

# Contextual Document Clustering

Vladimir Dobrynin

St.Petersburg State University and SOPHIA Search Ltd.

25 December 2009

# Outline

- 1 Contextual Document Clustering
- 2 Discussion
- 3 US patents abstracts. IPC G06: computing; calculating; counting
- 4 References

# Context

Compare contexts of these words and phrases:

- **computer**

It is hard to predict what it is about – text about computing or about health, art, sport, etc.

- **motherboard**

With high probability this text should be about computer hardware.

# Simple Ideas

- All authors (and texts by these authors) can be grouped into **communities** in such a way that members of one community share the same background in terms of education, profession, experience, interests and so on. They use the same language to express the same ideas.
- Members of any specific community usually use few but **very specific terms** that aren't used in any of other communities. These terms can be considered as **markers** that distinguish texts written by members of this community from texts written by members of other communities.

# Notations

- $\mathcal{X}$  – set of all documents in the document corpus
- $\mathcal{Y}$  – set of all words occurring in documents from  $\mathcal{X}$
- $tf(x, y)$  – term frequency – the number of occurrence of the word  $y \in \mathcal{Y}$  in the document  $x \in \mathcal{X}$
- $p(Y|z)$  – context of word  $z$   
A probability distribution of a set of words which co-occur with a given word  $z$  in a document from  $\mathcal{X}$ .

Context of word  $z$ :

$$p(y|z) = \frac{\sum_{x \in D_z} tf(x, y)}{\sum_{x \in D_z, y'} tf(x, y')},$$

where  $D_z$  is the set of all documents that contain the word  $z$ .

# Narrow context

Selection of words with narrow contexts is based on consideration of

- 1 the entropy

$$H(Y|z) = - \sum_y p(y|z) \log p(y|z)$$

- 2 and document frequency

$$df(z) = |\{x : tf(x, z) > 0\}|.$$

Main idea – given collection  $\mathcal{X}$ , extract set  $\mathcal{Z}$  of words with narrow contexts and use contexts of these words  $p(Y|z)$ ,  $z \in \mathcal{Z}$  as cluster attractors.

# Entropy and Context Narrowness

Let

$$p(y|z) = \frac{1}{|T(D_z)|}, \quad y \in T(D_z)$$

where  $T(D_z)$  is the set of words belonging to documents where  $z$  occurs. Then

$$H(p(Y|z)) = \log |T(D_z)|.$$

For any non-uniform distribution

$$H(p(Y|z)) < \log |T(D_z)|.$$

Small  $H(Y|z)$  means that context of the word  $z$  can be described by a relatively small set of words (concepts)

# Entropy and Document Frequency

Heaps law:

$$|\mathcal{Y}| = O(n^\beta)$$

where constant  $\beta < 1$  and  $n$  is size of collection in words.

Let  $n(z)$  – size of  $D_z$  in words and  $l_{avg}$  – average document size in  $\mathcal{X}$ . Then

$$n(z) \approx l_{avg} \cdot df(z)$$

and

$$\begin{aligned} |T(D_z)| &= O(n(z)^\beta) = O((l_{avg} \cdot df(z))^\beta) = \\ &= O(df(z)^\beta). \end{aligned}$$

Hence

$$\begin{aligned} H(Y|z) &= O(\log |T(D_z)|) = c + O(\log df(z)^\beta) = \\ &= c + O(\beta \log df(z)) = c + O(\log df(z)). \end{aligned}$$



# Contextual Document Clustering -1-

## Input

- Term frequencies  $tf(x, y)$ ,  $x \in \mathcal{X}$ ,  $y \in \mathcal{Y}$
- $N$  – number of clusters
- Parameters  $df_{min}, r, \alpha > 1$

## Output

- "Hard" clustering of all documents into  $N$  clusters

# Contextual Document Clustering -2-

## Context calculation

$$p(y|z) = \frac{\sum_{x \in D_z} tf(x, y)}{\sum_{x \in D_{z, y'}} tf(x, y')}, \quad z \in \mathcal{Y}$$

## Set $\mathcal{Z}$ of words with narrow contexts selection

For every  $i = 1, \dots, r$

$$\mathcal{Y}_i = \{z : df_i \leq df(z) < df_{i+1}\},$$

$$df_1 = df_{min}, \quad df_{i+1} = \alpha \cdot df_i,$$

$$\mathcal{Z}_i \subseteq \mathcal{Y}_i, \quad |\mathcal{Z}_i| = \frac{N \cdot |\mathcal{Y}_i|}{\sum_{j=1, r} |\mathcal{Y}_j|},$$

$$z_1 \in \mathcal{Z}_i, \quad z_2 \in \mathcal{Y}_i - \mathcal{Z}_i \rightarrow H(Y|z_1) \leq H(Y|z_2).$$

$$\mathcal{Z} = \bigcup_i \mathcal{Z}_i.$$

# Contextual Document Clustering -3-

## Document clustering

For every document  $x \in \mathcal{X}$  calculate word probability distribution

$$p(y|x) = \frac{tf(x, y)}{\sum_{y'} tf(x, y')}, y \in \mathcal{Y}$$

and document  $x$  will be assigned to the cluster with centroid  $p(Y|z)$  if

$$z = \operatorname{argmin}_{z'} JS_{0.5,0.5}[p(Y|z'), p(Y|x)]$$

where  $JS_{0.5,0.5}[p, q] = H(\frac{p+q}{2}) - 0.5H(p) - 0.5H(q)$  is Jensen-Shannon divergence of probability distributions  $p$  and  $q$ .

# Why we use $r$ and $\alpha$ parameters?

- If  $z \in \mathcal{Z}_i$ ,  $i = 1, \dots, r$  then

$$H(Y|z) = O(\log |T(D_z)|) = c + O(\log df(z))$$

$$= c + O(\log(df_{min}\alpha^i)) = c + O(i)$$

- Given a theme (Math) we would like to disclose all *narrow* topics of the theme (equations, algebra,...) presented by significant number of documents in  $\mathcal{X}$  and accumulate all other Math documents into a more *broad* topic.

# Complexity

CDC complexity is

$$O(K \cdot |S|),$$

where  $K$  is the number of clusters and  $S$  is the set of non-zero elements in document-term matrix. By the way, complexity of K-means algorithm is  $O(t \cdot K \cdot |S|)$  where  $t$  is the number of iterations.

# Clustering principle of CDC

Split document corpus into relatively large groups of documents that are covered by relatively small number of concepts.

# Comparing with K-means. Notations

- Document collection  $D$ ,  $|D| = N$
- Every document is represented by
  - Index  $n = 1, \dots, N$
  - Feature vector  $\mathbf{x}^{(n)} \in \mathbb{R}^M$  where  $M$  is a number of features
- Every cluster is represented by
  - Index  $k = 1, \dots, K$
  - centroid  $\mathbf{m}^{(k)} \in \mathbb{R}^M$  – average of feature vectors representing all documents from the cluster
- $d(\mathbf{x}, \mathbf{y})$  is a distance measure between two vectors  
 $\mathbf{x}, \mathbf{y} \in \mathbb{R}^M$

# Comparing with K-means. "Soft" K-means

## 1. Initialization

- Set centroids  $\{\mathbf{m}^{(k)}\}_{k=1,\dots,K}$  to random values.
- Set a value for parameter  $\beta$ .

## 2. Assignment step

Responsibility  $r_k^{(n)}$  determines degree to which  $\mathbf{x}^{(n)}$  is assigned to cluster  $k$ :

$$r_k^{(n)} = \frac{\exp(-\beta d(\mathbf{m}^{(k)}, \mathbf{x}^{(n)}))}{\sum_{k'} \exp(-\beta d(\mathbf{m}^{(k')}, \mathbf{x}^{(n)}))}$$

## 3. Update step

$$\mathbf{m}^{(k)} = \frac{\sum_n r_k^{(n)} \mathbf{x}^{(n)}}{R^{(k)}},$$

where

$$R^{(k)} = \sum_n r_k^{(n)}.$$



# Comparing with K-means. Mixture of two Gaussians

## -1-

Points from  $\{x_n\}_{n=1}^N$ ,  $x_n \in \Re$  are distributed according to mixture of two Gaussians:

$$P(x|\mu_1, \mu_2, \sigma) = \sum_{k=1}^2 p_k \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x - \mu_k)^2}{2\sigma^2}\right),$$

where  $p_k = 0.5$ .

# Comparing with K-means. Mixture of two Gaussians

## -2-

Given standard deviation  $\sigma$  we can find means  $\mu_1$  and  $\mu_2$  using iterations:

$$\mu_k = \frac{\sum_{n=1}^N \rho_k^{(n)} x_n}{\sum_{n=1}^N \rho_k^{(n)}}, \quad k = 1, 2,$$

$$\rho_1^{(n)} = \frac{1}{1 + \exp[-(\omega_1 x_n + \omega_2)]},$$

$$\rho_2^{(n)} = \frac{1}{1 + \exp[+(\omega_1 x_n + \omega_2)]},$$

$$\omega_1 = \frac{\mu_1 - \mu_2}{\sigma^2}, \quad \omega_2 = \frac{\mu_2^2 - \mu_1^2}{2\sigma^2}.$$

If in "soft" K-means we have  $d(x, y) = \frac{1}{2}(x - y)^2$  and  $\beta = \frac{1}{\sigma^2}$  then responsibility  $r_k^{(n)} = \rho_k^{(n)}$  and hence in this case "soft"

# Problems with K-means

This analysis shows that using k-means clustering we implicitly suppose (at least in case of special type of distance function) that:

- 1 points in a cluster are generated by a Gaussian generator;
- 2 size (number of points) of every cluster is the same ( $p_k = \frac{1}{K}$ , where  $K$  is the number of clusters);
- 3 every cluster is spherical in shape and has the same diameter (same  $\sigma$  for all generators).
- 4 initialization? number of iterations (complexity)?

# Clustering principle of K-means

Depends of choice of distance measure  $d(x, y)$ .

# Comparing with Information Bottleneck. Notations

- $X$  – discrete random variable with values from  $\mathcal{X}$  – documents
- $p(x)$  – probability distribution of  $X$  – is proportional to size of document  $x$
- $T$  – discrete random variable with values from  $\mathcal{T}$  – clusters
- $p(t|x)$  – "soft" clustering documents from  $\mathcal{X}$  into clusters  $\mathcal{T}$
- $\mathcal{Y}$  – set of documents features – words
- $Y$  – random variable with values in set  $\mathcal{Y}$
- $p(x, y)$ ,  $x \in \mathcal{X}$ ,  $y \in \mathcal{Y}$  – known joint probability distribution

$$p(x, y) = \frac{tf(x, y)}{\sum_{x', y'} tf(x', y')},$$

where  $tf(x, y)$  – a number of times term  $t$  occurs in document  $x$ .

# Comparing with Information Bottleneck. Optimization Criterion

Minimize

$$\mathcal{L}[p(t|x)] \equiv I(X; T) - \beta I(T; Y)$$

where Lagrange multiplier  $\beta > 0$ .

- Minimization of mutual information  $I(X; T)$  – maximization of **degree of compression** of information presented in document corpus.
- Maximization of mutual information  $I(T; Y)$  – all documents from the same cluster have the same features – in other words they are **similar to each another**.
- If  $\beta = 0$  then all documents are assigned to one cluster.
- If  $\beta \rightarrow \infty$  then we have one cluster per document.

# Comparing with Information Bottleneck. Optimal Clustering

## Theorem

Probability distribution  $p(t|x)$  that minimizes functional  $\mathcal{L}[p(t|x)]$  can be presented in the following form

$$p(t|x) = \frac{p(t)}{Z(x, \beta)} e^{-\beta D_{KL}[p(y|x) || p(y|t)]}, \quad t \in \mathcal{T}, y \in \mathcal{Y}.$$

Here  $Z(x, \beta)$  is a normalizer.

# Problems with Information Bottleneck

- How to select reasonable value for  $\beta$ ? This parameter should reflect user's idea about "good" clustering
- Number of clusters ( $|\mathcal{T}|$ ) should be less than 200 (Ron Bekkerman. Private communication)



# Clustering principle of Information Bottleneck

Depends on choice of parameter  $\beta$  value.

# Database statistics

Number of documents	553,792
Number of words	57,096
Number of clusters	1,074
Max cluster size	4,170

# Top 10 clusters

Context word	Cluster size	Stem	Stem document frequency	Stem rank
sdram	4170	sdram	368	7
gradation	3827	gradat	584	7
pll	3520	pll	298	7
epg	3246	epg	409	7
reproducing	2953	reproduc	4272	6
inode	2922	inod	120	7
uninstall	2910	uninstal	119	7
spooler	2696	spooler	144	7
mining	2694	mine	1434	6
metrology	2653	metrolog	319	7

# Top 10 words from context "mining"

Word	Stem	Probability
data	data	0.0616228538861579
mining	mine	0.0319750184037424
system	system	0.02626392154069
method	method	0.0210990002612143
information	inform	0.0132506945928617
model	model	0.0128944931253117
database	databas	0.0122652038659733
user	user	0.0108997649070314
includes	includ	0.0100805015316663
processing	process	0.00867944242596944

## Top 5 patents nearest to the context "mining"

Company	Date	Title
Oracle	2005-03-08	Data mining application programming interface
Oracle	2007-02-06	In-database clustering
Lucent Technologies	2002-05-07	System and method for analyzing and displaying telecommunications switch report output
IBM	2009-04-21	Computerized data mining system, method and program product
Oracle	2006-10-03	Enterprise web mining system and method

# Clusters for query "data mining"

Sophia Search - Windows Internet Explorer

http://services.sophiasearch.com/patents/cluster\_viewServlet?event=clickOnPage&pns=1&duration=8228

Google spiran Search + Sidewiki Bookmarks Check Translate AutoFill spiran Sign In

Demo Licence Dr Vladimir Dobrynin My Account My Sophia Sign Out Help

**Sophiasearch**  
deriving meaning from content

Keyword Search Search by Example

(G06) COMPUTING; CALCULATING; COUNTING; data mining Search in new category Search

Thematic Folders Execution time: 8.228 sec. Page 1 out of 19

- COMPUTERIZED METHOD, SYSTEM AND PROGRAM PRODUCT FOR GENERATING A DATA MINING MODEL** System and method for selecting a data mining modeling algorithm for data mining applications Computerized data mining system, method and program product

...system and program product for generating a data mining model... but which provides data mining functionality that is accessible to users having limited data mining expertise and which provides reductions in development times and costs for data mining projects... The present invention relates to a system and method for electronic and algorithmic data mining of an individual physician's prescribing history to determine the approximate distribution of diseases within their practice population for optimizing pharmaceutical sales and marketing...

This theme contains 662 key documents
- USER-CONTEXT-BASED SEARCH ENGINE** Method for extracting information utilizing a user-context-based search engine Locating, filtering, matching macro-context from indexed database for searching context where micro-context relevant to textual input by user

...A user may then query the tool for a desired type of information... The tool filters the database to provide a set of pinpoint site locations with information of the type requested in the query... The recommendation database is generated by performing the steps of: performing data mining using users search query logs...

This theme contains 22 key documents
- Patient Data Mining for Lung Cancer Screening** Early detection of disease outbreak using electronic patient data to reduce public health threat from bioterrorism Patient data mining for quality adherence

...The system includes a database including structured patient information for a patient population and a domain knowledge base including information about lung cancer... A list of persons for whom consent was obtained can be outputted and forwarded to the entity interested in performing the clinical trial and which requested the list... A method for the construction and utilization of a medical records system capable of providing a continuous data stream of epidemiological data to the records system via kits provided to the symptomatic population to obtain and record an epidemiological profile in a searchable database by applying data mining or automated intelligence techniques whereby...

This theme contains 20 key documents
- Investigating business processes** System and method of real estate data analysis and display to support business management Combining multidimensional expressions and data mining extensions to mine OLAP cubes

...A graphical user interface enables users to apply data warehousing and data mining techniques to business process execution data and to visualize process execution data along multiple configurable dimensions and at different configurable levels of granularity... A suite of software tools performs data mining on data periodically extracted from one or more Multiple Listing Service (MLS) databases to support the management of a real estate business...

Done Sophia Search - Windows Internet Explorer Internet | Protected Mode: On 100%

PNG... Sophia... Skype... 2 Wi... SPIRA... Visual... Micros... Docu... Micros... Untitle... EN 8:33 PM

# Patents from cluster 2

The screenshot displays the Sophia Search web application in a Windows Internet Explorer browser. The address bar shows the URL: [http://services.sophiasearch.com/patents/cluster\\_view\\_servlet?event=clickOnCluster&cid=251](http://services.sophiasearch.com/patents/cluster_view_servlet?event=clickOnCluster&cid=251). The page features a search bar with the text "(G06) COMPUTING; CALCULATING; COUNTING; data mining" and a "Search" button. Below the search bar, the "Documents View" section displays a list of five patent results under the heading "Key Documents". The results are numbered 1 through 5 and include titles, source links, and dates. The first result is "USER-CONTEXT-BASED SEARCH ENGINE [Source] [01/11/2007]". The second result is "Method for extracting information utilizing a user-context-based search engine [Source] [15/05/2007]". The third result is "Locating, filtering, matching macro-context from indexed database for searching context where micro-context relevant to textual input by user [Source] [14/03/2008]". The fourth result is "System and method for search and recommendation based on usage mining [Source] [15/08/2006]". The fifth result is "SYSTEM FOR EVOLVING EFFICIENT COMMUNICATION [Source] [19/02/2009]". The page also includes a "Documents View" sidebar on the left with a list of terms: database queries, database to search, modified query, program product, queries database, query responses, query result, query search, query to request, retrieving the document, search engines, search queries, search result, search terms. The bottom of the page shows the Windows taskbar with various application icons and the system clock indicating 7:45 PM.

Sophia Search - Windows Internet Explorer  
http://services.sophiasearch.com/patents/cluster\_view\_servlet?event=clickOnCluster&cid=251  
spiiran

Google Search by Example  
(G06) COMPUTING; CALCULATING; COUNTING; data mining Search in new category Search

Back to Themes

Dr Vladimir Dobrynin My Account My Sophia Sign Out Help

Sophiasearch deriving meaning from content

Documents View  
database queries, database to search, modified query, program product, queries database, query responses, query result, query search, query to request, retrieving the document, search engines, search queries, search result, search terms

Key Documents Related Documents

Execution time: 1.616 sec.

1 ☐ USER-CONTEXT-BASED SEARCH ENGINE [Source] [01/11/2007]  
...A data extraction tool is provided for cataloging information in an information source for searching by a user... Show neighbors

2 ☐ Method for extracting information utilizing a user-context-based search engine [Source] [15/05/2007]  
...A data extraction tool is provided for cataloging information in an information source for searching by a user... Show neighbors

3 ☐ Locating, filtering, matching macro-context from indexed database for searching context where micro-context relevant to textual input by user [Source] [14/03/2008]  
...A data extraction tool is provided for cataloging information in an information source for searching by a user... Show neighbors

4 ☐ System and method for search and recommendation based on usage mining [Source] [15/08/2006]  
...The method of searching comprises the steps of: receiving from a user a search query requesting information, retrieving at least one recommendation relating to the search query... Show neighbors

5 ☐ SYSTEM FOR EVOLVING EFFICIENT COMMUNICATION [Source] [19/02/2009]  
...The invention includes a database having one or more records... Show neighbors

Done Internet | Protected Mode: On 100% 7:45 PM

- Baeza-Yates and Ribeiro-Neto. *Morden Information Retrieval*. ACM Press, 1999
- Daniel Chandler. *Semiotics for Beginners*.  
<http://www.aber.ac.uk/media/Document/S4B>
- Thomas M. Cover and Joy A. Thomas. *Elements of Information Theory. Second Edition*. John Wiley & Sons, Inc., Hoboken, New Jersey, 2006
- Vladimir Dobrynin, David Patterson and Niall Rooney. Contextual Document Clustering. *Lecture Notes in Computer Science. Advances in Information Retrieval*, 2997,167–180, 2004
- David J.C. MacKay. *Information Theory, Inference, and Learning Algorithms*. Cambridge University Press, 2003



- N. Tishby, F. Pereira and W. Bialek. The Information Bottleneck Method. *Proc. 37th Allerton Conference on Communication and Computation*, 1999
- Noam Slonim. *The Information Bottleneck: Theory and Applications*. PhD thesis, Hebrew University, Israel, 2002