

Giving an Upperbound of the Number of Clusters and Relevant Words in Hierarchical Document Clustering Based on BIC

Yoshio Fukushige

fukushige.yoshio@jp.panasonic.com

Network Development Center, Matsushita Electric Industrial Co., Ltd.

Abstract

A new generative model based approach to automatic document clustering, using the BIC as the model selection criterion is described. A new method based on a graphical model is proposed to give an upperbound to the numbers of clusters and relevant words. The result of an experiment using the NTCIR web data collection is briefly reported.

Keywords: generative model, BIC

1. Introduction

1.1. Overview

When clustering documents, the number of the variables (dimension) is usually very large. In machine learning, high dimensionality is known to lead to high computational cost and low precision. Hence various approaches have been proposed to reduce the number of words used in the clustering; some are based on word frequency, some employ information gain, and some perform data transformation (e.g. LSI [2]). However those approaches lack the direct connection to the nature of clustering task.

And in the clustering problem, the number of the classes has to be estimated. Conventional clustering approaches solve this issue by fixing either the maximum number of the child clusters or minimum number of the elements in each cluster, without theoretical support.

Therefore the following issues are addressed in this work through a probabilistic approach:

- how to determine which words are relevant in clustering
- how to determine the maximum number of child clusters

The result of an experiment using the NTCIR data set is reported.

In this paper, as the problem setting, the generative model based clustering is quickly described in section 2 with a simple example

model. Section 3 describes a new model with word selection and section 4 describes the method for giving an upperbound of numbers of clusters. Section 5 describes the experiment with NTCIR data and section 6, the conclusion.

2. Generative model based clustering

2.1. Generative model based clustering

A generative model based clustering suppose a model with unknown parameters which can produce the given data set[3],[1]. Among those parameters are the "cluster indices" which denote which item belongs to which cluster. The clustering problem is then described as the parameter estimation problem or model selection problem. By the Bayes rule, we have:

$$P(\text{Model}|\text{Data}) \propto P(\text{Model})P(\text{Data}|\text{Model})$$
, so to maximize the posterior of the model (left hand), we could maximize the right hand instead. If we assume all models are equally probable in their prior, we can get the best model by searching for the best model which maximizes the likelihood $P(\text{Data}|\text{Model})$.

To balance the expressibility and generality, the BIC (Bayes Information Criterion)([5])below is often used as the melkmar.

$$BIC \equiv \log P(\text{Data}|\text{Model}) - \frac{\# \text{ of parameters}}{2} \log(\# \text{ of data})$$

(The base of the log in this paper is 2)

In that regard, in order to get "best" clustering result, one can search for the model which maximizes the BIC. Because the cluster indices are "hidden", one has to get the model through some indeterministic approach, for example, using the EM algorithm.

2.2. A simple example

Here we show a simple example of a generative model, which is called a bag-of-words model, where the whole data set is regarded as a collection of documents and each document is regarded as a bag of words (i.e. word order is neglected), and where each word token in the data set is regarded to have derived in the following manner:

choose the document d it appears in
 choose the cluster (= category) c it belongs to (a word token is considered to belong to a cluster)

choose its word type w (here, "word type" mean the string form of a word and "word token" means its occurrence in a document. We use hereafter "word" for "word type" when there is no confusion)

And we put an assumption that $p(w|d, c) = p(w|c)$

Therefore, "the probability of a word token to appear in document d , belong to cluster c , and have word type w " is $p(d)p(c|d)p(w|c)$

Figure 1 shows the graphical representation (called a Bayesian Network ([3])) of the model here. Each node represents a variable, and a directed edge between two nodes means a (direct) dependency between those variables. Here the node D corresponds to d (document), C to c , W to w . Each node has a CPT (conditional probability table), which describe the conditional probability of the corresponding variable given the variables corresponding to its parent nodes (node A is said to be a parent of node B when there is a directed edge heading to B). Then the joint probability of the data is

$$k! \prod_w \prod_c \prod_d \{p(d)p(c|d)p(w|c, d)\}^{N_{dwc}}$$

$$= k! \prod_d p(d)^{N_d} \cdot \prod_c \prod_w p(w|c)^{N_{wc}}$$

where k denotes the number of clusters, N_{dwc} denotes the number of tokens in document d , of word type w and belong to cluster c , and $N_d \equiv \sum_c \sum_w N_{dwc}$,

$$N_{wc} \equiv \sum_d N_{dwc}$$

($k!$ comes from the interchangeability of the hidden variable c 's)

$$\text{With } D \equiv |d|, \quad M \equiv |w|,$$

$$N \equiv \sum_d \sum_c \sum_w N_{dwc}$$

the number of the variables is $D + D(k-1) + k(M-1) = k(D+M-1)$, yielding

$$BIC = \sum_{i=1}^k \log i + \sum_d N_d \log p(d) + \sum_c \sum_w N_{wc} \log p(w|c) - \frac{k(D+M-1)}{2} \log N$$

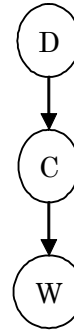


Figure 1 A simple generative model

2.3. Hierarchical clustering

Hierarchical clustering approaches are often claimed to lead to a good result. Our clustering method is also a hierarchical one. However, it is a greedy one, i.e. after choosing the best clustering in a certain level, say $l-1$, we proceed to search for the best clustering of the child clusters in level l . Some other approaches make the set of whole hierarchical models first and then chose the best from them ([6]).

3. A new model with word selection

Here we introduce a new generative model, as an extension to the simple model in 2.2, where a new variable r (relevance flag) which denotes whether a word token is relevant in clustering ($r=1$) or not ($r=0$). Each word token in the data set is regarded to have derived in the following manner:

1. choose the document d it appears in
2. choose the cluster c it belongs to (a word token is considered to belong to a cluster)
3. choose whether it is one of the relevant words (word types) to the clustering
4. choose its word type w

And we assume
 $p(r|d,c) = p(r|c)$, $p(w|d,c,r) = p(w|c,r)$.

Therefore, "the probability of a word token to appear in document d , belong to cluster c , have relevance flag r , and have word type w " is: $p(d)p(c|d)p(r|c)p(w|c,r)$

Figure 2 shows the graphical representation of the model here.

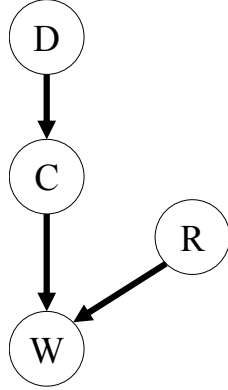


Figure 2: An extended graphical model

We further assume that word tokens of the same word type are either all relevant (saying the word type is relevant) or irrelevant (saying the word type is irrelevant), and w is independent from c given $r=0$.

Here we introduce some additional notations to those in 2.2:

$$W_1 \equiv \{w | p(r=1|w) = 1\}$$

$$W_0 \equiv \{w | p(r=0|w) = 1\}$$

$$N_{W_1c} \equiv \sum_{w \in W_1} N_{wc}$$

$$N_{W_0c} \equiv \sum_{w \in W_0} N_{wc}$$

$$N_c \equiv \sum_w N_{wc}$$

$$N_w \equiv \sum_c N_{wc}$$

$$N_{W_0} \equiv \sum_{w \in W_0} N_w$$

$$M_1 \equiv |W_1|$$

$$M_0 \equiv |W_0|$$

The probability of our having the data set is then,

$$k! \prod_d p(d)^{N_d} \cdot p(W_1)^{N_{W_1}} \cdot p(W_0)^{N_{W_0}} \\ \cdot \prod_c \prod_{w \in W_1} p(w|c, r=1)^{N_{wc}} \cdot \prod_c \prod_{w \in W_0} p(w|c, r=0)^{N_{wc}}$$

Here, the number of the variables is

$$D - 1 + D(k - 1) + 1 + k(M_1 - 1) + (M_0 - 1) \\ = k(D + M_1 - 1) + M_0 - 1$$

therefore

$$BIC = \sum_{i=1}^k \log i + \sum_d N_d \log p(d) \\ + N_{W_1} \log p(W_1) \\ + \sum_c \sum_{w \in W_1} N_{wc} \log p(w|c) \\ - \sum_c N_{W_1c} \log p(W_1|c) \\ + \sum_{w \in W_0} N_w \log p(w) \\ - \frac{k(D + M_1 - 1) + M_0 - 1}{2} \log N$$

When we turn a relevant word w^* to irrelevant, the BIC will be

$$BIC^* = \sum_{i=1}^k \log i + \sum_d N_d \log p(d) \\ + N_{W_1'} \log p(W_1') \\ + \sum_c \sum_{w \in W_1'} N_{wc} \log p(w|c) \\ - \sum_c N_{W_1'c} \log p(W_1'|c) \\ + \sum_{w \in W_0 \cup \{w^*\}} N_w \log p(w) \\ - \frac{k(D + M_1 - 2) + M_0}{2} \log N$$

where $W_1' \equiv W_1 - \{w^*\}$,
and the gain is

$$\begin{aligned} \Delta BIC &= N_{w'_1} \log p(W'_1) - N_{w_1} \log p(W_1) \\ &\quad - \sum_c N_{w^*_c} \log p(w^*|c) \\ &\quad - \sum_c N_{w'_1 c} \log p(W'_1|c) \\ &\quad + \sum_c N_{w_1 c} \log p(W_1|c) \\ &\quad + N_{w^*} \log p(w^*) + \frac{k-1}{2} \log N \end{aligned}$$

By substituting N_{w^*}/N for $p(w^*)$ and alike,

$$\begin{aligned} \Delta BIC &= \sum_c N_{w_1 c} \left[-\frac{N_{w'_1 c}}{N_{w_1 c}} \log \frac{N_{w'_1 c}}{N_{w_1 c}} \right. \\ &\quad \left. - \frac{N_{w^*_c}}{N_{w_1 c}} \log \frac{N_{w^*_c}}{N_{w_1 c}} \right] \\ &\quad + N_{w'_1} \log \frac{N_{w'_1}}{N_{w_1}} + N_{w^*} \log \frac{N_{w^*}}{N_{w_1}} \\ &\quad + \frac{k-1}{2} \log N \\ &\geq \sum_c \sum_{d \in c} N_{w_1 d} \left[-\frac{N_{w'_1 d}}{N_{w_1 d}} \log \frac{N_{w'_1 d}}{N_{w_1 d}} \right. \\ &\quad \left. - \frac{N_{w^*_d}}{N_{w_1 d}} \log \frac{N_{w^*_d}}{N_{w_1 d}} \right] \\ &\quad + N_{w'_1} \log \frac{N_{w'_1}}{N_{w_1}} \\ &\quad + N_{w^*} \log \frac{N_{w^*}}{N_{w_1}} \\ &\quad + \frac{k-1}{2} \log N \\ &= \sum_d N_{w_1 d} \left[-\frac{N_{w'_1 d}}{N_{w_1 d}} \log \frac{N_{w'_1 d}}{N_{w_1 d}} \right. \\ &\quad \left. - \frac{N_{w^*_d}}{N_{w_1 d}} \log \frac{N_{w^*_d}}{N_{w_1 d}} \right] \\ &\quad + N_{w'_1} \log \frac{N_{w'_1}}{N_{w_1}} + N_{w^*} \log \frac{N_{w^*}}{N_{w_1}} \\ &\quad + \frac{k-1}{2} \log N \end{aligned}$$

where

$$\begin{aligned} N_{wd} &\equiv \sum_c N_{wdc} , \\ N_{w_1 d} &\equiv \sum_{w \in W_1} N_{wd} , \\ N_{w'_1 d} &\equiv \sum_{w \in W'_1} N_{wd} , \\ N_w &\equiv \sum_d N_{wd} , \\ N_{w_1} &\equiv \sum_d N_{w_1 d} , \\ N_{w'_1} &\equiv \sum_d N_{w'_1 d} \end{aligned}$$

Therefore word w^* can be turned to be irrelevant regardless of the clustering result, if $k \geq 1$

$$\begin{aligned} &+ \frac{2}{\log N} \left[N_{w_1} \left\{ -\frac{N_{w'_1}}{N_{w_1}} \log \frac{N_{w'_1}}{N_{w_1}} - \frac{N_{w^*}}{N_{w_1}} \log \frac{N_{w^*}}{N_{w_1}} \right\} \right. \\ &\quad \left. - \sum_d N_{w_1 d} \left\{ -\frac{N_{w'_1 d}}{N_{w_1 d}} \log \frac{N_{w'_1 d}}{N_{w_1 d}} - \frac{N_{w^*_d}}{N_{w_1 d}} \log \frac{N_{w^*_d}}{N_{w_1 d}} \right\} \right] \end{aligned}$$

4. Giving an upperbound of numbers of clusters

Now we address the problem of giving an upperbound of numbers of clusters.

When there is no division, the variable r is not introduced and

$$\begin{aligned} BIC &= \sum_d N_d \log p(d) + \sum_w N_w \log p(w) \\ &\quad - \frac{D+M-2}{2} \log N \end{aligned}$$

Therefore, the gain of BIC by dividing the cluster is

$$\begin{aligned} \Delta BIC &= \sum_{i=1}^k \log i + N_{w_1} \log N_{w_1} \\ &\quad - \sum_{w \in W_1} N_w \log N_w \\ &\quad - \frac{(k-1)(D+M_1-1)}{2} \log N \\ &\quad + \sum_c \sum_{w \in W_1} N_{wc} \log \frac{N_{wc}}{N_{w_1 c}} \end{aligned}$$

What is wanted here is the minimum k s.t. "the gain of BIC is negative when the cluster is

divided into $(k+1)$ clusters no matter what the relevant word set is and the clustering is." Let us denote this k by k^* . Then

$$\begin{aligned}
 k^* &\leq \arg \min_k \forall \mathbf{c} \forall W_1 \left[\sum_{i=1}^{k+1} \log i \right. \\
 &\quad + N_{W_1} \log N_{W_1} - \sum_{w \in W_1} N_w \log N_w \\
 &\quad - \frac{k(D+M_1-1)}{2} \log N \\
 &\quad \left. - \sum_c \sum_{d \in c} N_{W_1 d} \left(- \sum_{w \in W_1} \frac{N_{wd}}{N_{W_1 d}} \log \frac{N_{wd}}{N_{W_1 d}} \right) \right. \\
 &\quad \left. < 0 \right] \\
 &= \arg \min_k \forall W_1 \left[\sum_{i=1}^{k+1} \log i \right. \\
 &\quad + N_{W_1} \log N_{W_1} - \sum_{w \in W_1} N_w \log N_w \\
 &\quad - \frac{k(D+M_1-1)}{2} \log N \\
 &\quad \left. + \sum_d \sum_{w \in W_1} N_{wd} \log N_{wd} - \sum_d N_{W_1 d} \log N_{W_1 d} \right. \\
 &\quad \left. < 0 \right]
 \end{aligned}$$

where \mathbf{c} denotes a clustering (= set of clusters)

When we introduce the function f , as

$$\begin{aligned}
 f(W_1, k) &\equiv \sum_{i=1}^{k+1} \log i + N_{W_1} \log N_{W_1} \\
 &\quad - \sum_{w \in W_1} N_w \log N_w \\
 &\quad - \frac{k(D+M_1-1)}{2} \log N \\
 &\quad - N_{W_1} \min_d \left(- \sum_{w \in W_1} \frac{N_{wd}}{N_{W_1 d}} \right. \\
 &\quad \left. \log \frac{N_{wd}}{N_{W_1 d}} \right)
 \end{aligned}$$

then

$$k^* \leq \arg \max_k \exists W_1 [f(W_1, k) \geq 0]$$

Letting $W_1(k)$ be the subset of W which doesn't include any w^* 's which satisfies

$$f(W_1', k) - f(W_1, k) \geq 0,$$

from $\forall W_1 f(W_1, k) \leq f(W_1(k), k)$, we have

$$\begin{aligned}
 k^* &\leq \arg \max_k \left[\sum_{i=1}^k \log i + N_{W_1(k)} \log N_{W_1(k)} \right. \\
 &\quad - \sum_{w \in W_1(k)} N_w \log N_w \\
 &\quad - \frac{(k-1)(D+M_1(k)-1)}{2} \log N \\
 &\quad \left. - N_{W_1} \min_d \left(- \sum_{w \in W_1(k)} \frac{N_{wd}}{N_{W_1(k)d}} \right) \geq 0 \right]
 \end{aligned}$$

where

$$M_1(k) \equiv |W_1(k)|$$

and

$$N_{W_1(k)d} \equiv \sum_{w \in W_1(k)} N_{wd}$$

5. Clustering with an EM algorithm

After choosing relevant words and getting the upperbound of number of the clusters, an EM algorithm based on the 3 node model described in 2.2 is applied to the clustering of the documents in the parent cluster.

Here, for the multinomial distributions, we adopt the Jeffrey's noninformative priors (*Dirichlet* ($\theta|0.5, \dots, 0.5$)), which yields the following EM steps:

$$\begin{aligned}
 \text{E step} &\quad p(c|d, w) \propto p(w|c)p(c|d) \\
 \text{M step} &\quad p(w|c) \propto \sum_d (N_{wd} + 0.5)p(c|d, w) \\
 &\quad p(c|d) \propto \sum_w (N_{wd} + 0.5)p(c|d, w)
 \end{aligned}$$

The steps are iterated until the sum of the difference of the $p(c|d)$ gets smaller than a fixed value. Then each document d is classified to the cluster c with the highest $p(c|d)$.

Because the result of the trial depends on the initial value at the beginning of the EM steps, trials are made with different initializations and the result of the trial with the highest BIC is adopted as the result of the clustering. In initialization, a document for each child cluster is selected at random as the seed for the cluster.

6. An experiment with NTCIR data

6.1. NTCIR Topical Classification Task

An experiment is carried through participating the NTCIR Topical Classification Task.

Participants are given 47 queries and a set of documents for each query. The document set for each query consists of 200 documents from NTCIR's NW100G-01 data set. Participants are then required to cluster documents in each set into arbitrary numbers of clusters.

6.2. Algorithm

The following is the outline of our algorithm applied to the task.

1. Extract words using a dictionary and counting their occurrences in each document.
2. Make the root cluster which includes all the documents and put it in the cluster queue.
3. Repeat while the cluster queue is not empty:
 - (a) Pick up the first cluster in the queue and name it the current cluster.
 - (b) Get the relevant words to the division of the current cluster into each number of children clusters
 - (c) Get the upperbound of the numbers of clusters (k_c^{\max}), where c is the current cluster. Go to next cluster if $k_c^{\max} < 2$
 - (d) Get the BIC of the current cluster without division ($\equiv BIC^{\max} = BIC_0$).
 - (e) For each N from 2 to k_c^{\max} , repeat until a stop condition is met
 - ① Get N child clusters by estimating the $p(c'|c, d)$'s where d is each document and c' is a child cluster of c using the EM algorithm. (Note: clusters with less than 5 documents are discarded and documents are re-classified into the rest of the child clusters)
 - ② Calculate the BIC for the resultant model (BIC_1).

- ③ if $BIC_1 > BIC_0$, (re)mark this result as the best clustering and $BIC^{\max} \leftarrow BIC_1$
 - (f) If $BIC^{\max} > BIC_0$, put the child clusters in the best clustering to the top of the cluster queue, else mark the current cluster undividable.
4. Output the resultant leaf clusters (= clusters without children) with relevant key words using a KL-based score

6.3. Heuristics for early quitting in division trials

In searching the "best" clustering given the number of the clusters, we applied some heuristics for early stopping.

First, we regard the process of dividing the parent cluster as a Bernoulli trial with $Beta(\theta|0.5, 0.5)$ as the prior for the success probability of the binomial distribution.

The expected success probability after n trials with no success is then $0.5/(n+1)$. Here "success" means that we found a division of the parent cluster with higher BIC than that for non-division model. The smallest n such that the probability is lower than 0.05 is 10, therefore, we quit after 10 trials with no success.

When we found a division with higher BIC than that of non-division model, the success probability in n trials to get a division with still higher BIC than the highest in the series of trials is, then $1.5/(n+1)$. The smallest n such that the probability is lower than 0.05 is 30, therefore, we quit after 30 trials to conclude that we won't get a division model with higher BIC.

Further when new division is got, the mean and standard deviation of the BICs are calculated and if the best BIC for the division of the parent cluster gets to be found far from the mean by a fixed number times the standard deviation, the trial gets quit. The coefficients for the standard deviation are fixed for each number of trials (see discussion).

When the BIC is not updated while in the trials with k clusters, trials with $k+1$ or more clusters are skipped.

6.4. Key word selection for the clusters

For each cluster c of documents, 20 words with smallest Kullback-Leibler divergence $KL(w, c)$ described below are chosen as the keywords representing the cluster.

$$KL(w, c) \equiv \sum_d p(w|d) \log \frac{p(w|d)}{p(c|d)}$$

where $p(w|d) = \frac{N_{wd} + 0.5}{N_d + 0.5M}$

Intuitively, the words whose distributions over the documents are most similar to those of the cluster are chosen.

Here again, a Jeffery's noninformational prior is adopted.

6.5. Overview of the result

By now the results for 11 queries are evaluated by the NTCIR staffs and the evaluation results are given to each participant.

The numbers of the relevant words to each clustering varies from 0 to 39,480 with the mean being 862.1 and standard deviation being 2740. Figure 3 shows the histogram of numbers of relevant words.

The numbers of the children turned out to be considerably smaller than the upperbounds above during preliminary experiments and a heuristic rule was employed to quickly abort the trial. Figure 4 shows the numbers of children found.

The number of the leaf clusters vary from 14 to 27, with the mean being 22.8 and standard deviation being 3.9. Figure 5 shows the numbers of leaf clusters for each query.

The class size (= number of documents in the class) varies from 5 to 61, with the mean being 7.98 and the standard deviation being 6.26. Figure 6 is the histogram of leaf cluster sizes. The minimum 5 is derived from a heuristic restriction setting of the minimum size to 5. Figure 7 is the box and whisker graph of cluster sizes for each query.

The result of the official evaluation (average base) by NTCIR standards are listed in Table 1.

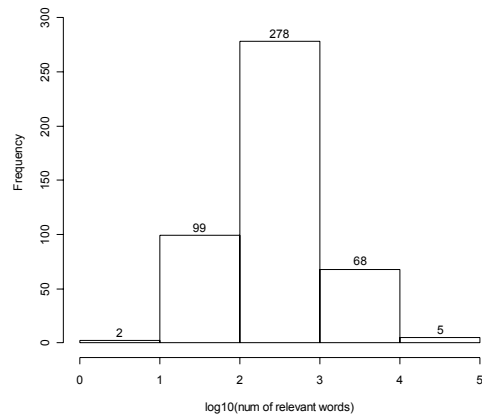


Figure 3: Numbers of relevant words

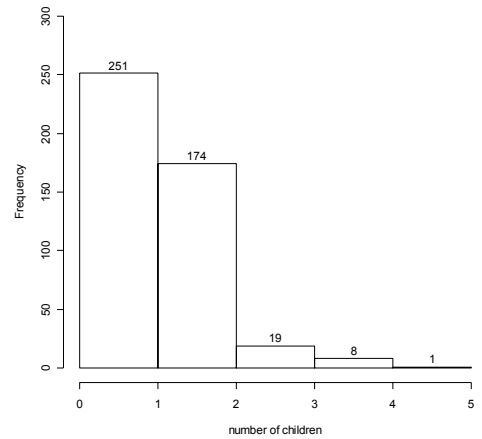


Figure 4 Histogram of numbers of children

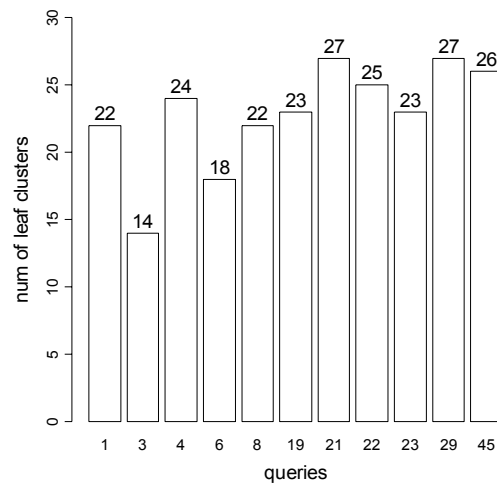


Figure 5 Number of leaf classes

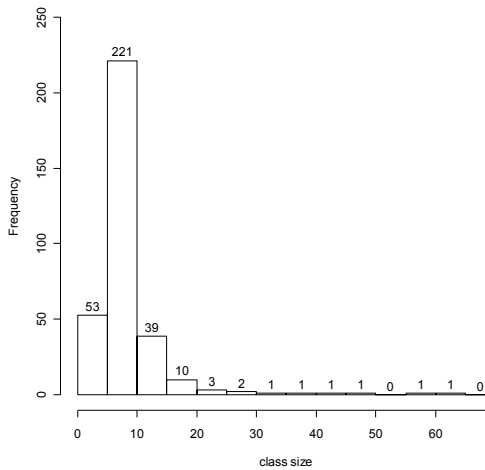


Figure 6 Histogram of (leaf) cluster sizes

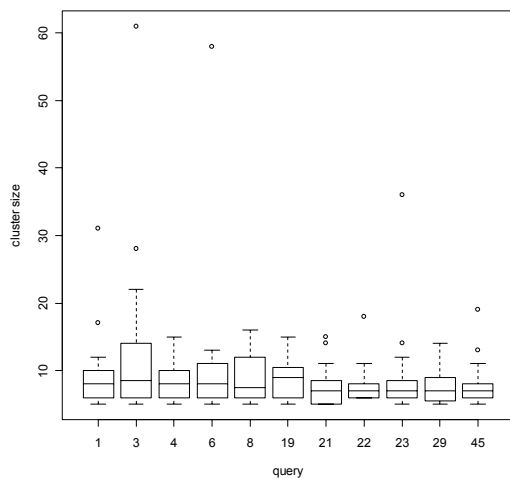


Figure 7 Cluster sizes for each query

Table 1 Evaluation results by NTCIR standards (average)

	relax	rigid
Avg. Precision	0.21	0.22
Precision @20	0.40	0.35
Recall @20	0.48	0.53
F-value@20	0.30	0.27
cg	8.09	6.91
dcg	3.01	2.64
mdcg1	2.80	2.45
mdcg2	7.25	6.14

7. Conclusion and Discussion

A new generative model based hierarchical clustering method is proposed, where the relevant words and upperbounds of number of clusters are calculated by introducing a class indices and using the BIC for model selection.

An EM algorithm is applied in estimating the parameters, with Jeffery's noninformational priors for smoothing.

Some heuristic stop conditions are used in the trials of cluster devision. Unfortunately, the stopping heuristics applied here was not theoretically justified enough. Especially the evaluation of the best BIC based on the mean and standard deviation was found to tend to underestimate the deviation of the data, and tend to cause premature stop, which means the possibility of missing better clustering with more child clusters.

Future work includes

More precise upperbounding of the number of the clusters. The upperbound here is just a rough upperbound and actual number of the child clusters were lower than the upperbound in the experimet.

Consideration of dependence between word types. Here, all word types are treated as independent from each other, which is not actually the case.

Use of variational methods for parameter estimation.

Comparison of BIC to other information criteria, such as AIC, stochastic complexity, etc.

References

- [1] D.M. Blei, A.Y. Ng and M.I. Jordan. Latent Dirichlet Allocation, *NIPS 14*, 2003.
- [2] S. Deerwester, S.T. Dumais, G.W. Furnas, T.K. Landauer, and R. Harshman. Indexing by latent semantic indexing. *Journal of the American Society for Information Science*, 41 (6), pp.391-407, 1990.
- [3] T. Hofmann. Probabilistic latent semantic indexing. *Proceedings of the Twenty-Second Annual International SIGIR Conference*, pp. 50-57, 1999.
- [4] J. Pearl: *Probabilistic Reasoning in Intelligent Systems*, Morgan Kaufmann, 1988.
- [5] G. Schwartz: Estimating the dimension of a model, *The Annals of Statistics*, 6, pp.461-464, 1978.
- [6] A. Vinokov and M. Girolami. A Probabilistic Framework for the Hierarchical Organisation and Classification of Document Collections, *Journal of Intelligent Information Systems*, 18 (2), pp.153-172, 2002.