

Session 4. Strategy 3: First Phrase Mining then Topic Modeling

Strategy 3: First Phrase Mining then Topic Modeling

- ❑ ToPMine [El-Kishky et al. VLDB'15]
 - ❑ First phrase construction, then topic mining
 - ❑ Contrast with KERT: topic modeling, then phrase mining
- ❑ The ToPMine Framework:
 - ❑ Perform **frequent *contiguous pattern* mining** to extract candidate phrases and their counts
 - ❑ Perform agglomerative merging of adjacent unigrams as guided by a significance score—This segments each document into a ***“bag-of-phrases”***
 - ❑ The newly formed 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

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

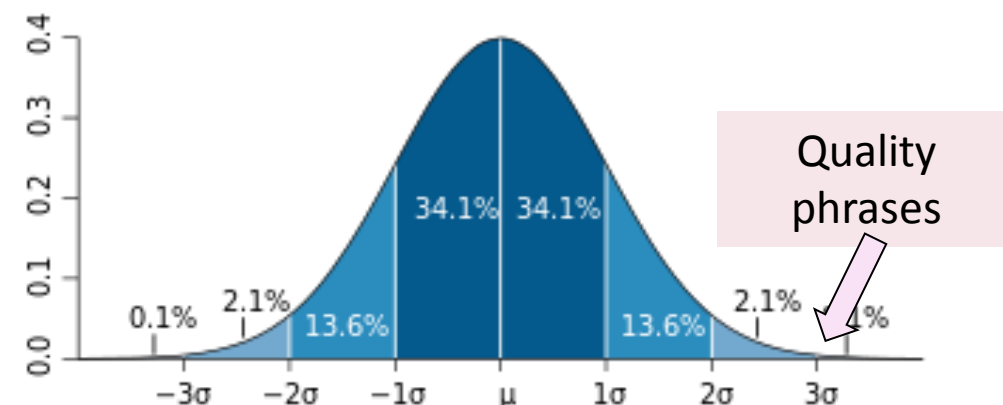
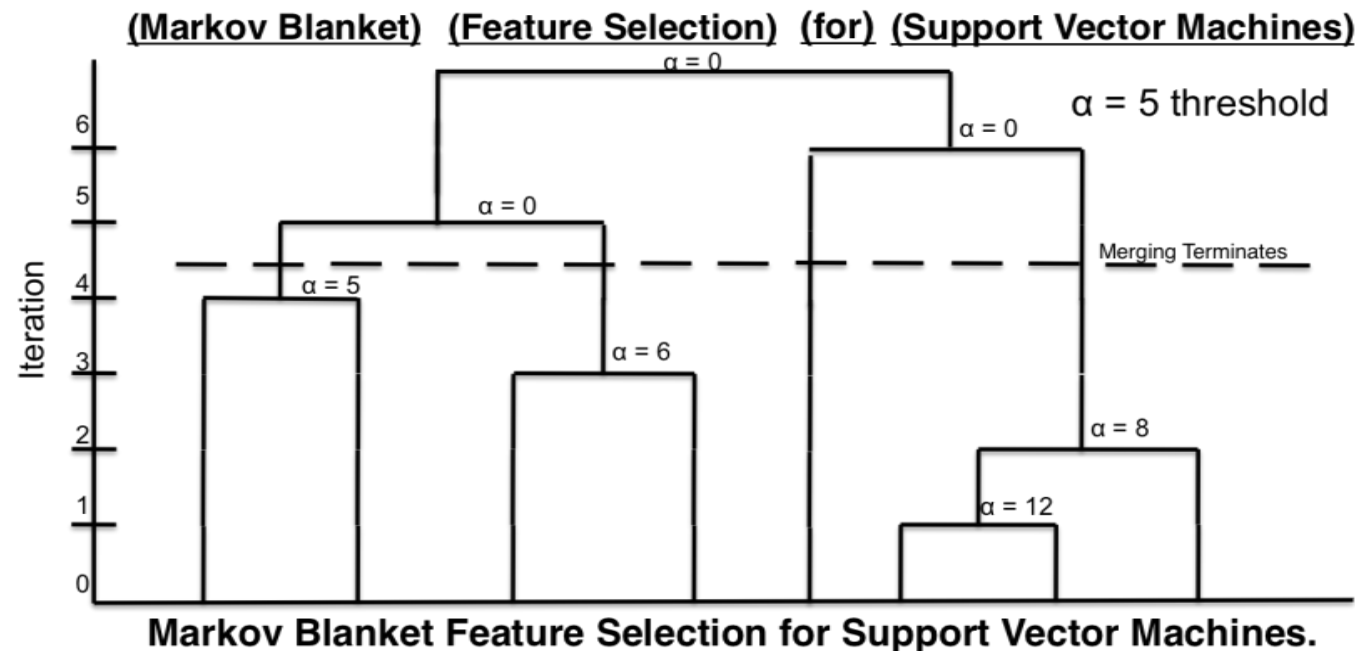
[knowledge discovery] using [least
squares] [support vector machine]
[classifiers] ...



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

- ❑ Techniques
 - ❑ Phrase mining and document segmentation
 - ❑ Topic model inference with phrase constraint

Phrase Mining: Frequent Pattern Mining + Statistical Analysis



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

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

[Markov blanket] [feature selection] for [support vector machines]
[knowledge discovery] using [least squares] [support vector machine] [classifiers]
...[support vector] for [machine learning]...

Phrase	Raw freq.	True freq.
[support vector machine]	90	80
[vector machine]	95	0
[support vector]	100	20

Collocation Mining

- Collocation: A sequence of words that occur more frequently than expected
 - Often “interesting” and due to their non-compositionality, often relay information not portrayed by their constituent terms (e.g., “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

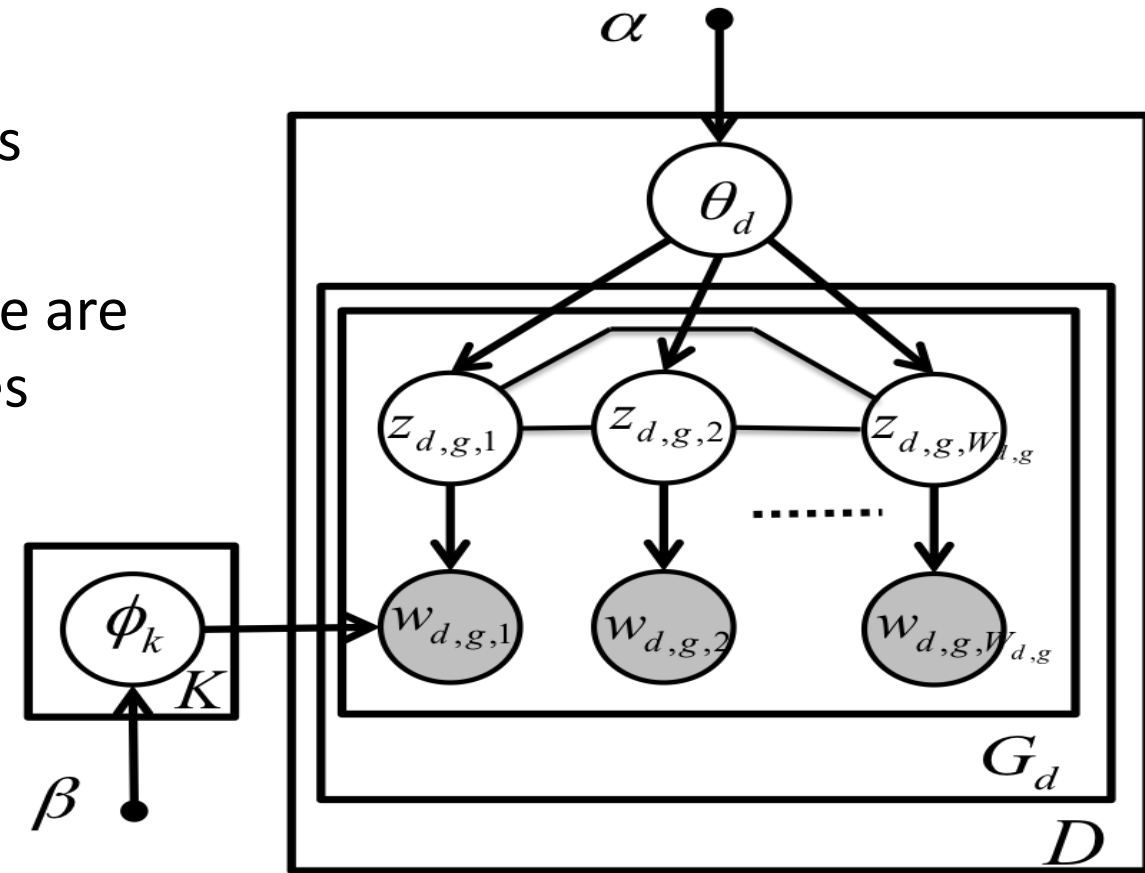
$$\text{PMI}(x, y) = \log \frac{p(x, y)}{p(x)p(y)} \quad \text{sig} = \frac{\text{count}(\text{phr}_{x+y}) - E[\text{count}(\text{phr}_{x+y})]}{\sqrt{\text{count}(\text{phr}_{x+y})}} \quad \chi^2 = \sum \frac{(O - E)^2}{E}$$

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

ToPMine: Phrase LDA (Constrained Topic Modeling)

- The generative model for PhraseLDA is the same as LDA
- Difference: the model incorporates constraints obtained from the “**bag-of-phrases**” input
- Chain-graph shows that all words in a phrase are constrained to take on the same topic values

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



Topic model inference with phrase constraints

Example Topical Phrases: A Comparison

social networks	information retrieval
web search	text classification
time series	machine learning
search engine	support vector machines
management system	information extraction
real time	neural networks
decision trees	text categorization
:	:
Topic 1	Topic 2

PDLDA [Lindsey et al. 12] Strategy 1
(3.72 hours)

information retrieval	feature selection
social networks	machine learning
web search	semi supervised
search engine	large scale
information extraction	support vector machines
question answering	active learning
web pages	face recognition
:	:
Topic 1	Topic 2

ToPMine [El-kishky et al. 14]
Strategy 3 (67 seconds)

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: Topics on Associate Press News (1989)

	<i>Topic 1</i>	<i>Topic 2</i>	<i>Topic 3</i>	<i>Topic 4</i>	<i>Topic 5</i>
unigrams	plant nuclear environmental energy year waste department power state chemical	church catholic religious bishop pope roman jewish rev john christian	palestinian israeli israel arab plo army reported west bank state	bush house senate year bill president congress tax budget committee	drug aid health hospital medical patients research test study disease
n-grams	energy department environmental protection agency nuclear weapons acid rain nuclear power plant hazardous waste savannah river rocky flats nuclear power natural gas	roman catholic pope john paul john paul catholic church anti semitism baptist church united states lutheran church episcopal church church members	gaza strip west bank palestine liberation prganization united states arab reports prime minister yitzhak shamir israel radio occupied territories occupied west bank	president bush white house bush administration house and senate members of congress defense secretary capital gains tax pay raise house members committee chairman	health care medical center united states aids virus drug abuse food and drug administration aids patient centers for disease control heart disease drug testing

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