

---

# Bayesian Clustering and Topic Discovery: Adventures with Gene Expression Data

---

Karren Dai Yang, Skanda Koppula

{KARREN, SKOPPULA}@MIT.EDU

Massachusetts Institute of Technology, Cambridge, MA 02139 USA

## 1. Introduction

Tumors cell lines are composed of different sub-populations of cells which often exhibit shared patterns of gene expression. Biologists are interested in two key questions: given gene expressions values, (1) can we identify cell clusterings, and (2) can we identify clusters of biologically-related genes? Our goal with this project was to answer both these two questions using Bayesian methods.

In brief, we explored the use of three Bayesian clustering methods (a vanilla mixture, integrated topic-mixture, and non-parametric models) to address the first question, and two topic models (vanilla and a dynamic-topic LDA) to address the second.

### 1.1. Description of Data

Using a contemporary gene sequencing machine, we obtained samples of single-cell RNA-sequencing data (scRNA-seq) taken from tumors in mice. Our data consisted of the expression values of 22,712 genes for each of 4645 cells. In total, the dataset amounted to 0.86 GB, presenting significant computational challenges when attempting posterior estimation. Apart from a few computational tricks (online inference, multi-core parallelization), for most methods, we pruned non-informative genes from the dataset, ranking based on the deviation in value across cells<sup>1</sup>. The Seurat biological toolkit in R was also used for the purpose of low-variance feature selection (Satija, 2012). We eschewed dimensionality reduction techniques such as PCA because of loss of its direct feature interpretability. Biologists have hand-labeled every cell in our dataset using existing tool; there are a total of 9 cell categories. The complete dataset, as well our preprocessed and pruned versions, are available at (Yang & Koppula).

<sup>1</sup>We recognize that this can bias towards noisy genes. We favor this method because it is simple and easy to implement, and a practice used in literature (Ling, 2012)

### 1.2. Prior work

Prior research has explored the use of various computational techniques to analyze gene expression data. Most commonly, Spearman and Pearson correlation metrics are frequently used infer sets of genes that cluster together (Xie, 2015; Borenszstein, 2017). Other techniques, including PCA followed by linear regression, has been used for expression-based cell clustering (Stegle, 2016). Yu et al. propose a unsupervised classifier ensemble as another approach to cell clustering (Yu, 2016).

More Bayesian approaches have also been tried in prior work from the Pe'er lab. (Prabhakaran, 2016) uses a Hierarchical Dirichlet Mixture Model to learn cell clusterings. (Azizi, 2017) builds on this to jointly learn optimal normalization pre-processing of the data. Bayesian networks have also been used in an attempt to learn gene dependencies from expression data (Pe'er, 2000).

### 1.3. Structure of Report

We first discuss our experiments using Bayesian topic models to discover related sets of genes in a 'topic': LDA in Section 2 and Dynamic Time Models in Section 4. We discuss our experiments in clustering: Dirichlet Mixtures in Section 3, Integrated Topic-Clustering in Section 5, and non-parametric models in Section 6. Finally, we conclude our paper with our observations from across all studies.

## 2. Latent Dirichlet Allocation

### 2.1. Model Description

In the generative process for LDA, the topic assigned to a word in a document is drawn from the document's topic distribution. The identity of the word is drawn from topic's word distribution. Assuming the reader is familiar with LDA, we relegate further details and formalization of the model to (Blei, 2003).

In the context of analyzing gene expression data, we are interested in discovering 'topics' that comprise of a set of top- $N$  genes within the topic distribution that are biologically related. For example, together the genes may direct a specific chemical function in a cell. Biologists denote such sets of genes as 'gene modules' which can be cross-referenced with existing gene module databases.

## 2.2. Implementation

We explored two implementations of posterior inference for LDA on large datasets: one using online mean-field variational bayes for posterior estimate (ovb, 2017; Hoffman), and a broken C++ Gibbs sampler for LDA that we modified to work (pld, 2015). We were able to get this latter sampler running parallel across multiple cores, with a burn-in of 100 iterations. We compare these two posterior estimation approaches using our entire dataset, with a 10% held-out testing partition. We experimented with  $k = 5, 10, 25, 50$  topics.

More details about the Gibbs sampler: We had a fixed burn-in of number of samples (200), a fixed number of sampling iterations after that (500). We didn't extensively explore varying these values, but trying out significantly more iterations (700) didn't seem to change the topic's word distributions significantly. We could explore this in more detail for the final report. There was one sampling chain per core, with four cores. Not sure exactly what you had in mind by reporting 'performance details', other than time for estimation/core, but maybe a somewhat more fair comparison, like test perplexity/instructions executed during posterior estimation?

## 2.3. Experiments

- posterior predictive checks - no intrusion testing meaningful - hypergeometric tests - likelihood - performance metrics

We did gene set enrichment analysis using the minimum hypergeometric test (Wagner, 2015) to compare our topics, which are ranked lists of genes, with existing collections of genes catalogued in the comprehensive gene module database, MSigDB (Institute, 2015). Figure 4 in the Appendix shows topic-to-MSigDB module matches for which the  $p$ -value of the match is at least less than 0.3. Models with more topics tended to have more matches with higher significance. These results suggest that we need to train the model with more topics. We note that the  $p$ -values do *not* factor for multiple hypothesis corrections, so at the moment

we are only using them as relative measures of model quality.

3. More details about biological significance: Good question about # of topics affecting the results for the hypergeometric test for biological relevance/correlation with existing biologically annotated. This is something I'm not too sure about, and why we experience the results we do. That said, # topics is something that we were experimented with, but soon realized we didn't consider a broad enough range, after seeing the results from our non-parametric model afterward. More investigation and thought will need to go into this for the final report.

Figure 2 shows the time to complete each estimation method. As expected, Gibbs scales poorly with the parameter dimensionality and is strictly worse than on-line variational Bayes across all studied topic counts. We also calculated perplexity and log-likelihood for the OVB parameters, and very surprisingly, we found that log-likelihood decreased as number of topics increased on a held-out test set (Figure 3). We are currently trying to understand whether this is because of programming error.

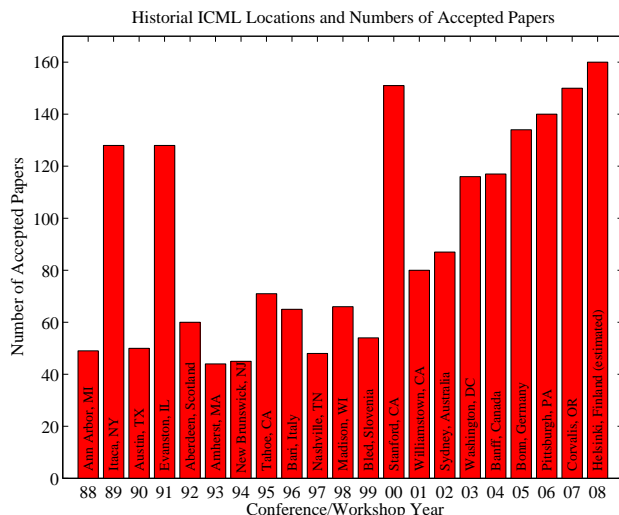


Figure 1. This is a demo figure.

## 3. Dirichlet Mixture Model

### 3.1. Model Description

Algorithm 1 describes the generative process for the mixture model.

**Algorithm 1** Mixture Model

---

**Input:** data  $x_i$ , size  $m$   
**repeat**  
  Initialize  $noChange = true$ .  
  **for**  $i = 1$  **to**  $m - 1$  **do**  
    **if**  $x_i > x_{i+1}$  **then**  
      Swap  $x_i$  and  $x_{i+1}$   
       $noChange = false$   
    **end if**  
  **end for**  
**until**  $noChange$  is  $true$

---

**3.2. Implementation****3.3. Experiments****4. Dynamic Time Model****4.1. Model Description****4.2. Implementation****4.3. Experiments****5. Integrated Topic-Clustering Model****5.1. Model Description****5.2. Implementation****5.3. Experiments****6. Non-parametric Models: IBP and HDP****6.1. Model Description****6.2. Implementation****6.3. Experiments****References**

C++ implementation of latent dirichlet allocation. <https://github.com/openbigdatagroup/plda/blob/master/lda.cc>, 2015.

Online latent dirichlet allocation with variational inference. [https://github.com/scikit-learn/scikit-learn/blob/master/sklearn/decomposition/online\\_lda.py](https://github.com/scikit-learn/scikit-learn/blob/master/sklearn/decomposition/online_lda.py), GitHub 2017.

Azizi, et al. Bayesian inference for single-cell clustering and imputing. <https://genomicscomputbiol.org/ojs/index.php/GCB/article/view/46/35>, Genomics 2017.

Blei, et al. Latent dirichlet allocation. <http://www.jmlr.org/papers/volume3/blei03a/blei03a.pdf>, JMLR 2003.

<http://www.jmlr.org/papers/volume3/blei03a/blei03a.pdf>, JMLR 2003.

Blei, et al. The IBP compound dirichlet process and its application to focused topic modeling. <http://www.cs.columbia.edu/~blei/papers/WilliamsonWangHellerBlei2010.pdf>, ICML 2010.

Borenszstein, et al. Xist-dependent imprinted x inactivation and the early developmental consequences of its failure. [http://www.nature.com/nsmb/journal/v24/n3/fig\\_tab/nsmb.3365\\_SF1.html](http://www.nature.com/nsmb/journal/v24/n3/fig_tab/nsmb.3365_SF1.html), Nature Structural & Molecular Biology 2017.

Hoffman, et al. Online learning for latent dirichlet allocation.

Institute, Broad. Molecular signatures database v6.0. <http://software.broadinstitute.org/gsea/msigdb>, PeerJ 2015.

Ling, et al. Improving relative-entropy pruning using statistical significance. [https://www.cs.cmu.edu/~awb/papers/coling2012/rep\\_coling2012.pdf](https://www.cs.cmu.edu/~awb/papers/coling2012/rep_coling2012.pdf), ACL 2012.

Pe'er, et al. Using bayesian networks to analyze expression data. <http://www.cs.huji.ac.il/~nir/Papers/FLNP1Full.pdf>, Genomics 2000.

Prabhakaran, et al. Dirichlet process mixture model for correcting technical variation in single-cell gene expression data. <http://www.c2b2.columbia.edu/danapeerlab/html/pub/prabhakaran16.pdf>, ICML 2016.

Satija, et al. Spatial reconstruction of single-cell gene expression data. <http://www.nature.com/nbt/journal/v33/n5/full/nbt.3192.html>, Nature Biotechnology 2012.

Stegle, et al. Computational and analytical challenges in single-cell transcriptomics. <https://www.nature.com/nrg/journal/v16/n3/full/nrg3833.html>, Nature Methods 2016.

Teh, et al. Stick-breaking construction for the indian buffet process. <http://mlg.eng.cam.ac.uk/zoubin/papers/TehGorGha07.pdf>, NIPS 2007.

Wagner, Florian. The xl-mhg test for enrichment: A technical report. <https://arxiv.org/pdf/1507.07905.pdf>, PeerJ 2015.

Xie, et al. SINCERA: A pipeline for single-cell rna-seq profiling analysis. <http://months.plos.org/ploscompbiol/article?id=10.1371/month.pcbi.1004575>, PLoS 2015.

Yang, Karren and Koppula, Skanda. Tumor cell melanoma data. [https://www.dropbox.com/sh/tdrplhiu3mzvmeu/AADF1ZITyPX\\_-I407jDuruJra?dl=0](https://www.dropbox.com/sh/tdrplhiu3mzvmeu/AADF1ZITyPX_-I407jDuruJra?dl=0).

Yu, et al. SC3 - consensus clustering of single-cell rna-seq data. <http://biorxiv.org/content/early/2016/09/02/036558>, Nature Methods 2016.

## 7. Appendix

### 7.1. Team’s Breakdown of Work

The implementation(s) and experiments involving LDA, Dirichlet Mixture, and Dynamic Time Model was completed by Skanda (Section 2, Section 3, and Section 4). Karren completed the implementations and experiments involved Integrated Topic-Clustering, Non-Parametric models, and the shared gene enrichment test implementation (Section 5 and Section 6. We believe the authors contributed equally in this work.

### 7.2. LDA vs. IBP-DP

One drawback to the HDP is that there tends to be a correlation between how frequently a topic appears across all documents and how prevalent this topic is within documents that it appears in. Williamson et al. (Blei, 2010) proposed the ‘focused topic model’ to overcome this drawback. In their model, each topic  $k = 1, 2, \dots$  has a relative prevalence  $\phi_k \sim \text{Gamma}(\gamma, 1)$  and a population frequency  $\pi_k = \prod_{j=1}^k \mu_k$ , where each  $\mu_k \sim \text{Beta}(\alpha, 1)$ . For each document  $m$ , whether topic  $k$  appears is sampled as  $b_{mk} \sim \text{Bernoulli}(\pi_k)$ , and the topic proportions are sampled as  $\theta_m \sim \text{Dirichlet}(b_m \cdot \phi)$ .

Since the code from the original IBP-DP paper was not available, we implemented an inference algorithm using collapsed Gibbs sampling (Blei, 2010). Due to the non-conjugacy of the model, sampling each latent variable from its full conditional required using another sampling method. To sample the topics parameters  $\pi$  and  $\phi$ , we used slice sampling based on the semi-ordered stick-breaking representation of the model (Teh, 2007).

We tested our code on a subset of the Reuters-21578 dataset, using several different values of the concentration hyper-parameter  $\alpha$ , which influences the number of clusters. Although higher values of  $\alpha$  yielded better log-likelihood values, we found that it resulted in a large number of very small topics, which are not very

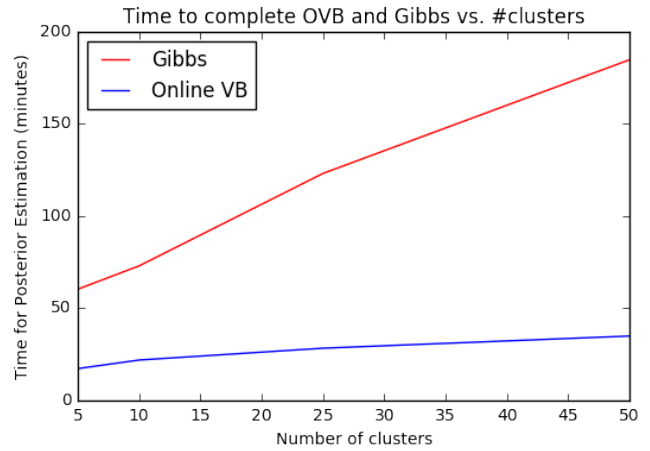


Figure 2. Comparison of the running times of each of the posterior estimation methods across various numbers of topics.

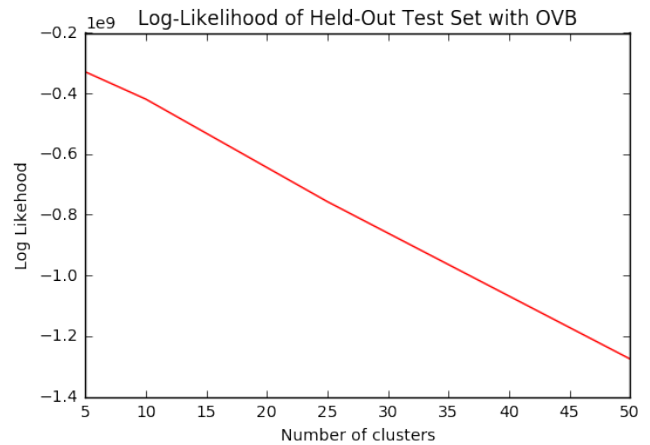


Figure 3. Log-likelihood of the held-out testing set, across various numbers of topics.

useful (Figure 7). Qualitatively, we did not find the topics from IBP-DP (Figure 8) to be more coherent than topics than LDA (Figure 9). The most prevalent topics from the IBP-DP each corresponded to similar topics from LDA; less prevalent topics tended to consist of a few unrelated words. These results discouraged us from optimizing the code to train this model on our scRNA-seq dataset, as we do not think it would yield more coherent gene modules than LDA. We emphasize that these results are not completely unexpected, as the authors of the IBP-DP paper did not show any topics from their model, nor did they assess the quality of their model with metrics other than perplexity.

## GIBBS TOPICS

```

cluster_size: 5
cluster_size: 10
  match: topic 1 , KINSEY TARGETS OF EWSR1 FLII FUSION UP , pval: 0.041958041958
  match: topic 1 , PILON KLF1 TARGETS DN , pval: 0.137062937063
  match: topic 6 , GRAESSMANN APOPTOSIS BY DOXORUBICIN DN , pval: 0.263736263736
cluster_size: 25
  match: topic 8 , PILON KLF1 TARGETS DN , pval: 0.201398601399
  match: topic 22 , KINSEY TARGETS OF EWSR1 FLII FUSION UP , pval: 0.041958041958
  match: topic 22 , PILON KLF1 TARGETS DN , pval: 0.137062937063
  match: topic 22 , GRAESSMANN APOPTOSIS BY DOXORUBICIN DN , pval: 0.263736263736
cluster_size: 50
  match: topic 2 , ROME INSULIN TARGETS DN , pval: 0.193006993007
  match: topic 3 , PILON KLF1 TARGETS DN , pval: 0.137062937063
  match: topic 6 , PILON KLF1 TARGETS DN , pval: 0.137062937063
  match: topic 7 , BLALOCK ALZHEIMERS DIS , pval: 0.263736263736
  match: topic 17 , KINSEY TARGETS OF EWSR1 FLII FUSION UP , pval: 0.041958041958
  match: topic 17 , PILON KLF1 TARGETS DN , pval: 0.137062937063
  match: topic 22 , PILON KLF1 TARGETS DN , pval: 0.137062937063
  match: topic 24 , WEI MYCN TARGETS WITH FOXO1 , pval: 0.083916083916

```

## OVb TOPICS

```

cluster_size: 5
cluster_size: 10
  match: topic 5 , PILON KLF1 TARGETS DN , pval: 0.263736263736
cluster_size: 25
  match: topic 5 , DIAZ CHRONIC MYELOGENOUS LEUKEMIA UP , pval: 0.263736263736
  match: topic 8 , KINSEY TARGETS OF EWSR1 FLII FUSION UP , pval: 0.041958041958
  match: topic 15 , PUJANA ATM PCC NETWORK , pval: 0.263736263736
  match: topic 15 , PUJANA BRCA1 PCC NETWORK , pval: 0.263736263736
cluster_size: 50
  match: topic 6 , GRAESSMANN APOPTOSIS BY DOXORUBICIN DN , pval: 0.263736263736
  match: topic 15 , PUJANA ATM PCC NETWORK , pval: 0.263736263736
  match: topic 22 , MARSON BOUND BY FOXP3 UNSTIMULATED , pval: 0.263736263736
  match: topic 28 , PILON KLF1 TARGETS DN , pval: 0.201398601399
  match: topic 38 , PUJANA BRCA1 PCC NETWORK , pval: 0.263736263736
  match: topic 41 , MARSON BOUND BY FOXP3 UNSTIMULATED , pval: 0.263736263736
  match: topic 42 , KINSEY TARGETS OF EWSR1 FLII FUSION UP , pval: 0.193006993007

```

Figure 4. Matches between the gene collections found in LDA topics and published gene sets in MSigDB. Cluster size refers to the number of topics in the model

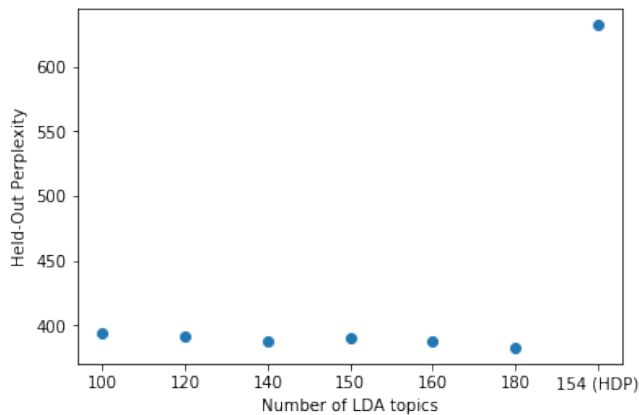


Figure 5. Comparison of log-likelihood of the held-out testing set, under various LDA models and the HDP model.

Figure 6. Comparison of distributions of p-values from gene set enrichment analysis between LDA models and the HDP model.

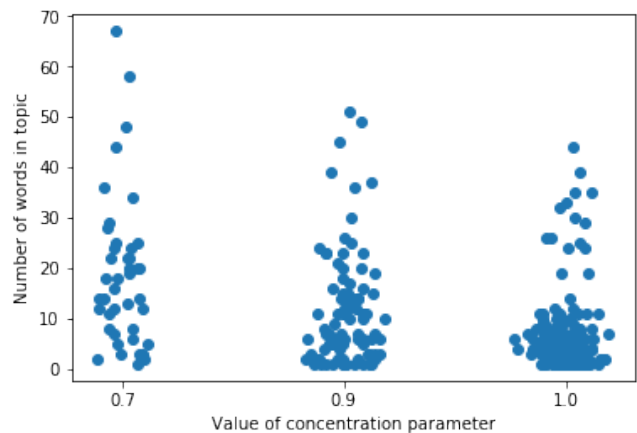


Figure 7. Beeswarm plot of number of words per topic, for 3 different IBP-DP models with different concentration parameters. Each point represents one topic from its model



Topic 0: charles diana prince royal parker bowles camilla queen family marriage public princess love britain england  
 Topic 1: church told during year time world very years made life saying last became take while  
 Topic 2: catholic n't bishop son women father love mother television day cardinal britain woman leaders days  
 Topic 3: sunday including against day won group held known last police official roman few come late  
 Topic 4: pope france visit john french paul first both pontiff trip church religious catholics king against  
 Topic 5: mother teresa doctors home heart charity hospital order tuesday people work told peace world house  
 Topic 6: set years ceremony led german germany second made called around rights married time few south  
 Topic 7: pope health vatican mass trip reporters surgery during saturday death left past spokesman people monday  
 Topic 8: people president state last later took around government good since say percent age won head  
 Topic 9: media later catholic former own head newspaper next known leader against wednesday taken n't called  
 Topic 10: official added first health ago under among around wednesday paul few minister monday children mother  
 Topic 11: former end throne century taken reports children newspaper england home french paris 1992 international state  
 Topic 12: french members visit leaders family germany national statement say three work against thursday year  
 Topic 13: four heart clinton percent age  
 Topic 14: life say wednesday left while month tuesday times later official take peace paris doctors including  
 Topic 15: says 1992 princess love left government political under german married took since saturday house around  
 Topic 16: bernardin cardinal among death u.s own surgery told minister doctors great several news end until  
 Topic 17: reports roman us told president times  
 Topic 18: paul white show reports end son union month monday bishop long  
 Topic 19: service end million since wednesday roman home wife spokesman son reports city  
 Topic 20: first local place princess british peace white saturday taken expected several united married made century  
 Topic 21: part say sunday members year leaders church days  
 Topic 22: former ago n't three century rome funeral group year saying led family  
 Topic 23: france home news south whose president work east take first both united women country officials  
 Topic 24: four week throne saying became former died members camilla made long show  
 Topic 25: statement french week thursday war christian born vatican house heart leader britain 1992 set three  
 Topic 26: minister prime expected group officials union died children times michael whose around american church  
 Topic 27: years president  
 Topic 28: diana funeral service princess reuters hospital son died  
 Topic 29: monday u.s around several four children throne year since john tuesday members churchill told statement  
 Topic 30: day friday party private wednesday later british officials former family until throne white  
 Topic 31: against  
 Topic 32: city international good prize second won paris take since died house off local  
 Topic 33: times children world  
 Topic 34: paris work  
 Topic 35: michael king paul local father political show ceremony next part private war german week whose  
 Topic 36: rights government prize service million become officials held police head reuters political us party  
 Topic 37: television reporters show n't political times clinton several off own government son wednesday very during  
 Topic 38: known news  
 Topic 39: part great mass later women past says  
 Topic 40: told charles called off born  
 Topic 41: taken ceremony president  
 Topic 42: thursday long leaders

Figure 8. Topics from IBP-DP model trained on subset of Reuters dataset.

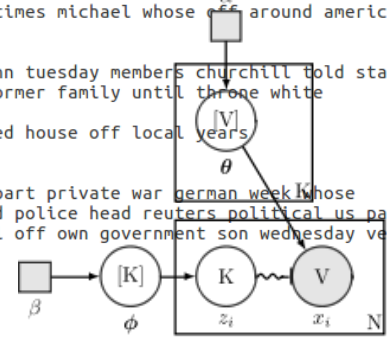


Figure 10. Finite K-sized mixture model currently implemented.  $\theta$  is the parameter for every cluster component, represented from a categorical draw of over all genes.  $z_i$  is the cluster assignment, and  $\phi$  is the distribution of clusters. As usual,  $\alpha, \beta$  are hyper-parameters.

Topic 0 : france home news years work mother president during take whose women east country love told  
 Topic 1 : church years cardinal bishop take england against million vatican past news british told sunday ceremony  
 Topic 2 : pope health mass during visit saturday trip told john paul pontiff people church service spokesman  
 Topic 3 : mother teresa heart sunday home hospital tuesday told doctors order people catholic charity peace house  
 Topic 4 : pope france french visit church trip first paul catholic pontiff both john state including paris  
 Topic 5 : teresa mother doctors charity official hospital home work first around told world during under saying  
 Topic 6 : church media michael paul former marriage princess england never n't love very told public years  
 Topic 7 : bishop church catholic son father n't told women love mother woman roman years leaders ago  
 Topic 8 : royal family queen prince charles throne church princess century britain first british media 1992 head  
 Topic 9 : order day city friday group during own doctors monday very reuters prize last people roman  
 Topic 10 : television told show n't reporters president later day own times political clinton off year years  
 Topic 11 : president rights government people last church says state life died told political country group catholic  
 Topic 12 : diana charles princess britain time wednesday ago family monday million camilla church newspaper bowles parker  
 Topic 13 : charles parker bowles prince camilla diana royal marriage public queen love king church woman family  
 Topic 14 : pope bernardin vatican church surgery health year time left told life say death cardinal made

Figure 9. Topics from LDA model with 15 topics trained on subset of Reuters dataset.