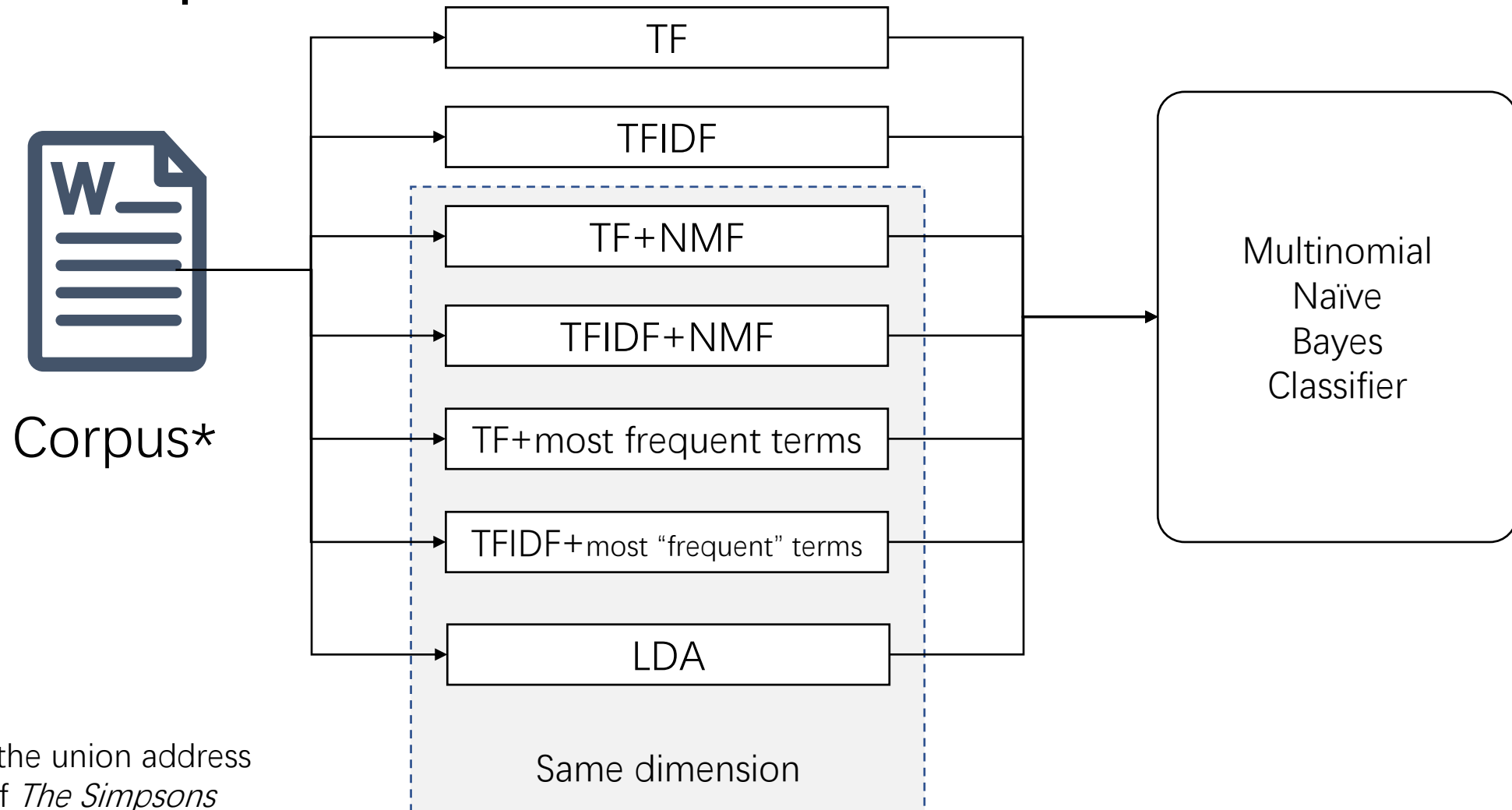


Weather  Climate  Temperature

Large $\alpha$

Climate  Geometry  Season

Small $\alpha$

- **Experiment Pipeline**



Corpus*

| TF |
| TFIDF |
| TF+NMF |
| TFIDF+NMF |
| TF+most frequent terms |
| TFIDF+most "frequent" terms |
| LDA |

Same dimension

Multinomial Naïve Bayes Classifier

* 1. State of the union address
  2. Scripts of *The Simpsons*

- **Experiment Pipeline**



Corpus*

TF

TFIDF

TF+NMF

TFIDF+NMF

TF+most frequent terms

TFIDF+most "frequent" terms

LDA

Dataset size | $\alpha$ | Dimen-sion

Multinomial Naïve Bayes Classifier

* 1. State of the union address
  2. Scripts of *The Simpsons*

Average 10-rerun Accuracies, dimension=5, alpha=0.0001, test_size=0.3

# Experiment Result  Accuracies with $\alpha = 0.0001, \mathbf{0.01}, 1, 10$



Average 10-rerun Accuracies, dimension=5, alpha=0.01, test_size=0.3

Average 10-rerun Accuracies, dimension=5, alpha=1, test_size=0.3

Average 10-rerun Accuracies, dimension=5, alpha=10, test_size=0.3

Average 10-rerun Accuracies, dimension=5, alpha=10, test_size=0.3

## Conclusions

- As $\alpha$ increases, LDA accuracy decreases

- When $\alpha$ is small, LDA works better than other vectorization methods with same dimension

Average 10-rerun Accuracies, dimension=5, alpha=0.0001, test_size=0.3

Average 10-rerun Accuracies, dimension=5, alpha=0.0001, test_size=0.5

Average 10-rerun Accuracies, dimension=5, alpha=0.0001, test_size=0.8

Average 10-rerun Accuracies, dimension=5, alpha=0.0001, test_size=0.8

## Conclusions

- With size of training set decreasing, Accuracies of TF and TF-IDF decrease (of course)

- Accuracies of methods with fixed dimension other than LDA keep nearly unchanged.

- Surprisingly, LDA works better than others only with larger dataset!

Average 10-rerun Accuracies, dimension=5, alpha=0.01, test_size=0.8

Average 10-rerun Accuracies, dimension=10, alpha=0.01, test_size=0.8

Average 10-rerun Accuracies, dimension=20, alpha=0.01, test_size=0.8

Average 10-rerun Accuracies, dimension=20, alpha=0.0001, test_size=0.3

1.0

Accuracy

10-fold Cross Validation
Test Accuracy

## Conclusions

- With dimension increasing, accuracies of other dimension-fixed methods increase (of course)

- However, that of LDA keep nearly unchanged.

- LDA doesn't work better than others.

.69

37.65
32.84
34.69    34.31
39.28    38.91

0.0

tf    tfidf    tf_dim    tfidf_dim    tf_nmf    tfidf_nmf    LDA

Model

## Conclusions

- **When should LDA work better (than other document vectorization methods)?**

    - With smaller document-topic-prior ($\alpha$)?

      Yes!

    - With small dataset?

      No!

    - With lower vector dimension?

      Not better than using other dimension reduction methods.

# Questions?