
Latent Dirichlet Allocation

Si Chen and Yufei Wang

Department of Electrical and Computer Engineering
University of California San Diego
{sic046, yuw176}@ucsd.edu

Abstract

Latent Dirichlet allocation(LDA) is a generative topic model to find latent topics in a text corpus. It can be trained via collapsed Gibbs sampling. In this project, we train LDA models on two datasets, Classic400 and BBCSport dataset. We discuss possible ways to evaluate goodness-of-fit and to detect overfitting problem of LDA model, and we use these criteria to choose proper hyperparameters, observe convergence, and evaluate the models, the criteria we use include perplexity, VI-distance, visualization of clustering results, and highest-probability words.

1 Introduction

Latent Dirichlet allocation introduced by [1] is a generative probabilistic model for collection of discrete data, such as text corpora. It assumes each word is a mixture over an underlying set of topics, and each topic is a mixture over a set of topic probabilities. Evaluating the models is a tough issue. There are several types of methods that people use: The models can be applied to some tasks such as document classification, where the performance can be easily evaluated; Several methods estimate the likelihood of held-out documents; Subjective judgement can be made by examine word and document similarities. In this project, we learn the models of two datasets, Classic400 and BBCSport dataset, by collapsed Gibbs sampling, and use several methods to evaluate the models, including perplexity, VI-distance, visualizing result and highest-probability words.

This report is organized as follows. Section 2 gives a brief overview of LDA model and training process, and introduces several methods to evaluate goodness-of-fit and check overfitting. We describe some implementation details in Section 3. In Section 4, we describe the design of our experiments. Section 5 shows the experiment results with discussions. Finally, we draw some conclusions in Section 6.

2 Theoretical Overview

2.1 Latent Dirichlet Allocation

LDA is a mixture model. It assumes that each document contains various topics, and words in the document are generated from those topics. All documents contain a particular set of topics, but the proportion of each topic in each document is different.

The generative process of the LDA model can be described as follows:

Given: Dirichlet distribution with parameter vector α of length K

Given: Dirichlet distribution with parameter vector β of length V

for topic number 1 to topic number K

draw a word distribution, i.e. a multinomial with parameter vector ϕ_k

according to β . $\phi \sim \text{Dirichlet}(\beta)$
for document number 1 to topic number M
draw a topic distribution, i.e. a multinomial with parameter vector θ
according to α . $\theta \sim \text{Dirichlet}(\alpha)$
for each word in the document
draw a topic z according to θ . $z \sim \text{Multinomial}(\theta)$
draw a word w according to ϕ_z . $w \sim \text{Multinomial}(\phi_z)$

Not that V is the cardinality of the vocabulary. The words of a document is observed, but the topic assignment is latent.

2.2 Training via collapsed Gibbs sampling

Training data in this project is the words in all documents. The goal of training is to infer the multinomial parameters θ for each document, and ϕ_k for each topic.

We use collapsed Gibbs sampling for learning. First it infers the hidden value z_{nm} for each word occurrence in each document: $p(\vec{z}|\vec{w})$. Note z_{nm} is the mixture indicator that chooses the topic for n th word in document m . \vec{w} is the vector of words making up the entire corpus, and \vec{z} is the corresponding topics.

The idea of Gibbs sampling is assuming that we know the value of \vec{z} for every word occurrence in the corpus except occurrence number i , then we draw a z_i value for i according to its distribution. Then we assume that this value is known to be true value, and draw a z_j for another word, and so on. Eventually this process will converge to a correct distribution $p(\vec{z}|\vec{w})$. Therefore, we need the conditional distribution for a word token with index $i = (m, n)$:

$$p(z_i = k | \vec{z}_{-i}, \vec{w}) \propto \frac{n_{k,-i}^{(i)} + \beta}{\sum_{t=1}^V (n_{k,-i}^{(t)} + \beta)} (n_{m,-i}^{(k)} + \alpha) \quad (1)$$

where z_i is the topic association of i th word. \vec{z}_{-i} is \vec{z} with i th topic removed. $n_k^{(i)}$ is the number of times that word t has been observed with topic k . $n_m^{(k)}$ refers to the number of times that topic k occurs with a word of document m . $n_{\cdot,-1}^{(i)}$ indicates that the word i is excluded from the corresponding document or topic.

Finally, we can obtain the multinomial parameters θ and ϕ :

$$\phi_{k,t} = \frac{n_k^{(t)} + \beta}{\sum_{t=1}^V (n_k^{(t)} + \beta)} \quad (2)$$

$$\theta_{m,k} = \frac{n_m^{(k)} + \alpha}{\sum_{k=1}^K (n_m^{(k)} + \alpha)} \quad (3)$$

where $\phi_{k,t}$ is the multinomial parameter for word t drawn from topic k . $\theta_{m,k}$ is the parameter of topic k drawn from document m .

2.3 Goodness-of-fit of the model

There are two typical ways to evaluate LDA model.

2.3.1 Likelihood of held-out data

One way is to calculate the probability of held-out documents that are not used for training. The likelihood $p(\vec{w}|\alpha, \beta)$ is intractable, and there are several ways to estimate it.

The probability of the held-out documents \vec{w} given training data \vec{w}' can be written as

$$p(\vec{w}|\vec{w}') = \int d\phi d\alpha p(\vec{w}|\phi, \alpha) p(\phi, \alpha|\vec{w}') \quad (4)$$

[2] approximates this integral by evaluating at a point estimate, and the problem becomes how to evaluate $p(\vec{w}|\phi, \alpha)$. [2] summaries several estimating methods: importance sampling methods, harmonic mean method, and annealed importance sampling. [2] also proposes two alternative methods: a Chib-style estimator and a “left-to-right” evaluation algorithm.

For example, the frequently used method, Harmonic mean method, estimates the probability as follows

$$p(\vec{w}|\phi, \alpha) \simeq \frac{1}{\frac{1}{s} \sum_s \frac{1}{p(\vec{w}|\vec{z}^{(s)}, \phi)}} \quad (5)$$

where $\vec{z}^{(s)}$ is drawn from $p(\vec{z}|\vec{w})$.

[3] approximates the held-out likelihood with empirical likelihood(EL). It produces a large set of pseudo documents with parameter α and β , and the pseudo documents are then used to train a tractable model. The true likelihood of the test set is then estimated as its likelihood under the tractable model.

Perplexity is another way to calculate the likelihood. It is defined as the reciprocal geometric mean of the token likelihoods in the test corpus given the model:

$$p(\vec{W}|M) = \exp - \frac{\sum_{m=1}^M \log p(\vec{w}_{\tilde{m}}|M)}{\sum_{m=1}^M N_m} \quad (6)$$

where M is the trained model, and $\vec{w}_{\tilde{m}}$ is the word vector in document \tilde{m} . Lower values of perplexity indicate lower misrepresentation of the words of the test documents by the trained topics.

The log-likelihood can be expressed as:

$$\log p(\vec{w}_{\tilde{m}}|M) = \sum_{t=1}^V n_{\tilde{m}}^{(t)} \log \left(\sum_{k=1}^K \phi_{k,t} \cdot \theta_{\tilde{m},k} \right) \quad (7)$$

where $n_{\tilde{m}}^{(t)}$ is the number of times word t occurs in document \tilde{m} .

2.3.2 Evaluation of clustering result

Another way to evaluate the model is to measure its performance on some secondary tasks, for example, document classification or clustering. The LDA model provides a soft clustering of the documents. The evaluation of clustering quality can be done in many ways:

1. We can simply check if the most frequent words generated from each topic are semantically related.
2. We can also check if intra-class document similarity is higher than inter-class document similarity.
3. When we have the label of each document, and the number of topics K equals the number of given labels, we can convert the soft clustering result into a hard result, and calculate the categorization accuracy. Converting soft clustering to a hard one is simply done by assigning topic with highest probability to each document.
4. We can also directly compare the soft clustering result with the priori categorization([4]). Variation of Information distance (VI-distance) can be used to compare the two clusterings. It assumes two distributions over class for each document: $p(c = j|d_m)$ and $p(z = k|d_m)$, where $j \in [1, J]$ and $k \in [1, K]$ are the class labels or topics of the two distribution. Note that K doesn't need to be the same with J . d_m is m th document. The class probabilities are obtained by averaging over the corpus: $p(c = j) = \frac{1}{M} \sum_m p(c = j|d_m)$ and $p(z = k) = \frac{1}{M} \sum_m p(z = k|d_m)$. The joint probability of co-occurring pairs is $p(c = j, z = k) = \frac{1}{M} \sum_m p(c = j|d_m)p(z = k|d_m)$. Then the VI-distance measure is defined as follow:

$$D_{VI}(C, Z) = H(C) + H(Z) - 2I(C, Z) \quad (8)$$

with

$$I(C, Z) = \sum_{j=1}^J \sum_{k=1}^K p(c = j, z = k) [\log_2 p(c = j, z = k) - \log_2 p(c = j)p(z = k)] \quad (9)$$

$$H(C) = - \sum_{j=1}^J p(c = j) \log_2 p(c = j) \quad (10)$$

$$H(Z) = - \sum_{k=1}^K p(z = k) \log_2 p(z = k) \quad (11)$$

VI-distance is always nonnegative, and smaller VI-distance indicates the two clustering are more similar.

2.4 Overfitting monitoring

In Section 2.3.1, we introduce several ways to calculate the likelihood of held-out data. In addition to using likelihood or perplexity to evaluate the goodness-of-fit, we can also use it to monitor overfitting. By calculating likelihood/perplexity of the training data, and comparing it with likelihood/perplexity of test data, we can get the idea whether overfitting occurs. When no overfitting occurs, the difference between two types of likelihood should remain low.

3 Implementation of algorithm

Algorithms are realized in Python.

Data The vocabulary of words is very large. However, in one document, only a small part of the vocabulary occurs, and there is no need to traverse every word in vocabulary for each document. Therefore, word counts of all documents are stored in `scipy.sparse.csc_matrix`, a compressed sparse column matrix. Non-zero values can be efficiently located, so the training process can be speeded up.

Initialization The hyperparameters α , β , and K are predefined. The initial topics z_{mn} associating with words w_{mn} are random number ranging from 1 to K . For words in all documents, the counts $n_k^{(t)}$ and $n_m^{(k)}$ are calculated.

Gibbs sampling The counts $n_k^{(t)}$ and $n_m^{(k)}$ are stored in matrices. Instead of looping over K to calculate the conditional probability $p(z_i = k | \vec{z}_{-i}, \vec{w})$ for each topic k , we use matrix operation to calculate unnormalized conditional probability of all topics. This is faster than simply looping over K topics for each word.

Visualization of clustering result To evaluate the clustering result in an intuitive manner, we draw a mapping from a document to a 3D point according to each document's multinomial parameter θ . When there are 3 topics or less, we can simply assign each topic to one axis. In this way, each coordinate corresponds to one topic distribution (θ). When there are more than 3 topics, we use principle component analysis (PCA) to reduce the dimension. Three principle components are extracted, and then each documents can be mapped to a 3D point. We use points with different colors and shapes to represent different classes of documents. Therefore we can observe classification result as well as clustering result in the 3D image.

4 Design of experiments

4.1 Datasets

We use two datasets for experiments. The first is Classic400 dataset ¹. It consists of 400 documents, with a vocabulary of 6205 words. There are 3 categories of the documents. The second dataset we use is the BBC Sports dataset ². It consists of 737 documents with a vocabulary of 4613 words. There

¹ <http://cseweb.ucsd.edu/users/elkan/151/classic400.mat>

² <http://mlg.ucd.ie/datasets/bbc.html>

are 5 classes of the documents. For both datasets, we randomly choose 10% from each category as test data, and 90% as training data.

4.2 Choice of hyperparameters

We use perplexity and VI-distance as the criteria to choose the three hyperparameters: α , β , and K . Best choices of α and K have strong correlation, as is shown in [5], which chooses hyperparameters as $\alpha = 50/K$, $\beta = 0.1$. For both datasets, we have the true label of every document. Therefore, an intuitive guess of K would be the number of classes. Grid search for all three hyperparameters requires too much work, so we first fix K as the number of classes, and apply grid search of α and β . After deciding best α and β for fixed K , we then search for best K with fixed β and changing α , under the assumption that $\alpha \propto \frac{1}{K}$.

4.3 Evaluating results

We use several ways to evaluate the models:

1. We look at the 10 highest-probability words for each topic, to see if topics we learn are meaningful.
2. We calculate perplexity of training data and perplexity of test data for every 5 epochs, to check convergence and check overfitting.
3. We calculate VI-distance of training data to monitor clustering result.
4. We visualize the clustering result in 3 dimensional space.
5. We also compare the running time of each Gibbs sampling epoch of the two models.

5 Experimental results

5.1 Classic400 dataset

5.1.1 Best hyperparameters

By fixing K at 3, we use grid search for pairs (α, β) . We choose from $\alpha = \{0.01, 0.1, 1\} = \{0.03/K, 0.3/K, 3/K\}$, and $\beta = \{0.1, 1, 2\}$. Figure 1 shows the visualizing result of the 9 pairs of parameters, and Table 1 shows the corresponding perplexity and VI-distance.

From Figure 1, we can see the role of α . For smaller α , documents in each cluster are more scattered; for larger α , different clusters are more separated, and the documents are prone to locate in three corners of the triangle. This is also illustrated in Table 1, where VI-distance increases with α . This observation is intuitive: α is the pseudo count of topic per document. Larger α will bring an averaging effect of different topics.

From Table 1, we can observe the role of β . For smaller β , perplexity of training data is smaller. However, perplexity of test data doesn't change much with β . β is the pseudo count of word per document. Larger β brings a larger averaging effect of all the words, therefore making the model not so fit for training data.

We choose $(\alpha = 0.3/K, \beta = 1)$ as the best parameter pairs, to achieve low perplexity and VI-distance.

Then, we fix $\beta = 1$, $\alpha = 0.3/K$, and search for best K . We check perplexity, VI-distance and the mapping image for $K = \{2, 3, 4, 5\}$. Figure 2 shows visualization of clustering results with different K . For different K , location of points forms different 3D shapes: for $K = (2, 3, 4, 5)$, the shapes are line, triangle, tetrahedron, tetrahedron with an extra point in the middle, respectively. Points reside in vertex of those shapes. Table 2 shows the evaluation of goodness-of-fit. $K = 3$ achieves best perplexity on test data and VI-distance. Training perplexity decreases with larger K , because when we have larger topic number we have more parameters to fit the training data, however, bigger K causes overfitting: Test perplexity increases when K is big, this is because when K becomes too large, the number of parameters are too large for the training data. Therefore, we choose $K = 3$ as the best parameter.

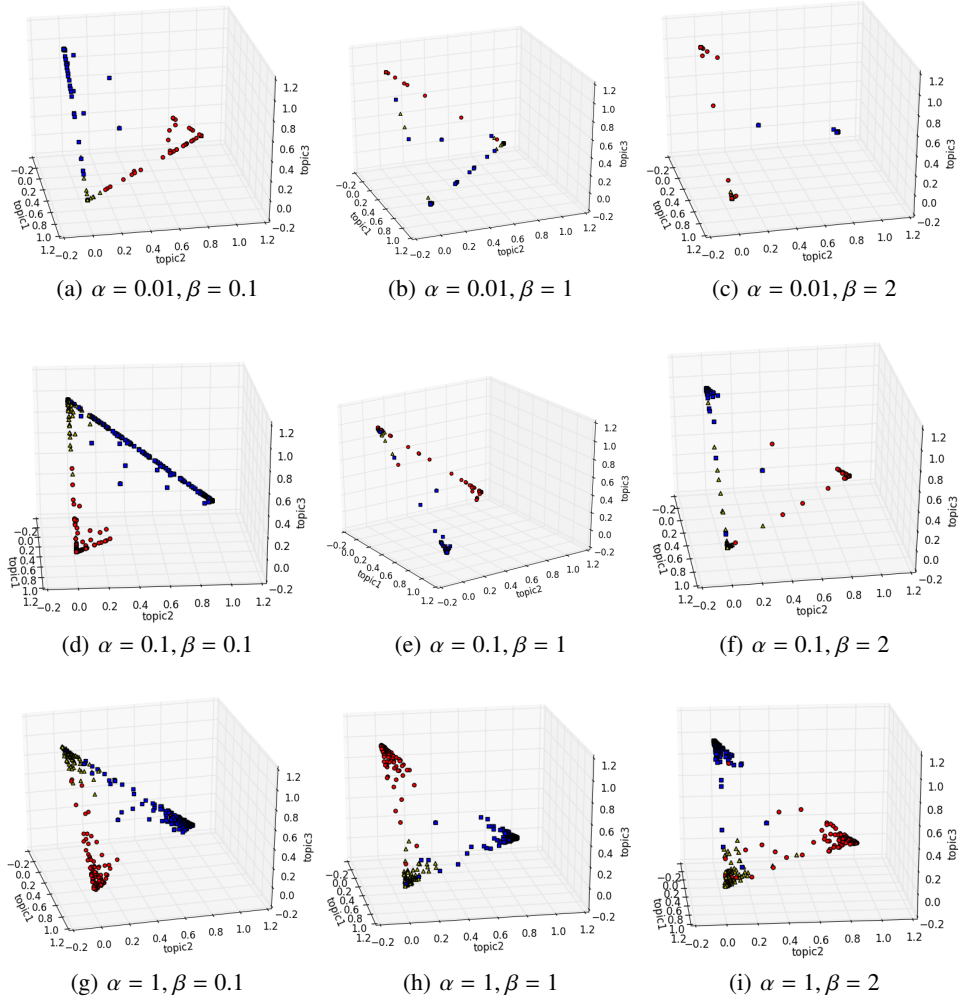


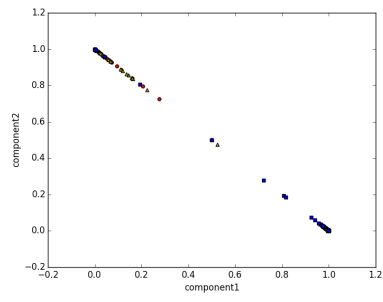
Figure 1: Plots of Classic400 dataset with $K = 3$ and varying α and β .

Parameter(α, β)	Perplexity (training)	Perplexity (test)	VI-distance (training)
(0.01, 0.1)	1376.024120	2500.267304	0.796688
(0.01, 1)	1803.455444	2303.252994	1.381389
(0.01, 2)	2088.023795	2506.341132	0.909851
(0.1, 0.1)	1451.690149	2646.197142	1.387363
(0.1, 1)	1753.193891	2256.134722	0.801898
(0.1, 2)	2071.336636	2462.702027	0.704664
(1, 0.1)	1484.825815	2536.642742	1.567750
(1, 1)	1802.399166	2310.176229	1.260563
(1, 2)	2128.084925	2492.720476	1.362268

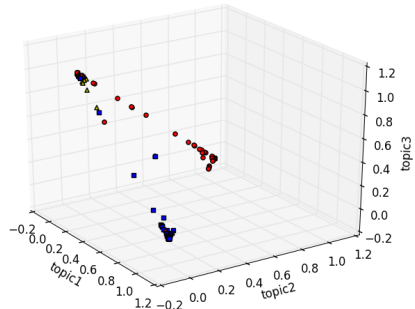
Table 1: Perplexity and VI-distance of Classic400 dataset, with $K = 3$ and varying α, β .

Parameter(K)	Perplexity (training)	Perplexity (test)	VI-distance(training)
2	1931.919028	2415.384221	0.864250
3	1753.193891	2256.134722	0.801898
4	1718.810470	2283.469816	1.339418
5	1704.840676	2364.641784	1.864304

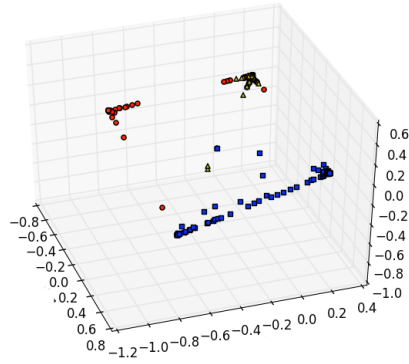
Table 2: Perplexity and VI-distance of Classic400 dataset, with $\beta = 1, \alpha = 0.3/K$ and varying K



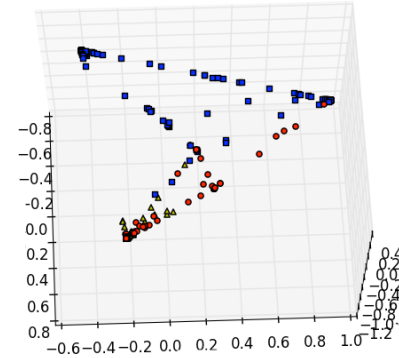
(a) $K = 2, \alpha = 0.15$



(b) $K = 3, \alpha = 0.1$



(c) $K = 4, \alpha = 0.075$



(d) $K = 5, \alpha = 0.06$

Figure 2: Plots of Classic400 dataset with $\beta = 1, \alpha = 0.3/K$ and varying K .

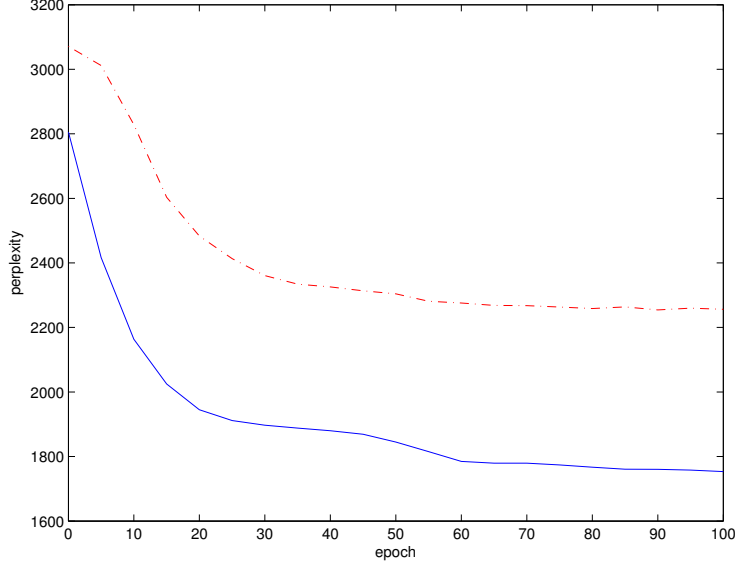


Figure 3: Training and test perplexity of LDA model for Classic400. Blue hard line: training data; Red dash-dot line: test data.

Topic	Ten most frequent words
1: 'Medical'	patients fatty acids nickel ventricular aortic cases left glucose septal
2: 'Scientific methods'	system scientific retrieval research language science methods systems subject journals
3: 'aero-physics'	boundary layer wing supersonic velocity mach wings ratio jet plate

Table 3: Top 10 frequent words for each topic of LDA model for Classic400.

5.1.2 Evaluation results

We have decided the parameters to be $K = 3, \alpha = 0.3/K, \beta = 1$. Then, we evaluate the model.

Figure 3 shows training and test perplexity of the model with number of epochs. Perplexity on test data is higher than that on training data, and the dropping rate of test perplexity is slightly lower. However, the trend of the two perplexities are similar. Therefore, we could infer there may exist slight overfitting, this may be caused by lack of data (we only have 360 training documents). According to the figure, the model converges after 60 epochs.

Table 3 shows 10 words that most frequently occur in each topic. In each topic, words are semantically related, and they are related to the true topic assigned by human. This means the model we train has a good clustering for the topics, with semantic meaning.

We calculate the average running time for every epoch of Gibbs sampling. The average running time of one epoch is 3.3649 seconds, which is shown in Table 7.

5.2 BBCSport dataset

5.2.1 Best hyperparameters

First we fix $K = 5$, which is the number of labels given. With the experience in Classic400 dataset, this time we do the grid search for $\alpha = \{0.3/K, 3/K\}, \beta = \{1, 2\}$. Figure 4 shows the visualization result of four models with each parameter pair. Table 4 shows the perplexity and VI-distance of the corresponding models.

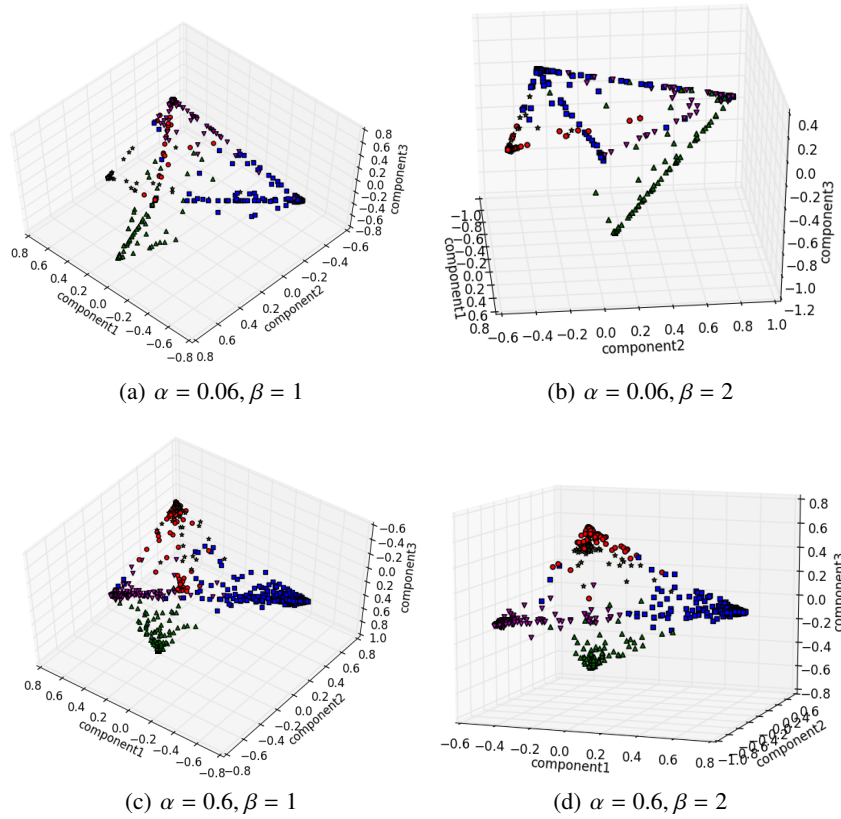


Figure 4: Plots of BBCSports dataset with $K = 5$ and varying β and α .

Parameters (α, β)	Perplexity(training)	Perplexity(test)	VI-distance
(0.06,1)	1194.730739	1413.515161	1.836619
(0.06,2)	1278.352876	1499.750074	1.735083
(0.6,1)	1207.601372	1434.633843	2.127757
(0.6,2)	1282.181586	1491.711059	2.013811

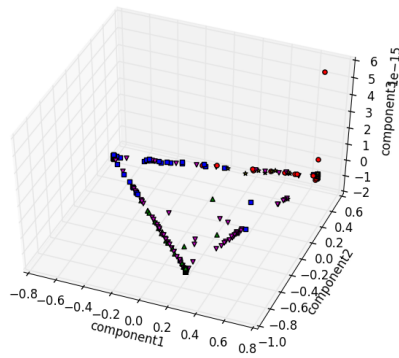
Table 4: Perplexity and VI-distance of BBCSport dataset, with $K = 5$ and varying α, β .

From Figure 4, we can see that when α increases, the documents are more scattered. Note that there are 5 topics, and PCA cannot separate them entirely on the 3D image. From Table 4, we observe that larger β leads to larger perplexity on both training and test data. VI-distance is larger with larger α . These trends are similar with what we observe at Classic400 dataset. For best perplexity and VI-distance, we choose the pair $(\alpha = 0.3/K, \beta = 1)$. Then, we fix $(\alpha = 0.3/K, \beta = 1)$ and change K , with $K = \{3, 4, 5, 6\}$. From Figure 5 and Table 5, we choose $K = 5$, with lowest perplexity in both training and test data, and low VI-distance.

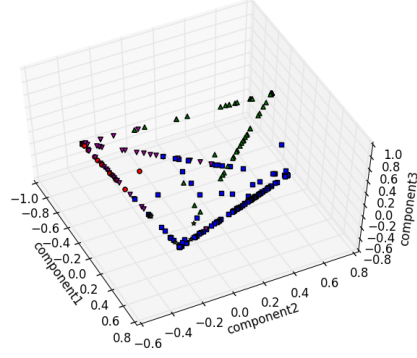
Therefore, the hyperparameters we choose are $K = 5, \alpha = 0.3/K, \beta = 1$.

Parameter (K)	Perplexity(training)	Perplexity(test)	VI-distance
3	1347.412626	1610.204673	1.691689
4	1272.958416	1498.080381	1.801365
5	1163.215412	1395.538324	1.671155
6	1155.549721	1402.170199	2.183145

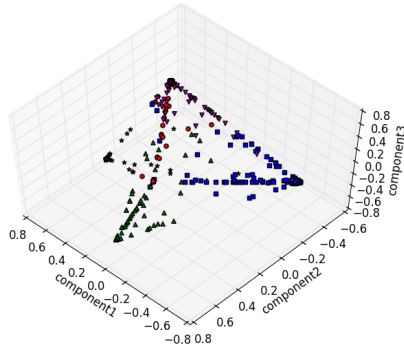
Table 5: Perplexity and VI-distance of BBCSport dataset, with $\beta = 1, \alpha = 0.3/K$ and varying K



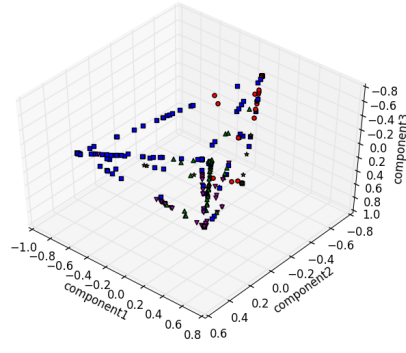
(a) $K = 3, \alpha = 0.1$



(b) $K = 4, \alpha = 0.075$



(c) $K = 5, \alpha = 0.06$



(d) $K = 6, \alpha = 0.05$

Figure 5: Plots of BBCSport dataset with $\beta = 1, \alpha = 0.3/K$ and varying K .

Topic	Ten most frequent words
1: Rugby	england wale game ireland rugbi against nation plai six player
2: Athletics	olymp world athlet race year indoor athen test champion win
3: Cricket	test cricket plai england first seri south match run australia
4: Football	game player club plai chelsea arsen unit leagu goal footbal
5: Tennis	plai open win match first set year final game roddick

Table 6: Top 10 frequent words for each topic of LDA model for BBCSport.

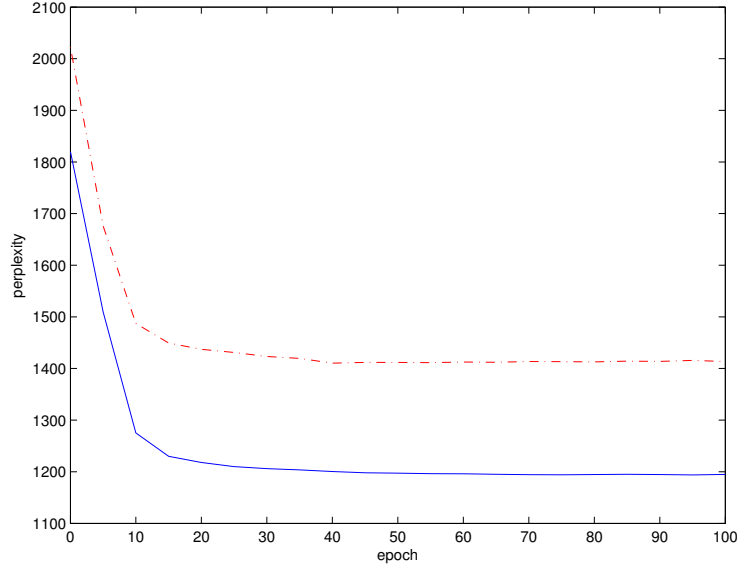


Figure 6: Training and test perplexity of LDA model for BBCSport. Blue hard line: training data; Red dash-dot line: test data.

5.2.2 Evaluation results

Figure 6 shows training and test perplexity of the model as number of epochs increases. Perplexity on test data is still higher than that on training data, and the dropping rate of test perplexity is slightly lower. Compared with Figure 6, we can see that model of BBCSport has less overfitting than Classic400’s model. The reason may be there being more training data in this BBCSport dataset: BBCSport training set has 109579 words in total, while Classic400 training set has only 27415 words. Although BBCSport dataset needs to train more parameters, it still has larger dataset-parameter ratio. From the figure, we can also observe rate of convergence of the model. The model converges after 40 epochs.

Table 6 shows ten words that most frequently occur in each topic. The words in the dataset are pre-processed before training: Stemming is performed. This avoids a word of different tense or a word in plural form being regarded as a different word. Since we only care about the topic related to the words, eliminating the impact of tense or plural form is sensible. Note that $K = 5$ is the true label number, therefore we can assign each topic to a label, judging by the semantic meaning of the top 10 words. Some words are very closely related to the label, for example, “chelsea”, “arsen” (which stands for arsenal) are strongly related to football. There are other words that are shared by some topics, for example, “year”, “win” and “plai” are shared by two or three topics. This is because the five topics are all about sports. They are semantically close to each other, and share some words that associate with all kinds of sports.

From one training process, we calculate the average running time for every epoch of Gibbs sampling. The average running time of one epoch is 12.3923 seconds, as is shown in Table 7. The normal time

Datasets	Number of topics: K	Total number of words:N	Running time (s)
Classic400	3	27415	3.3649
BBCSport	5	109579	12.3923
BBCSport	6	109579	12.3834
BBCSport	10	109579	12.1788
BBCSport	20	109579	11.9438
BBCSport	1000	109579	18.0922

Table 7: Comparison of running time of two datasets.

to perform one epoch of Gibbs sampling is $O(NK)$, where N is the total number of words in all documents. However, as is described in Section 3, we can make the time much less relevant to K by using matrix operation, making the time comparable to $O(N)$ when K is not too large. This is also illustrated in Table 7: For BBCDataset, increasing K from 5 to 20 doesn't affect running time much. When K becomes 1000, there is a noticeable time increase. Therefore, one-epoch running time for BBCSport dataset is approximately 4 times of Classic400 dataset, which is proportional to their word number ratio.

6 Conclusions

In this project, we realize the LDA model by Gibbs sampling, and evaluate it on two datasets. We use multiple types of evaluation criteria to check performance of the model. From the experiment results, we find that hyperparameters have different roles in the model: increasing α makes document clusters more scattered; smaller β results in smaller perplexity when no overfitting occurs; and larger topic number K makes training perplexity smaller but may cause overfitting. The perplexity results of two datasets suggest probable overfitting, which is caused by insufficient data. However, we can observe that latter dataset with larger data gets less overfitting. One problem of this implementation of LDA model is that it is not well scalable. Training time for one epoch is proportional to scale of dataset. Therefore training will become very slow with a large dataset.

References

- [1] Blei, D.M., Ng, A.Y., Jordan, M.I.: Latent dirichlet allocation. J. Mach. Learn. Res. **3** (2003) 993–1022
- [2] Wallach, H.M., Murray, I., Salakhutdinov, R., Mimno, D.: Evaluation methods for topic models. In: Proceedings of the 26th Annual International Conference on Machine Learning, ACM (2009) 1105–1112
- [3] Doyle, G., Elkan, C.: Accounting for burstiness in topic models. In: Proceedings of the 26th Annual International Conference on Machine Learning, ACMv (2009) 281–288
- [4] Heinrich, G.: Parameter estimation for text analysis. Technical report (2004)
- [5] Griffiths, T.L., Steyvers, M.: Finding scientific topics. Proceedings of the National Academy of Sciences (2004)