

The background features a complex network of thin, light-colored lines forming a mesh-like structure. Scattered throughout are small, colored dots in shades of green, blue, and orange. A prominent, darker, reddish-brown geometric shape, resembling a stylized dome or a complex polygon, is centered in the upper half. The overall color palette is muted, with earthy tones and soft pastels.

# **ToPMine: Phrase Mining without Training Data**



# Strategy 3: First Phrase Mining then Topic Modeling

- Why first Phrase Mining then Topic Modeling?
  - With Strategy 2, tokens in the same phrase may be assigned to different topics
    - Ex. *knowledge discovery* using *least squares support vector machine classifiers*...
    - *Knowledge discovery* and *support vector machine* should have coherent topic labels

- Solution: switch the order of phrase mining and topic model inference

[knowledge discovery] using [least squares] [support vector machine] [classifiers] ...



[knowledge discovery] using [least squares] [support vector machine] [classifiers] ...

- Techniques for this strategy
  - ① □ Phrase mining, document segmentation, and phrase ranking
  - ② □ Topic model inference with phrase constraint



# ToPMine: Phrase Mining before Topic Modeling

- ToPMine [El-Kishky et al. VLDB'15]: Phrase mining, then phrase-based topic modeling

## Phrase mining

- ① Frequent **contiguous pattern** mining: Extract candidate phrases and their counts
- ② Agglomerative merging of adjacent unigrams as guided by a **significance score**
- ③ Document segmentation to count phrase occurrence
  - Calculate rectified (i.e., true) phrase frequency
- ④ Phrase ranking (using the criteria proposed in KERT)
  - Popularity, concordance, informativeness, completeness

Phrase	Raw frequency	Rectified frequency
[support vector machine]	90	80
[vector machine]	95	0
[support vector]	100	20

## Phrase-based topic modeling

- The mined bag-of-phrases are passed as input to PhraseLDA, an extension of LDA, that constrains all words in a phrase to each sharing the same latent topic

# Collocation Mining

- Collocation: A sequence of words that occur more frequently than expected
  - Often “interesting”, relay information not portrayed by their constituent terms
    - Ex. “made an exception”, “strong tea”
- Many different measures used to extract collocations from a corpus [Dunning 93, Pederson 96]
  - E.g., mutual information, t-test, z-test, chi-squared test, likelihood ratio

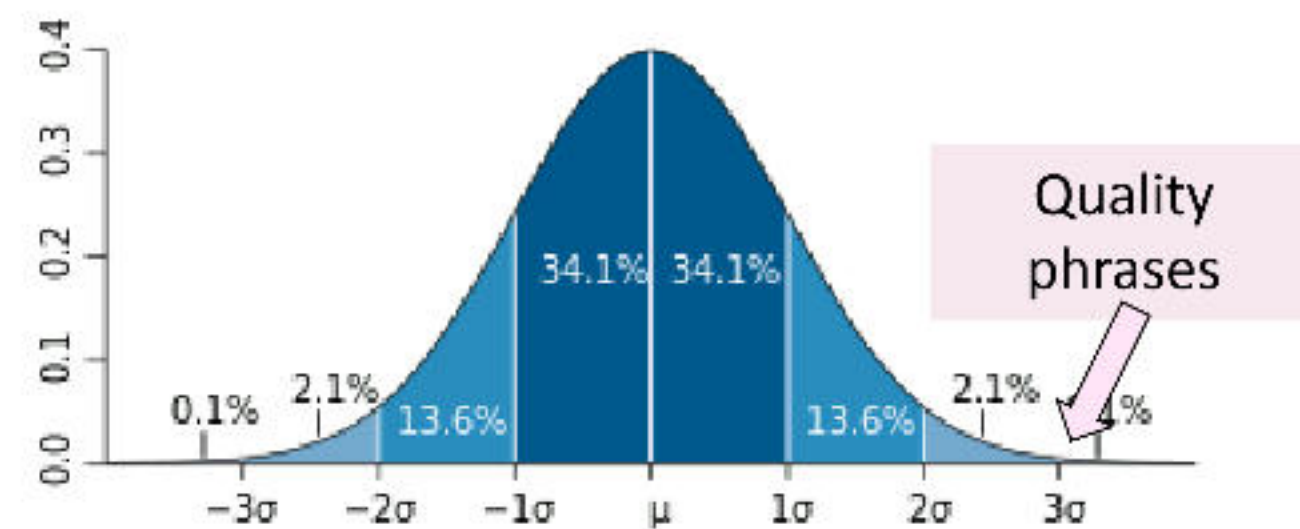
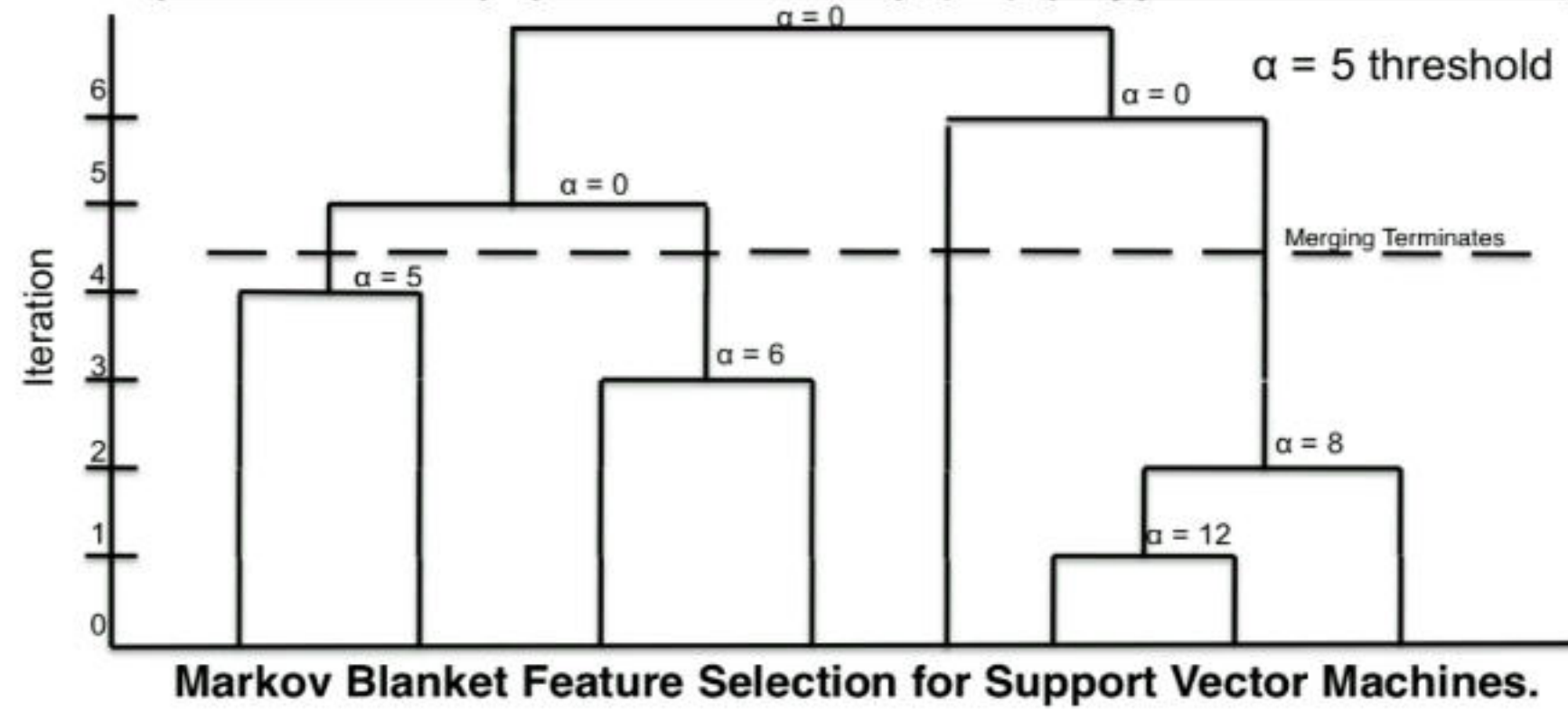
$$\text{PMI}(x, y) = \log \frac{p(x, y)}{p(x)p(y)} \quad \text{sig} = \frac{\text{count}(\text{phr}_{x+y}) - E[\text{count}(\text{phr}_{x+y})]}{\sqrt{\text{count}(\text{phr}_{x+y})}} \quad \chi^2 = \sum \frac{(O - E)^2}{E}$$

- Many of these measures can be used to guide the agglomerative **phrase-segmentation** algorithm



# Phrase Candidate Generation: Frequent Pattern Mining + Statistical Analysis

(Markov Blanket) (Feature Selection) (for) (Support Vector Machines)



Based on significance score [Church et al.'91]:

$$\alpha(P_1, P_2) \approx (f(P_1 \bullet P_2) - \mu_0(P_1, P_2)) / \sqrt{f(P_1 \bullet P_2)}$$

Note for the first title:

- [feature selection] forms phrase but not [selection for] based on the significant scores computed
- [support vector machine] does not contribute to the counts of [support], [vector], [support vector], [vector machine]

[Markov blanket] [feature selection] for [support vector machines]

[knowledge discovery] using [least squares] [support vector machine] [classifiers]

...[support vector] for [machine learning]...



# ToPMine: Experiments on DBLP Abstracts

	<i>Topic 1</i>	<i>Topic 2</i>	<i>Topic 3</i>	<i>Topic 4</i>	<i>Topic 5</i>
unigrams	problem algorithm optimal solution search solve constraints programming heuristic genetic	word language text speech system recognition character translation sentences grammar	data method algorithm learning clustering classification based features proposed classifier	programming language code type object implementation system compiler java data	data patterns mining rules set event time association stream large
n-grams	genetic algorithm optimization problem solve this problem optimal solution evolutionary algorithm local search search space optimization algorithm search algorithm objective function	natural language speech recognition language model natural language processing machine translation recognition system context free grammars sign language recognition rate character recognition	data sets support vector machine learning algorithm machine learning feature selection paper we propose clustering algorithm decision tree proposed method training data	programming language source code object oriented type system data structure program execution run time code generation object oriented programming java programs	data mining data sets data streams association rules data collection time series data analysis mining algorithms spatio temporal frequent itemsets

ToPMine is efficient and generates high-quality topics and phrases without any training data

特性.



# ToPMine: Experiments on Yelp Reviews

	<i>Topic 1</i>	<i>Topic 2</i>	<i>Topic 3</i>	<i>Topic 4</i>	<i>Topic 5</i>
unigrams	coffee ice cream flavor egg chocolate breakfast tea cake sweet	food good place ordered chicken roll sushi restaurant dish rice	room parking hotel stay time nice place great area pool	store shop prices find place buy selection items love great	good food place burger ordered fries chicken tacos cheese time
n-grams	ice cream iced tea french toast hash browns frozen yogurt eggs benedict peanut butter cup of coffee iced coffee scrambled eggs	spring rolls food was good fried rice egg rolls chinese food pad thai dim sum thai food pretty good lunch specials	parking lot front desk spring training staying at the hotel dog park room was clean pool area great place staff is friendly free wifi	grocery store great selection farmer's market great prices parking lot wal mart shopping center great place prices are reasonable love this place	mexican food chips and salsa food was good hot dog rice and beans sweet potato fries pretty good carne asada mac and cheese fish tacos

ToPMine works well for phrase and topic mining in social media data